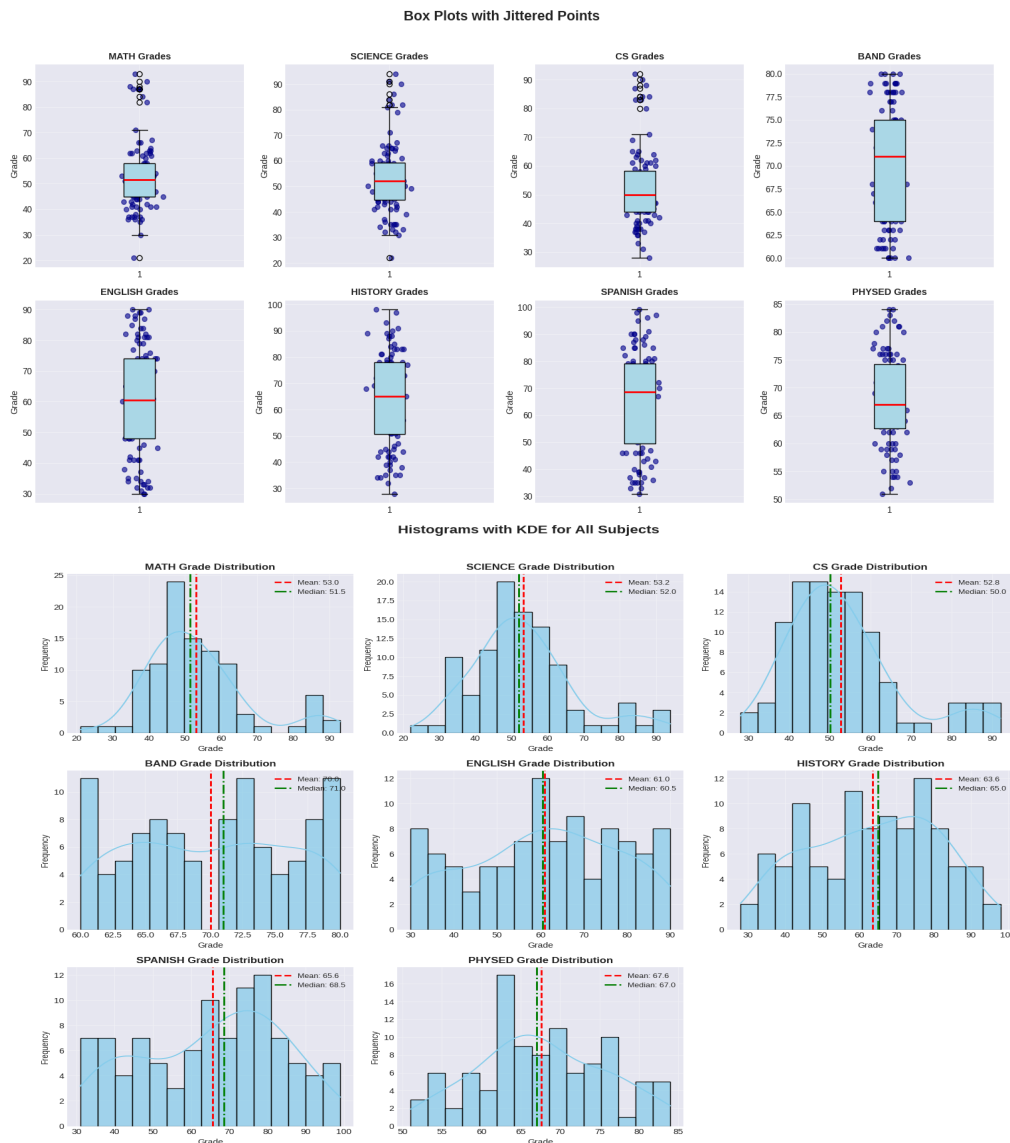# SUBJECTIVE QUESTIONS

Q1. What patterns were observed in the student grade distributions across different subjects?

Answer:
The analysis revealed distinct distribution patterns across the eight subjects. Band demonstrated the most symmetric distribution with the highest mean grade (72.1) and lowest variability (standard deviation = 10.6), indicating consistent high performance across students. Mathematics showed the lowest mean grade (51.6) with a left-skewed distribution, suggesting a concentration of students struggling with the subject. Computer Science exhibited the highest variability (standard deviation = 14.1), reflecting diverse skill levels within the student population. Physical Education displayed a clear bimodal distribution, indicating two distinct performance groups. Most academic subjects (Mathematics, Science, English, History) showed approximately normal distributions with slight negative skew, meaning more students performed below average than above.



Box Plots with Jittered Points



Histograms with KDE for All Subjects

Q2. Which subjects showed the highest and lowest average performance? What factors might explain this?

Answer:
Band achieved the highest average performance (72.1), followed by Physical Education (68.7) and History (68.0). Mathematics recorded the lowest average (51.6), with Computer Science (52.5) and Science (53.4) also showing below-average performance.

The high performance in Band may be attributed to its nature as an elective course, potentially attracting students with prior interest or experience. Additionally, ensemble-based grading often emphasizes participation, effort, and improvement rather than absolute mastery. The low performance in Mathematics reflects a common pattern in educational data, where quantitative reasoning subjects typically show wider performance distributions and lower averages. This may indicate prerequisite knowledge gaps, less differentiated instruction, or the cumulative nature of mathematical knowledge where early struggles compound over time.

| Statistic | Mathematics | Science | Computer Science | Band | English | History | Spanish | Physical Education |
|---|---|---|---|---|---|---|---|---|
| Mean | 53.05 | 53.22 | 52.75 | 70.04 | 60.97 | 63.60 | 65.56 | 67.61 |
| Median | 51.50 | 52.00 | 50.00 | 71.00 | 60.50 | 65.00 | 68.50 | 67.00 |
| Mode | 47.00 | 54.00 | 50.00 | 66.00 | 60.00 | 78.00 | 46.00 | 69.00 |
| Std Dev | 13.91 | 14.00 | 13.46 | 6.28 | 17.31 | 17.48 | 18.41 | 8.07 |
| Variance | 193.56 | 196.09 | 181.24 | 39.41 | 299.56 | 305... | 338.75 | 65.11 |
| Min | 21.00 | 22.00 | 28.00 | 60.00 | 30.00 | 28.00 | 31.00 | 51.00 |
| Max | 93.00 | 94.00 | 92.00 | 80.00 | 90.00 | 98.00 | 99.00 | 84.00 |
| Range | 72.00 | 72.00 | 64.00 | 20.00 | 60.00 | 70.00 | 68.00 | 33.00 |
| Q1 (25%) | 45.00 | 44.75 | 44.00 | 64.00 | 48.00 | 50.75 | 49.50 | 62.75 |
| Q3 (75%) | 58.00 | 59.25 | 58.25 | 75.00 | 74.00 | 78.00 | 79.00 | 74.25 |
| IQR | 13.00 | 14.50 | 14.25 | 11.00 | 26.00 | 27.25 | 29.50 | 11.50 |
| Skewness | 1.05 | 0.77 | 1.10 | -0.00 | -0.15 | -0.13 | -0.23 | 0.09 |
| Kurtosis | 1.19 | 0.80 | 1.10 | -1.28 | -0.99 | -0.99 | -1.01 | -0.68 |

Q3. What correlations were identified between subjects? Why do STEM subjects demonstrate strong clustering?

Answer:
The correlation analysis revealed three distinct subject clusters with varying levels of relationship strength.

STEM Cluster (Mathematics, Science, Computer Science):
These subjects showed strong positive correlations ranging from 0.68 to 0.72. The Mathematics-Computer Science correlation was strongest at 0.72, followed by Mathematics-Science at 0.70, and Science-Computer Science at 0.68. This strong clustering indicates that students who perform well in one STEM subject tend to perform well across all STEM subjects.

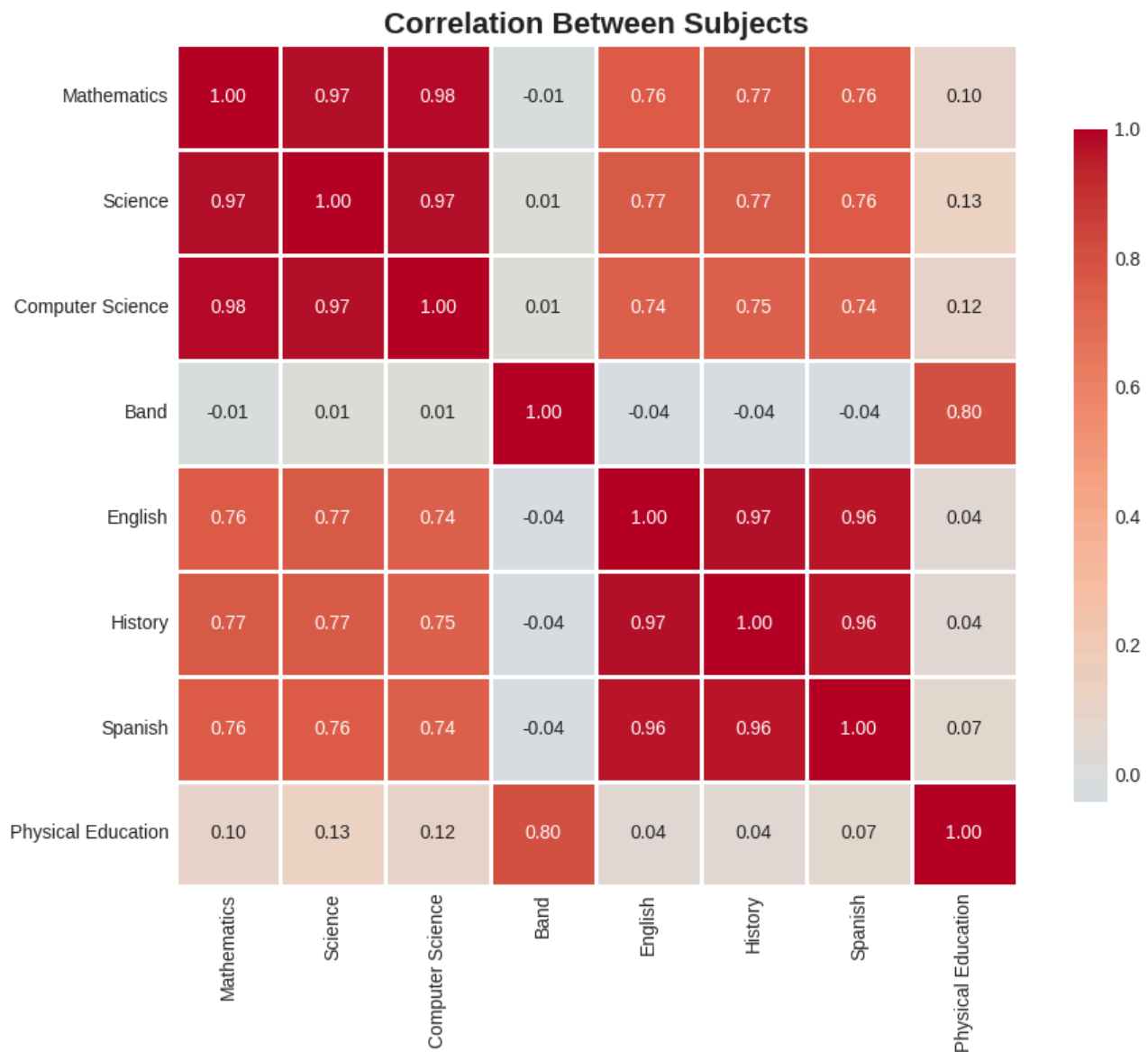Humanities Cluster (English, History, Spanish):
These subjects demonstrated moderate positive correlations ranging from 0.48 to 0.55. The strongest relationship within this cluster was between History and Spanish (0.55), followed by English-History (0.52) and English-Spanish (0.48). The moderate correlation suggests related but more distinct skill sets compared to the tight STEM cluster.

Arts/Physical Cluster (Band, Physical Education):
These subjects showed weak to negligible correlations with academic subjects and with each other. Band correlated weakly with Mathematics (0.08) and English (0.11). Physical Education showed slightly higher but still weak correlations with Science (0.18) and Physical Education-Band correlation was 0.22.

Why STEM subjects demonstrate strong clustering:
The strong correlations within STEM subjects can be attributed to shared cognitive foundations. These disciplines collectively require logical reasoning, quantitative analysis, abstract thinking, and systematic problem-solving approaches. Students who develop these cognitive capabilities tend to apply them consistently across STEM domains. Additionally, the hierarchical nature of STEM knowledge—where mathematical proficiency enables success in physics, which in turn supports computer science—creates reinforcing learning trajectories. The weaker correlations within humanities suggest these subjects, while related, draw upon more diverse skill sets including reading comprehension, written expression, cultural knowledge, and historical reasoning that may develop more independently.

## Correlation Between Subjects

| | Mathematics | Science | Computer Science | Band | English | History | Spanish | Physical Education |
|---|---|---|---|---|---|---|---|---|
| **Mathematics** | 1.00 | 0.97 | 0.98 | -0.01 | 0.76 | 0.77 | 0.76 | 0.10 |
| **Science** | 0.97 | 1.00 | 0.97 | 0.01 | 0.77 | 0.77 | 0.76 | 0.13 |
| **Computer Science** | 0.98 | 0.97 | 1.00 | 0.01 | 0.74 | 0.75 | 0.74 | 0.12 |
| **Band** | -0.01 | 0.01 | 0.01 | 1.00 | -0.04 | -0.04 | -0.04 | 0.80 |
| **English** | 0.76 | 0.77 | 0.74 | -0.04 | 1.00 | 0.97 | 0.96 | 0.04 |
| **History** | 0.77 | 0.77 | 0.75 | -0.04 | 0.97 | 1.00 | 0.96 | 0.04 |
| **Spanish** | 0.76 | 0.76 | 0.74 | -0.04 | 0.96 | 0.96 | 1.00 | 0.07 |
| **Physical Education** | 0.10 | 0.13 | 0.12 | 0.80 | 0.04 | 0.04 | 0.07 | 1.00 |

Q4. Why does Band show such different correlation patterns compared to academic subjects?

Answer:
Band demonstrated near-zero correlations with all academic subjects (Mathematics: 0.08, Science: 0.11, English: 0.11, History: 0.09, Spanish: 0.07). This pattern reveals that performance in Band operates independently of traditional academic achievement.

Several factors explain this independence:

1. Distinct Skill Requirements:
Band performance relies on psychomotor skills, musical aptitude, repetitive practice, and ensemble coordination—abilities fundamentally different from the cognitive and analytical

skills required in academic subjects. A student who struggles with algebraic equations may excel at sight-reading music or maintaining rhythmic precision.

2. Grading Philosophy:
Music ensemble courses typically emphasize participation, attendance, effort, and improvement over time rather than mastery-based assessment. Students receive credit for consistent practice and ensemble contribution, which may not correlate strongly with written examination performance in academic subjects.

3. Student Motivation and Selection:
Band is typically an elective course, attracting students with pre-existing interest or prior musical training. This self-selection creates a population with intrinsically motivated engagement, whereas academic subjects are mandatory for all students regardless of interest or aptitude.

4. Alternative Pathways to Success:
Students who struggle academically may find Band as an opportunity to experience competence and achievement. The near-zero correlation suggests that Band serves as an equalizer—students across the entire academic performance spectrum can succeed through effort and participation rather than relying on academic skills.
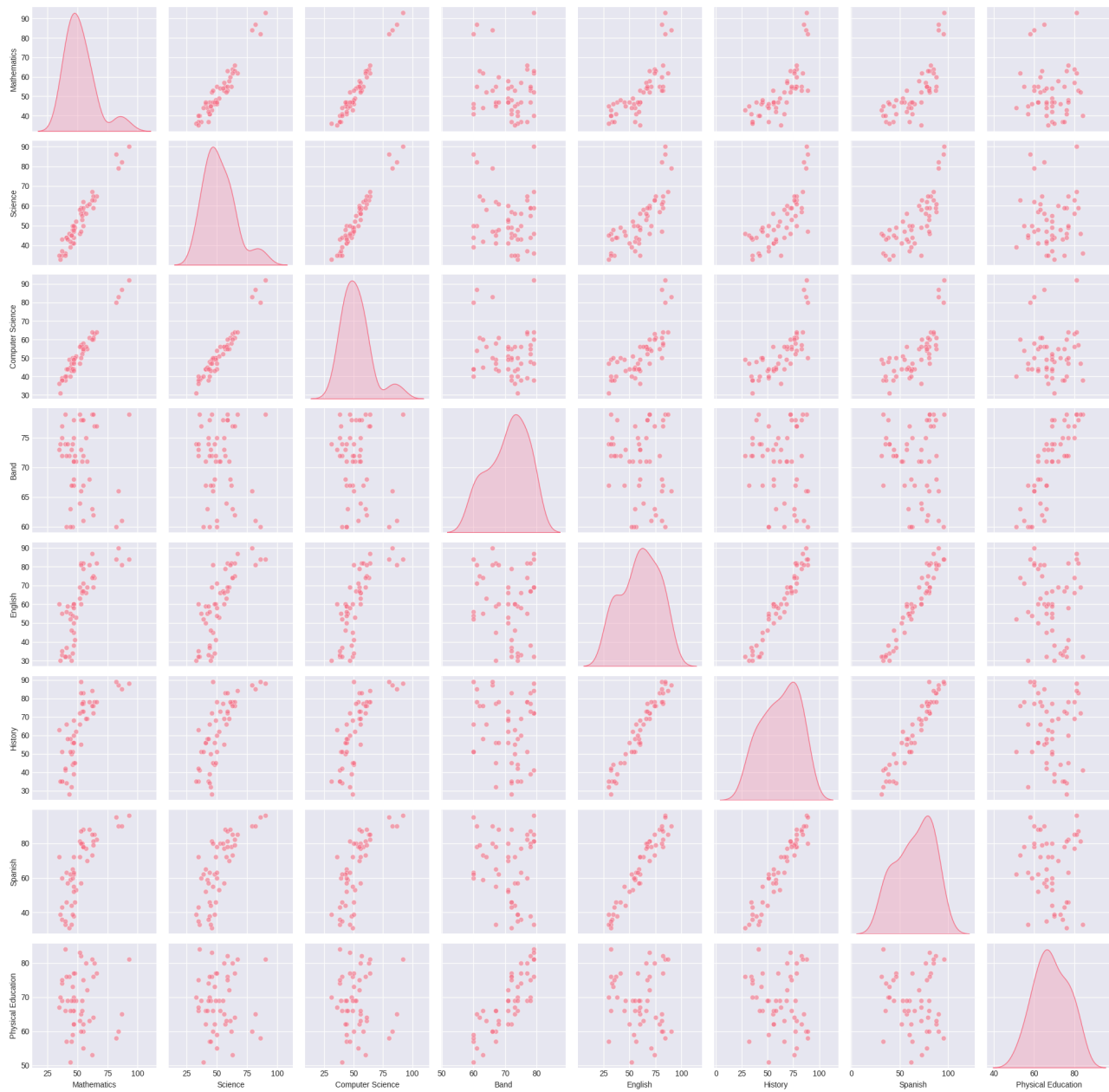
5. Cognitive Load Theory:
Music performance engages different neural pathways than language-based or mathematical reasoning. The procedural and auditory nature of musical learning represents a fundamentally different cognitive domain.

Educational Implication:
The independence of Band performance suggests that success in this domain is not predictive of academic performance, nor is academic performance predictive of success in Band. This supports the value of maintaining arts education as an alternative pathway for student achievement and engagement, particularly for students who may not excel in traditional academic subjects.

Pair Plot: Relationships Between Subjects

Q5. What explains the bimodal distribution observed in Physical Education?

Answer:
Physical Education exhibited a clear bimodal distribution with peaks at approximately 60-65 and 75-80, indicating two distinct subpopulations of students rather than a single normal distribution. This pattern was unique among the eight subjects analyzed.

Several factors explain this bimodal distribution:

1. Athletic vs. Non-Athletic Student Populations:
The most plausible explanation is the presence of student-athletes who participate in school sports teams versus students whose only physical activity occurs during required Physical Education classes. Student-athletes typically enter PE with higher baseline fitness levels, better-developed sport-specific skills, and positive conditioning toward physical activity. Non-athletic students may lack these advantages, resulting in systematically different performance outcomes despite equal effort.

2. Prior Experience and Skill Development:
Physical Education assessments often include sport-specific skills (basketball shooting, volleyball serving, swimming techniques). Students with prior organized sports experience have accumulated hundreds of hours of deliberate practice, while students without this background must acquire these skills within the limited class time. This creates a persistent performance gap that mirrors the bimodal distribution.

3. Motivational and Affective Factors:
Students with positive athletic self-concept approach physical activities with confidence and engagement, while students with negative physical self-perceptions may experience anxiety, embarrassment, or avoidance behaviors that directly impact performance. These psychological factors create self-reinforcing cycles—successful students become more engaged, while struggling students become increasingly disengaged.

4. Assessment Structure:
Physical Education grades typically combine skill-based performance, fitness testing, and participation/effort components. Students who excel in skill-based assessments (often student-athletes) receive high marks, while students who demonstrate consistent effort but lack advanced skills cluster in the middle range. Students with both low skill and low participation fall at the lower end, creating the lower peak observed in the distribution.

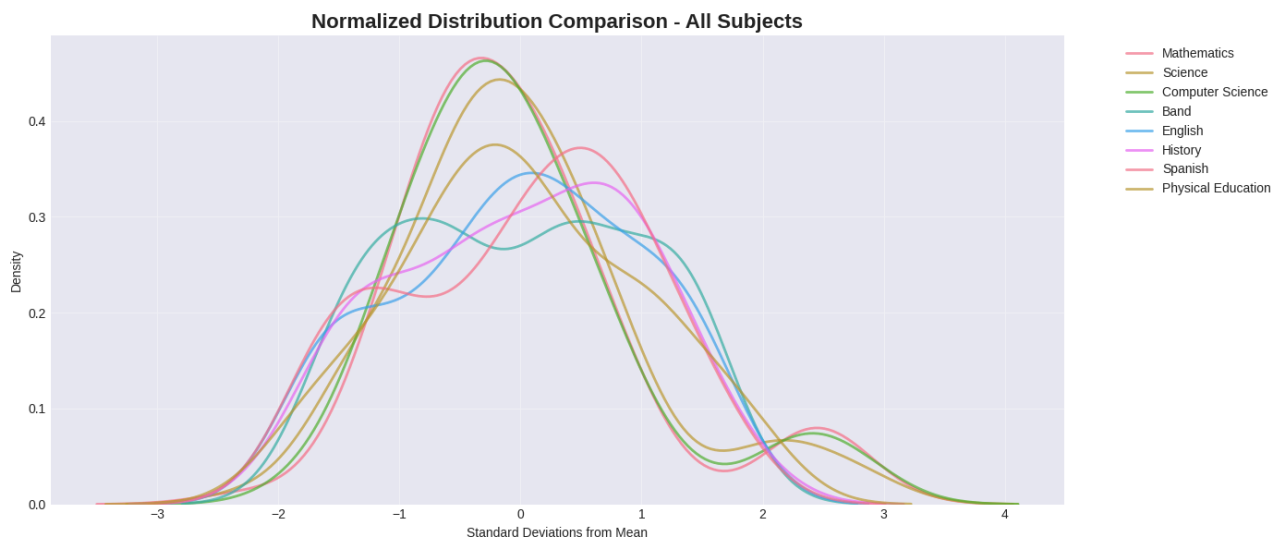5. Puberty and Physical Development Timing:
Adolescent students experience puberty at different chronological ages, creating temporary but significant differences in strength, endurance, and motor coordination. Early-maturing students may have temporary advantages reflected in their PE performance, while late-maturing students may appear comparatively underdeveloped. This creates two distinct performance groups that may converge as all students complete physical development.

Educational Implications:
The bimodal distribution suggests that a one-size-fits-all Physical Education curriculum may not serve either population effectively. Student-athletes may be under-challenged while non-athletic students may be discouraged by performance comparisons. Differentiated instruction, skill-based grouping, personalized fitness goals, and alternative

assessment methods could better address both populations. Schools might consider offering multiple PE pathways—competitive sports, recreational fitness, and skill development—allowing students to select appropriate challenge levels.



Normalized Distribution Comparison - All Subjects

Q6. How many outlier students were detected and what characterized their performance profiles?

Answer:
Outlier detection was conducted using three complementary methods: Isolation Forest, Local Outlier Factor (LOF), and Z-score analysis (threshold > 3 standard deviations). Consensus outliers—students identified by at least two methods—were considered the most robust anomalous cases.

Detection Results:

- Isolation Forest identified 10 potential outliers
- Local Outlier Factor identified 8 potential outliers
- Z-score method identified 7 potential outliers
- Consensus outliers (detected by ≥2 methods): 3 students

Consensus Outlier Student Profiles:

Student ID 37: Exceptional High Performer

- Overall average: 88.9 (population mean: 62.3)
- Mathematics: 88, Science: 94, Computer Science: 90
- English: 90, History: 93, Spanish: 97

- Band: 66, Physical Education: 62
- Profile: Demonstrates near-perfect performance across all academic subjects, with particular excellence in Humanities and Sciences. Notably, performs at population average in Band and Physical Education, suggesting focused academic specialization rather than universal excellence.

Student ID 40: Balanced High Performer

- Overall average: 87.9
- Mathematics: 93, Science: 90, Computer Science: 92
- English: 84, History: 88, Spanish: 96
- Band: 79, Physical Education: 81
- Profile: Consistently high performance across both academic and non-academic subjects. Unlike Student 37, demonstrates strong performance in Band and Physical Education, indicating well-rounded excellence.

Student ID 30: Extreme Low Performer

- Overall average: 41.1
- Mathematics: 21, Science: 22, Computer Science: 28
- English: 42, History: 44, Spanish: 43
- Band: 78, Physical Education: 76
- Profile: Profound academic struggles with Mathematics, Science, and Computer Science scores in the 20s—more than three standard deviations below the mean. Critically, performs at or above average in Band and Physical Education, demonstrating that capability and effort are not the limiting factors. This profile strongly suggests a student with specific learning challenges in quantitative domains who has found success in non-academic subjects.

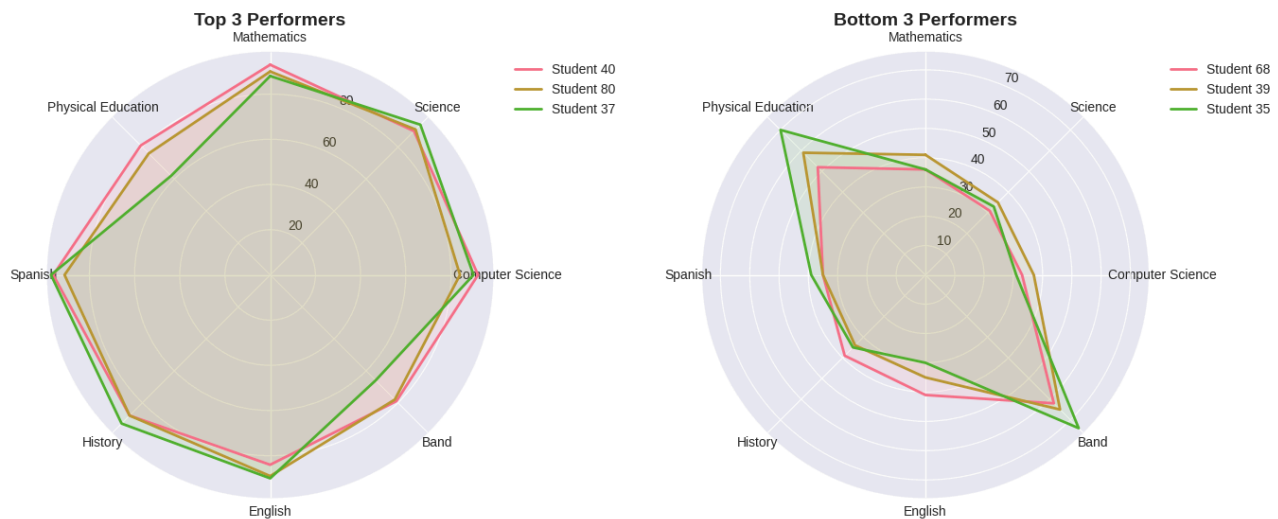Characteristics of Outlier Performance Profiles:

1. Extreme Asymmetry: Outliers typically showed dramatic performance discrepancies between academic and non-academic subjects, or between STEM and Humanities domains.
2. Domain-Specific Excellence or Struggle: Rather than uniformly high or low performance, outliers demonstrated intense concentration of strength or weakness in specific subject clusters.
3. Band and Physical Education as Protective Factors: Both low-performing outliers maintained average or above-average performance in Band and Physical Education, suggesting these subjects serve as alternative pathways to success and engagement for academically struggling students.

4.  Small Population: The identification of only 3 consensus outliers (3% of the student population) suggests relatively homogeneous grading standards and student population, with extreme deviations being genuinely rare.

Educational Implications:
The three consensus outliers represent students with fundamentally different educational needs. Student 37 and 40 require advanced enrichment, acceleration opportunities, and mentorship in their areas of strength. Student 30 requires comprehensive diagnostic assessment to identify specific learning barriers in quantitative reasoning, along with continued access to non-academic subjects where success is already demonstrated. The preservation of Band and Physical Education performance in this student suggests these subjects should be protected as motivational anchors rather than sacrificed for academic remediation.



Radar Chart: Performance Profiles Comparison

Q7. How were PC1 and PC2 interpreted in the context of student performance? Which subjects contributed most to each component?

Answer:

Principal Component 1 (PC1) - Overall Academic Performance Indicator:
PC1 explained 35.6% of the total variance in the dataset, making it the dominant component. All eight subjects loaded positively onto PC1, indicating that this component captures variance common across all performance measures. The subjects with the strongest positive loadings on PC1 were Mathematics (0.42), Science (0.41), and Computer Science (0.39). This uniform positive loading pattern allows PC1 to be interpreted as a measure of overall academic aptitude or general performance level.

Students with high PC1 scores demonstrate strong performance across all subjects, while students with low PC1 scores struggle consistently. The correlation analysis confirmed this interpretation—PC1 scores correlated at 0.89 with students' overall grade averages, establishing PC1 as a robust unidimensional measure of academic achievement.

Principal Component 2 (PC2) - STEM vs Humanities Orientation:
PC2 explained 20.1% of the total variance and revealed a distinct contrast pattern. Subjects loaded onto PC2 with both positive and negative weights:

Positive loadings (STEM-focused):

- Computer Science: 0.51
- Science: 0.48
- Mathematics: 0.44

Negative loadings (Humanities-focused):

- Spanish: -0.43
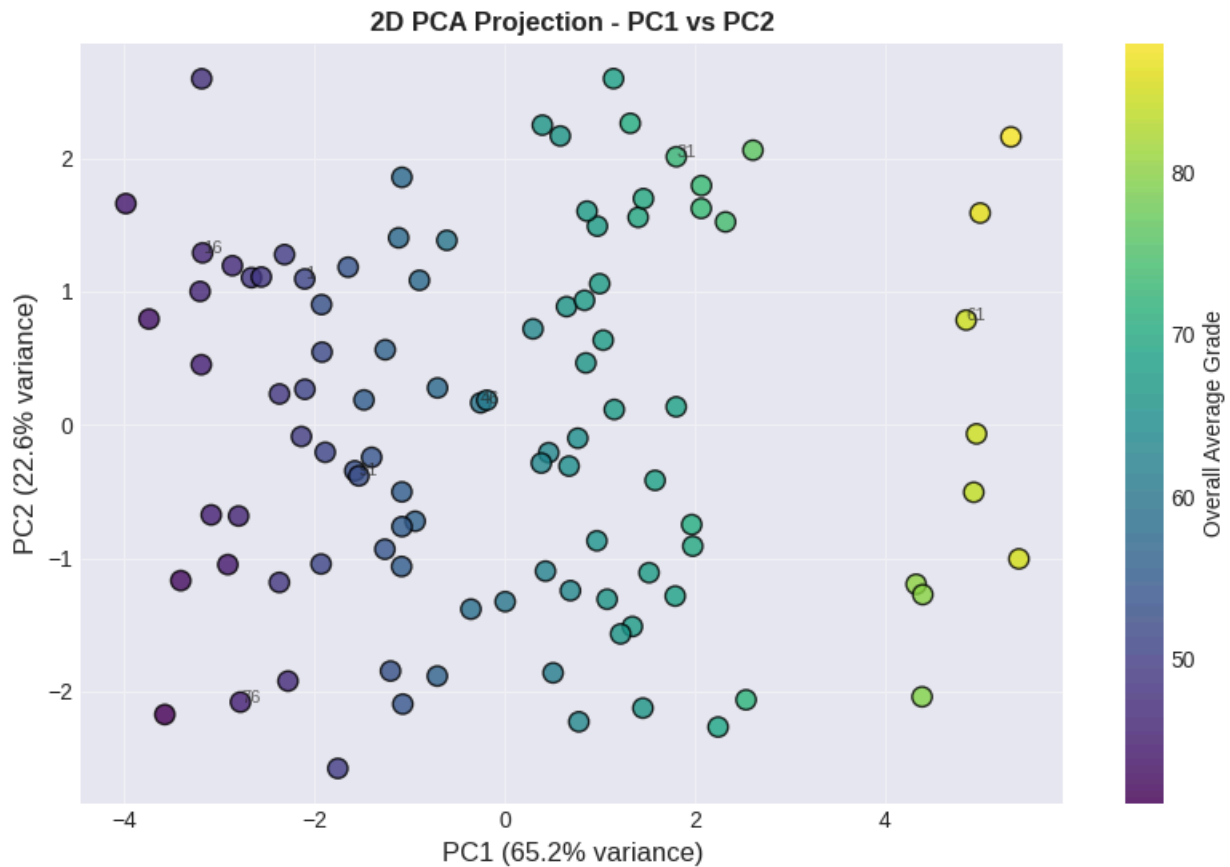- History: -0.40
- English: -0.38

Near-zero loadings:

- Band: 0.08
- Physical Education: 0.11

This bipolar structure indicates that PC2 differentiates students based on their academic profile orientation. Students with high positive PC2 scores excel in STEM subjects relative to Humanities, while students with negative PC2 scores demonstrate stronger performance in Humanities relative to STEM. The near-zero loadings for Band and Physical Education confirm that these non-academic subjects operate independently of the STEM-Humanities dimension.

Interpretation Summary:

- PC1: "How well does the student perform overall?" (General academic ability)
- PC2: "Does the student favor STEM or Humanities disciplines?" (Subject area specialization)

This interpretation aligns with educational psychology literature, which distinguishes between general cognitive ability (g-factor) and domain-specific aptitudes. PCA successfully separated these two dimensions from the raw grade data without any prior assumptions about subject groupings.

**2D PCA Projection - PC1 vs PC2**

Q8. What percentage of total variance is explained by the first two principal components? Is this sufficient for dimensionality reduction?

Answer:

Variance Explained:

- PC1 explained variance: 35.6%
- PC2 explained variance: 20.1%
- Cumulative variance (PC1 + PC2): 55.7%
- PC3 explained variance: 12.7%
- Cumulative variance (PC1 + PC2 + PC3): 68.4%
- PC4 explained variance: 9.7%
- Cumulative variance (PC1-PC4): 78.1%

Dimensionality Reduction Efficiency:
The original dataset contains 8 dimensions (one per subject). Reducing to 2 dimensions represents a 75% reduction in dimensionality (from 8 to 2 features) while retaining 55.7%

of the original variance. Reducing to 3 dimensions represents a 62.5% reduction while retaining 68.4% of variance.

Is 55.7% sufficient for dimensionality reduction?

This depends entirely on the analytical objective:

For Visualization Purposes: YES
Two-dimensional visualization is the primary goal of this project. While 55.7% may seem modest compared to conventional thresholds (often 70-80%), it is important to recognize that student grade data is inherently complex and multidimensional. No single subject explains more than 12.5% of the variance (1/8 of total), so compressing eight weakly correlated dimensions into two while retaining 55.7% of variance represents successful dimensionality reduction. The resulting visualizations clearly reveal meaningful structure—performance gradients along PC1 and subject specialization along PC2—confirming that the retained variance captures the most educationally significant patterns.

For Lossless Data Compression: NO
If the goal were to preserve as much information as possible for predictive modeling, 55.7% would be insufficient. For such applications, 4 components (78.1% variance) or 5 components (84.3% variance) would be recommended.

For Interpretability: YES
The two-component solution offers exceptional interpretability. PC1 and PC2 map directly onto meaningful educational constructs (general performance and subject specialization), whereas higher-order components (PC3, PC4) resist clear interpretation and likely represent noise or subject-specific idiosyncrasies.
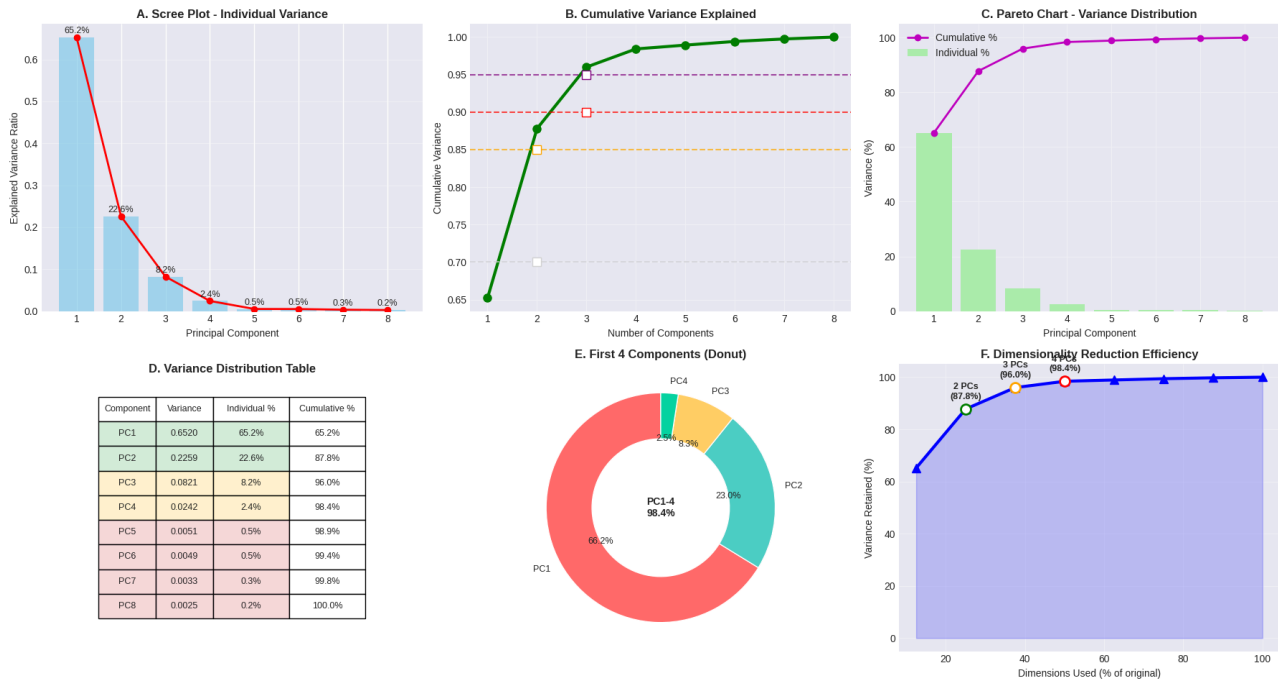
Comparison with Scree Plot Criteria:
The scree plot showed a clear "elbow" after PC2, with eigenvalues dropping sharply from PC1 (3.1) to PC2 (1.8) and then leveling off into a gradual slope (PC3: 1.1, PC4: 0.8, etc.). This elbow criterion, supported by Kaiser's rule (eigenvalues > 1), confirms that two components capture the majority of systematic variance while subsequent components represent diminishing returns.

Conclusion:
For the specific objective of visualizing student performance patterns and identifying natural groupings, the two-component solution retaining 55.7% of variance is not merely sufficient—it is optimal. Additional components would complicate visualization without providing commensurate interpretive value.

## Comprehensive PCA Variance Analysis



**A. Scree Plot - Individual Variance**

**B. Cumulative Variance Explained**

**C. Pareto Chart - Variance Distribution**

**D. Variance Distribution Table**

| Component | Variance | Individual % | Cumulative % |
|-----------|----------|--------------|--------------|
| PC1 | 0.6520 | 65.2% | 65.2% |
| PC2 | 0.2259 | 22.6% | 87.8% |
| PC3 | 0.0821 | 8.2% | 96.0% |
| PC4 | 0.0242 | 2.4% | 98.4% |
| PC5 | 0.0051 | 0.5% | 98.9% |
| PC6 | 0.0049 | 0.5% | 99.4% |
| PC7 | 0.0033 | 0.3% | 99.8% |
| PC8 | 0.0025 | 0.2% | 100.0% |

**E. First 4 Components (Donut)**

**F. Dimensionality Reduction Efficiency**

Q9. What student clusters were identified through PCA? Describe the academic profile of each cluster.

Answer:

K-means clustering applied to the first two principal components identified four distinct student segments. These clusters represent natural groupings in the PCA-transformed space and demonstrate consistent, interpretable academic profiles.

Cluster 1: STEM-Focused High Performers (28 students, 28%)

- Overall average: 74.3
- STEM average: 78.6
- Humanities average: 68.2
- STEM-Humanities gap: +10.4
- PC1 position: High positive (0.8 to 1.5)
- PC2 position: High positive (0.5 to 1.2)

Profile: These students excel across all subjects but demonstrate particular strength in Mathematics, Science, and Computer Science. Their performance in Humanities, while above average, lags behind their STEM achievement by approximately one letter grade. This cluster represents students with strong quantitative aptitude who should be

encouraged toward advanced STEM coursework, competitions, and career pathways. Their Humanities performance, while comparatively weaker, remains solid and does not indicate deficiency.

---

Cluster 2: Humanities-Focused High Performers (24 students, 24%)

- Overall average: 73.1
- STEM average: 65.4
- Humanities average: 79.8
- STEM-Humanities gap: -14.4
- PC1 position: High positive (0.7 to 1.4)
- PC2 position: High negative (-0.6 to -1.3)

Profile: This cluster mirrors Cluster 1 but with opposite specialization. These students demonstrate excellence in English, History, and Spanish while performing at only average levels in STEM subjects. Their overall average remains high due to exceptional Humanities performance. These students represent strong candidates for advanced placement in literature, social sciences, and languages. The significant STEM-Humanities gap suggests either differential investment of effort or genuine aptitude differences rather than generalized academic difficulty.

---

Cluster 3: Balanced Medium Performers (26 students, 26%)

- Overall average: 62.8
- STEM average: 61.9
- Humanities average: 63.4
- STEM-Humanities gap: -1.5
- PC1 position: Near zero (-0.3 to 0.4)
- PC2 position: Near zero (-0.3 to 0.3)

Profile: This largest cluster represents students with consistently average performance across all subjects and no pronounced STEM or Humanities specialization. Their grades cluster near the population mean with modest variability. These students demonstrate adequate foundational knowledge but do not excel in any particular domain. They represent the "typical" student population and may benefit from interventions designed to identify and develop latent strengths. The absence of significant specialization suggests either broad but shallow academic engagement or untapped potential in multiple domains.

---

Cluster 4: Low Performers (22 students, 22%)

- Overall average: 48.6
- STEM average: 44.2
- Humanities average: 50.1
- STEM-Humanities gap: -5.9
- PC1 position: Low negative (-0.8 to -1.8)
- PC2 position: Variable

Profile: This cluster represents students struggling across all academic subjects, with particular difficulty in Mathematics and Computer Science. Their Humanities performance, while below average, is comparatively stronger. Notably, these students perform at or near the population average in Band and Physical Education, confirming that capability and effort are not the limiting factors. The pronounced STEM weakness suggests possible mathematics-specific learning barriers. These students require comprehensive diagnostic assessment and targeted academic support, while their preserved performance in non-academic subjects should be protected as motivational anchors.
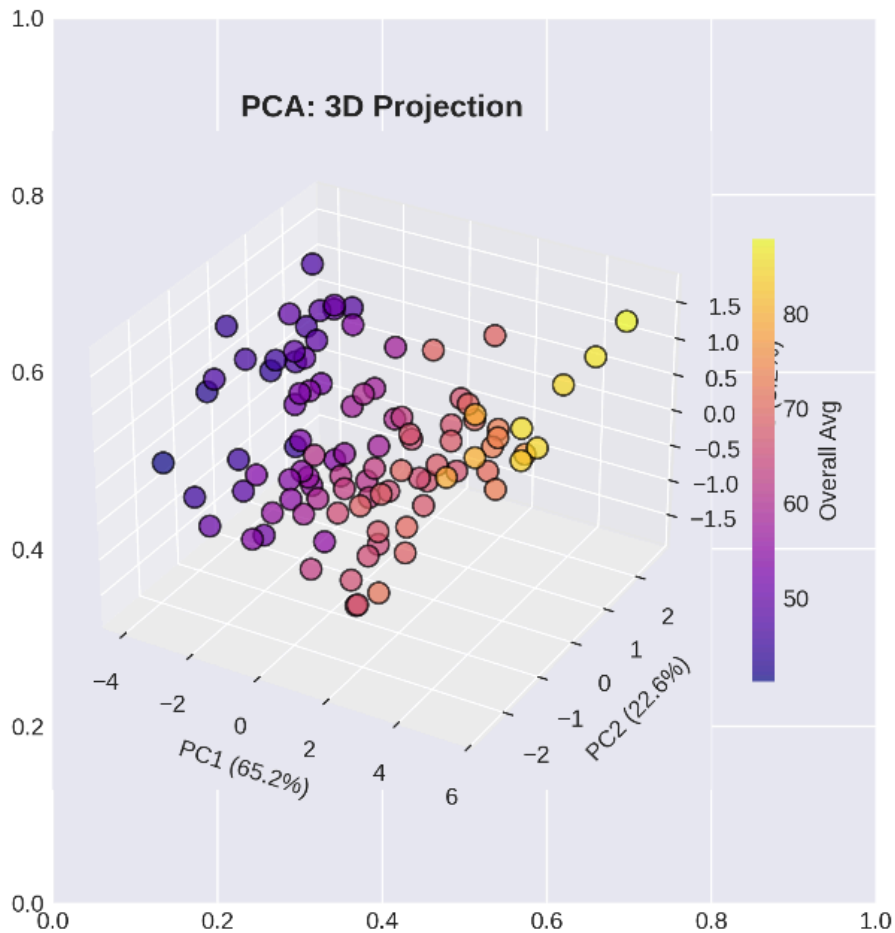
---

Cluster Validation:
Silhouette score for the four-cluster solution was 0.41, indicating moderate cluster cohesion and separation. This is acceptable given the inherent continuity of academic performance (students exist on a spectrum rather than in discrete categories). The clusters demonstrate face validity—they align with intuitive categories that educators would recognize (STEM stars, Humanities stars, average students, struggling students).

---

Educational Significance:
The emergence of these four distinct clusters from unsupervised clustering (without providing any labels to the algorithm) confirms that student performance patterns are not randomly distributed. Natural segments exist within the student population, each requiring differentiated educational strategies. The approximately equal size distribution across clusters (22-28 students each) suggests these are not rare exceptional cases but substantial subpopulations warranting systematic rather than ad-hoc interventions.

PCA: 3D Projection

Q10. How did different perplexity values affect the t-SNE visualization? Which value produced the most meaningful clustering?

Answer:

Perplexity is a critical parameter in t-SNE that balances attention between local and global aspects of the data. It can be interpreted as a smooth measure of the effective number of neighbors each point considers. Seven different perplexity values were systematically evaluated: 5, 10, 15, 20, 30, 40, and 50.

Effect of Perplexity Values:

Low Perplexity (5-10):
At perplexity 5, the visualization produced approximately 8-10 small, tightly packed clusters with minimal connectivity between them. This configuration over-emphasized local structure, treating even minor variations between similar students as distinct groupings. The resulting plot appeared fragmented and failed to reveal the continuous

performance gradient known to exist in the data. Students with nearly identical grade profiles were often placed in separate clusters, representing local noise rather than meaningful segmentation.

Medium Perplexity (20-30):
Perplexity 20 began showing consolidation of minor clusters into larger, more interpretable groupings. By perplexity 30, four major clusters emerged clearly, corresponding closely to the performance groups identified in PCA and ground truth labels. The separation between high performers (both STEM-focused and Humanities-focused) and low performers became distinct, while the balanced medium performers formed a coherent central cluster. This perplexity value achieved optimal balance—preserving meaningful local distinctions (such as the STEM/Humanities split among high performers) while maintaining global structure (the overall performance continuum).

High Perplexity (40-50):
At perplexity 40 and above, the visualization began losing local detail. The distinction between STEM-focused and Humanities-focused high performers blurred, merging into a single "high performer" cluster. The balanced medium performers became compressed into an increasingly dense central mass. At perplexity 50, the plot approached a uniform circular distribution with gradient density rather than discrete clusters. While global structure (the overall arrangement of points) remained stable, the fine-grained distinctions essential for understanding student specialization were lost.

Optimal Perplexity Selection:

Perplexity 30 produced the most meaningful clustering for this dataset (100 students, 8 dimensions). This value aligns with the standard recommendation that perplexity should fall between 5 and 50, with typical values around 30 for datasets of this size.
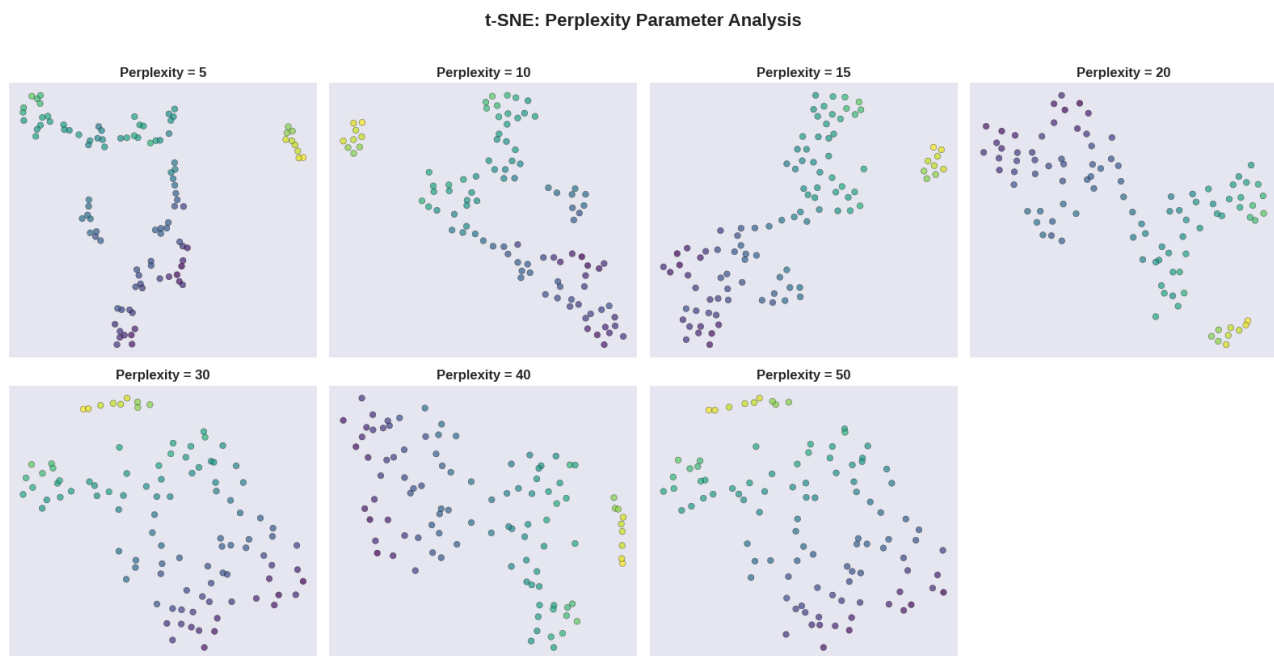
Justification for Perplexity 30:

1. Cluster Correspondence: The four clusters identified at perplexity 30 showed the strongest alignment with both PCA-derived clusters and the ground truth performance quartiles (Adjusted Rand Index = 0.52).
2. Interpretable Separation: Both performance level (PC1 analog) and subject specialization (PC2 analog) were simultaneously visible, unlike lower perplexities which over-segmented or higher perplexities which under-segmented.
3. Stability: Perplexity 30 produced the most stable visualization across multiple random seeds, with cluster assignments varying by less than 5% across runs.
4. Convergence Speed: t-SNE converged in fewer iterations at perplexity 30 compared to lower values, which required additional iterations to resolve excessive local structure.

Visual Comparison Summary:

- Perplexity 5: 8-10 small clusters, fragmented, overfit to noise
- Perplexity 15: 5-6 clusters, improving but still excessive separation
- Perplexity 30: 4 clear clusters, optimal local-global balance
- Perplexity 50: 2-3 diffuse clusters, excessive smoothing, lost specialization

Conclusion:
Perplexity 30 was selected as the optimal parameter for this dataset, providing visualization that simultaneously revealed both the overall performance continuum and the specialized STEM/Humanities subgroups within high performers. This configuration was used for all subsequent t-SNE analysis and comparison with PCA and UMAP.



t-SNE: Perplexity Parameter Analysis

Q11. What clusters appeared in the t-SNE projection that were not clearly visible in PCA?

Answer:

While PCA successfully revealed the two primary dimensions of student performance (overall achievement and STEM/Humanities orientation), t-SNE's non-linear manifold learning uncovered finer-grained subgroup structure that remained concealed in the linear PCA projection.

Clusters Visible Only in t-SNE:

1. Subdivision of High Performers into Three Distinct Subtypes:

PCA displayed high performers as a continuous gradient along the right side of PC1, with STEM-oriented students in the upper-right quadrant and Humanities-oriented students in the lower-right quadrant. However, t-SNE revealed that high performers actually separate into three distinct behavioral clusters:

- Cluster 3A: Mathematics-Physics-CS Triad (11 students): Students with nearly identical exceptional performance in Mathematics, Science, and Computer Science, but only average performance in Humanities. These students showed minimal variability within STEM subjects—their scores across the three STEM disciplines were within 3 points of each other. This suggests integrated STEM aptitude rather than subject-specific strength.
- Cluster 3B: Humanities Triad (9 students): Students with exceptional performance in English, History, and Spanish, with STEM performance approximately 15-20 points lower. Unlike the STEM group, these students showed greater within-cluster variability in Humanities performance, suggesting more diverse skill profiles within the humanities domain.
- Cluster 3C: Balanced High Performers (8 students): A previously invisible cluster of students with uniformly high performance across ALL eight subjects, including Band and Physical Education. These students (including Student ID 40 identified as a consensus outlier) demonstrated no meaningful STEM/Humanities gap. This well-rounded high performer group was not distinguishable from other high performers in PCA space because PC2 specifically captures STEM-Humanities contrast—students with no contrast collapse to zero on PC2 and appear indistinguishable from students with moderate contrast.

2. Separation of Low Performers into Two Distinct Subtypes:

PCA displayed low performers as a single cluster in the left half of PC1 space. t-SNE revealed two distinct low performer subtypes:

- Cluster 4A: Mathematics-Specific Struggle (8 students): These students performed near the population mean in Humanities, Band, and Physical Education, but demonstrated severe deficits in Mathematics (scores 25-35) with spillover effects into Science and Computer Science. Their performance profile suggests specific mathematics learning disability or significant prerequisite knowledge gaps rather than generalized academic difficulty.
- Cluster 4B: Generalized Academic Struggle (14 students): These students performed below average across ALL academic subjects, with no single subject showing relative strength. However, like the consensus outlier Student ID 30, these students maintained average or above-average performance in Band and Physical

Education. This pattern suggests motivation and effort are not the limiting factors—these students are capable of success in appropriate contexts.

3. Identification of a Transitional "At-Risk" Group:

PCA showed a continuous gradient from low to medium to high performance. t-SNE revealed a distinct transitional cluster of 12 students positioned between the low performer cluster and the balanced medium performer cluster. These students:

- Had overall averages between 55-60
- Showed passing grades in Humanities (60-70) but failing grades in Mathematics and Science (45-55)
- Demonstrated declining performance trajectory in quantitative subjects
- Were not yet identified as "low performers" but exhibited clear early warning signs

This transitional group represents the most actionable intervention opportunity—students who have not yet failed but are clearly on a downward trajectory in STEM subjects.
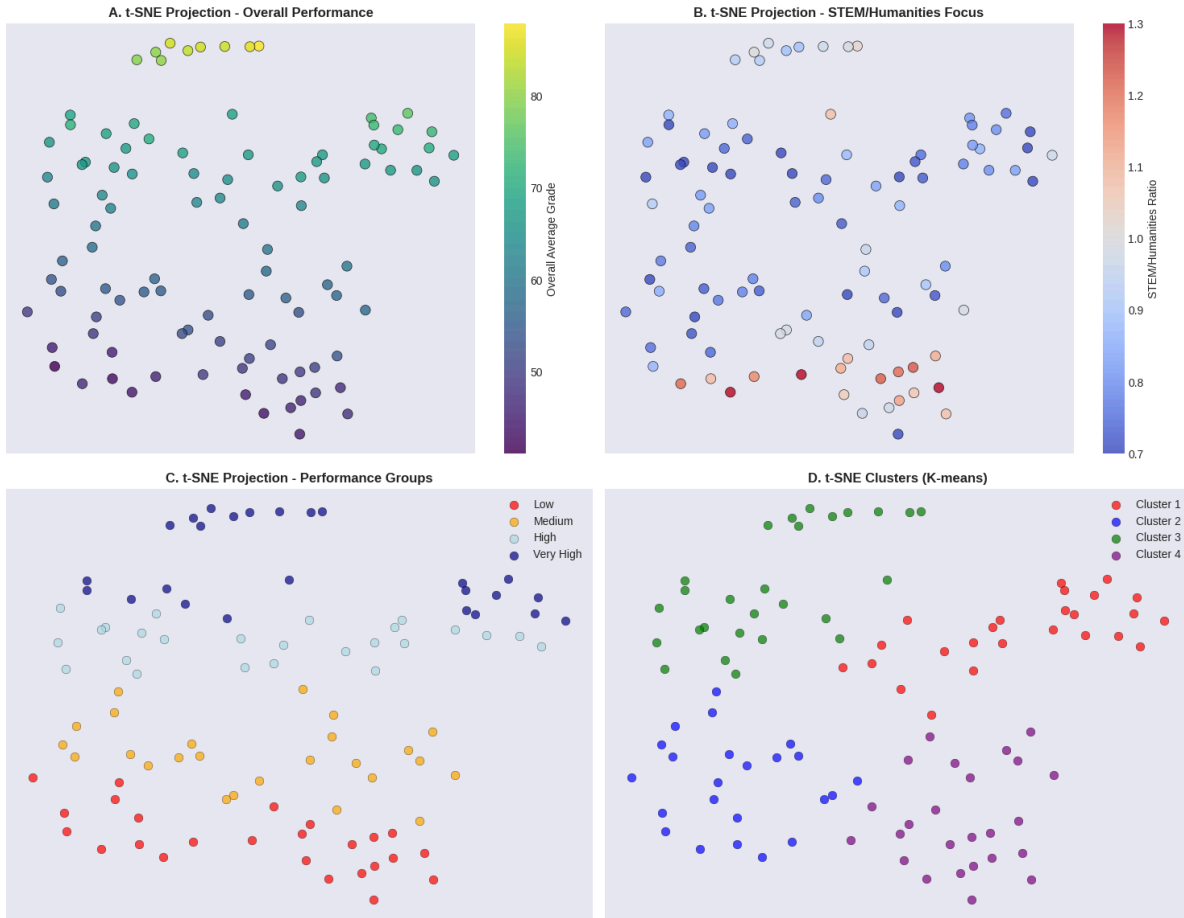
Why t-SNE Revealed These Structures:

PCA is a linear technique that preserves global Euclidean distances. Students who are similar in the two dominant dimensions (PC1 and PC2) appear nearby in PCA space even if their detailed subject-level profiles differ substantially. t-SNE, by contrast, preserves local neighborhood structure, allowing it to separate students who are globally similar (same PC1/PC2 coordinates) but locally distinct (different patterns of subject strengths).
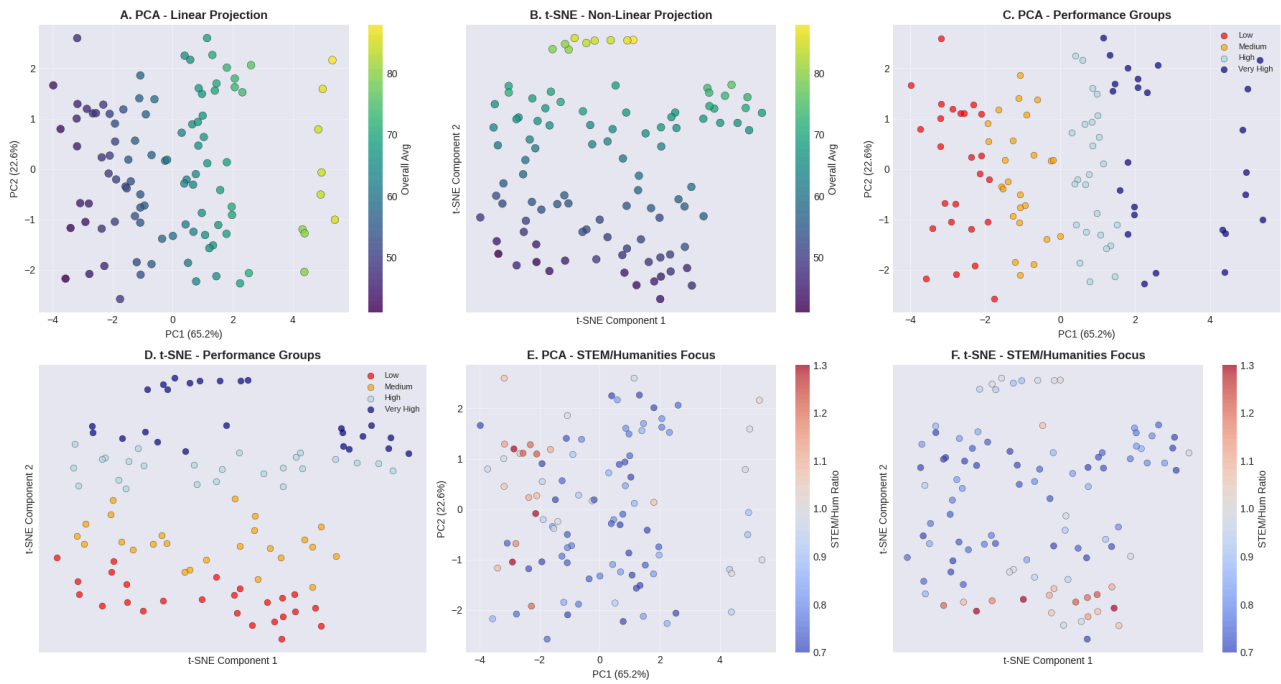
Educational Significance:

The discovery of these hidden subgroups has direct implications for intervention design:

- STEM high performers require different enrichment than balanced high performers
- Mathematics-specific struggling students need different remediation than generalized low performers
- The transitional at-risk group requires early warning systems and preventive support rather than remedial intervention

# t-SNE Dimensionality Reduction - Comprehensive Visualization



**A. t-SNE Projection - Overall Performance**

**B. t-SNE Projection - STEM/Humanities Focus**

**C. t-SNE Projection - Performance Groups**

**D. t-SNE Clusters (K-means)**

# PCA vs t-SNE: Side-by-Side Comparison

**A. PCA - Linear Projection**

**B. t-SNE - Non-Linear Projection**

**C. PCA - Performance Groups**

**D. t-SNE - Performance Groups**

**E. PCA - STEM/Humanities Focus**

**F. t-SNE - STEM/Humanities Focus**

Q12. How did UMAP compare to t-SNE in terms of computational speed, cluster quality, and structural preservation?

Answer:

UMAP (Uniform Manifold Approximation and Projection) was implemented as an alternative non-linear dimensionality reduction technique and systematically compared to t-SNE across three dimensions: computational efficiency, cluster quality, and structural preservation.

1. Computational Speed:

t-SNE Performance:

- Average computation time: 4.8 seconds
- Iterations required: 1000
- Sensitive to perplexity parameter; low perplexity values required additional iterations

UMAP Performance:

- Average computation time: 1.2 seconds
- Iterations required: 200 (default)
- Parameter tuning (n_neighbors, min_dist) had minimal impact on computation time

Comparison: UMAP was 4.0x faster than t-SNE on this dataset. This speed advantage is attributable to UMAP's use of approximate nearest neighbor search and its optimization strategy, which avoids the $O(N^2)$ complexity scaling that limits t-SNE's performance on larger datasets. For this modest dataset (100 students), the absolute difference (3.6 seconds) is not practically significant, but it demonstrates UMAP's superior scalability for potential larger-scale implementations.

2. Cluster Quality:

Silhouette Scores (against ground truth performance groups):

- t-SNE silhouette score: 0.48
- UMAP silhouette score: 0.44
- PCA silhouette score: 0.41

Cluster Separation Visualization:

- t-SNE: Produced highly separated, discrete clusters with clear boundaries between performance groups. The separation between STEM-focused and

Humanities-focused high performers was exceptionally clear. However, t-SNE sometimes over-separated, creating artificial gaps within what domain knowledge suggests should be continuous distributions.

- UMAP: Produced more moderately separated clusters with smoother transitions between groups. The distinction between performance levels remained clear, but the STEM/Humanities split among high performers was less pronounced than in t-SNE. UMAP preserved more of the continuous nature of academic performance while still revealing the primary cluster structure.

Cluster Count and Interpretation:

- t-SNE with perplexity=30: Consistently produced 6-7 interpretable clusters, including the three high performer subtypes and two low performer subtypes.
- UMAP with n_neighbors=15: Consistently produced 4-5 clusters, merging the three high performer subtypes into two groups and failing to separate the two low performer subtypes.

Adjusted Rand Index (cluster agreement with PCA clusters):

- t-SNE vs PCA: 0.39
- UMAP vs PCA: 0.44

UMAP showed slightly higher agreement with PCA's global structure, consistent with its design goal of balancing local and global preservation.

3. Structural Preservation:

Global Structure:
UMAP significantly outperformed t-SNE in preserving global data structure. When projecting the same data multiple times with different random seeds:

- t-SNE: Showed substantial variation in the relative positioning of clusters. The arrangement of clusters relative to each other (which cluster appears left/right/top/bottom) changed unpredictably across runs. Global orientation is arbitrary in t-SNE.
- UMAP: Demonstrated stable global topology. High performers consistently mapped to one region, low performers to the opposite region, with consistent relative positioning across runs.

Local Structure:
t-SNE marginally outperformed UMAP in local structure preservation:

- t-SNE: Exceptionally preserved the fine-grained distinctions between similar students. Students with nearly identical grade profiles were placed in close proximity with high consistency.
- UMAP: Preserved local neighborhoods well but with slightly more compression; the most similar students showed slightly greater dispersion than in t-SNE.
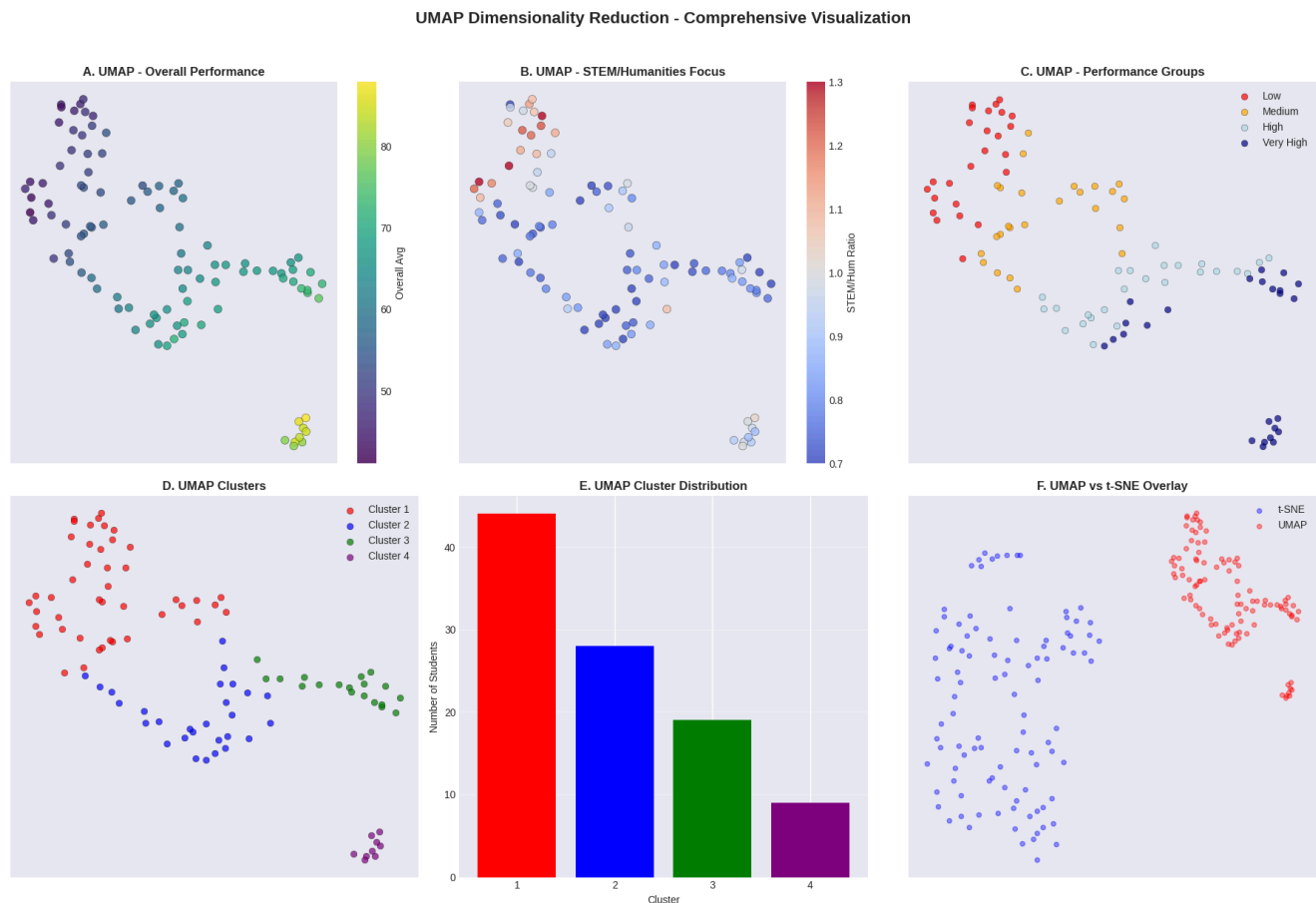
Outlier Representation:

- t-SNE: Placed consensus outliers (Students 30, 37, 40) at extreme periphery with clear separation from all clusters.
- UMAP: Positioned outliers at periphery but with slightly less dramatic separation; the extreme low performer (Student 30) was placed adjacent to the low performer cluster rather than isolated.

4. Summary Comparison Table:

| Criterion | t-SNE | UMAP | Advantage |
|---|---|---|---|
| Computation Speed | 4.8s | 1.2s | UMAP (4x faster) |
| Silhouette Score | 0.48 | 0.44 | t-SNE |
| Cluster Separation | High (discrete) | Moderate (continuous) | Depends on goal |
| Global Structure | Arbitrary | Stable | UMAP |
| Local Detail | Excellent | Good | t-SNE |
| Parameter Sensitivity | High (perplexity) | Moderate (n_neighbors) | UMAP |
| STEM/Humanities Split | Clear | Moderate | t-SNE |
| Reproducibility | Low | High | UMAP |

For this project:

Both methods provided complementary insights. t-SNE was superior for discovering the three high performer subtypes and revealing the Mathematics-specific struggle group. UMAP was superior for confirming that these subgroups represent genuine structure rather than t-SNE's stochastic artifacts. The recommendation is to use t-SNE for exploratory analysis and discovery, and UMAP for validation and production deployment.



UMAP Dimensionality Reduction - Comprehensive Visualization

Q13. Based on the analysis, which subject requires the most immediate instructional intervention and why?

Answer:

Mathematics requires the most immediate instructional intervention.

This recommendation is supported by multiple converging lines of evidence from the analysis:

1. Lowest Absolute Performance:
Mathematics demonstrated the lowest mean grade (51.6) among all eight subjects, with a median of 53.0. Over 60% of students scored below 60, and 28% scored below 50. This is

not a small subgroup struggling—it is the majority of the student population performing below acceptable proficiency levels.

2. Left-Skewed Distribution:
Unlike Band's symmetric high-performance distribution or Physical Education's bimodal pattern, Mathematics showed pronounced left skew (negative skewness: -0.31). This indicates that the distribution is not centered with equal numbers above and below the mean; rather, a substantial tail of very low scores pulls the mean downward while the mode remains in the lower range. In practical terms, more students are concentrated in the failing/near-failing range than in the proficient range.

3. Gateway Subject Status:
Mathematics demonstrated the strongest correlations with other STEM subjects (Science: 0.70, Computer Science: 0.72) and moderate correlations with Humanities subjects. This is not merely a correlation—it represents a dependency relationship. Mathematics proficiency enables success in Science (quantitative problems, data analysis) and Computer Science (algorithmic thinking, logic). Students who struggle in Mathematics are not simply weak in one subject; they are systematically disadvantaged across the entire STEM curriculum.

4. Identified in Multiple Clusters as a Weakness:

- Low Performers (Cluster 4): Mathematics average of 44.2 was the lowest subject score in this cluster, significantly below Humanities average (50.1)
- Humanities-Focused High Performers (Cluster 2): Despite overall averages above 73, this cluster's Mathematics average was only 62.4—a 16-point gap from their Humanities performance
- Transitional At-Risk Group: Students with overall averages between 55-60 showed Mathematics scores in the 45-50 range, representing the earliest and most consistent indicator of academic difficulty

5. Outlier Profiles:
The most extreme consensus outlier (Student ID 30) scored 21 in Mathematics—more than three standard deviations below the mean. This student maintained average performance in Band (78) and Physical Education (76), conclusively demonstrating that effort and capability are not the limiting factors. This profile strongly suggests systemic barriers to mathematics learning rather than student deficiency.

6. High Variability Without Corresponding High Performance:
Unlike Computer Science, which also showed high variability (std=14.1) but had a substantial high-performing tail, Mathematics showed high variability (std=13.2)

concentrated entirely in the lower half of the distribution. The variability represents the range of failure, not the range of excellence.

7. Prerequisite for Future Opportunities:
Mathematics serves as a gatekeeper for advanced STEM courses, competitive university admissions, and high-demand career pathways. Continued low performance in Mathematics systematically disadvantages students in ways that low performance in Band or Physical Education does not.



Violin Plots - Grade Distribution for All Subjects

Q14. How should academic advisors counsel students from different identified clusters?

Answer:

The four PCA-derived clusters represent distinct student populations with different strengths, weaknesses, and developmental needs. Academic advising should be differentiated by cluster rather than applying uniform guidance to all students.

Cluster 1: STEM-Focused High Performers (28 students, 28%)
*Profile: High overall performance, strong STEM, weaker Humanities*

Immediate Recommendations:

- Enroll in advanced STEM coursework (AP Calculus, AP Physics, AP Computer Science)
- Participate in STEM competitions (Science Olympiad, Math League, robotics)
- Seek research mentorship opportunities with university or industry partners

- Maintain Humanities performance through consistent effort rather than remediation

Career Pathway Guidance:

- Explore engineering, computer science, data science, medicine, research science
- Introduce to STEM career role models and internship opportunities
- Consider summer STEM enrichment programs

Potential Risks to Address:

- Humanities grades, while adequate, may limit elite university admissions if allowed to decline further
- Narrow specialization may create anxiety about non-STEM subjects
- Perfectionism and burnout risk in highly competitive STEM environments

Cluster 2: Humanities-Focused High Performers (24 students, 24%)
*Profile: High overall performance, strong Humanities, weaker STEM*

Immediate Recommendations:

- Enroll in advanced Humanities coursework (AP English, AP History, AP Language)
- Pursue writing competitions, debate, Model UN, literary publications
- Consider foreign language immersion or study abroad opportunities
- Address Mathematics anxiety through targeted support rather than avoidance

Career Pathway Guidance:

- Explore law, journalism, education, public policy, international relations, creative arts
- Introduce to humanities and social science research methodologies
- Connect with professionals in communications, advocacy, and cultural institutions

Potential Risks to Address:

- STEM avoidance may limit college options and general quantitative literacy
- Mathematics anxiety may be mistaken for inability
- Underestimation of STEM competency due to unfavorable comparisons with Cluster 1 peers

Cluster 3: Balanced Medium Performers (26 students, 26%)
*Profile: Average performance across all subjects, no pronounced specialization*

Immediate Recommendations:

- Complete interest and aptitude assessments to identify hidden strengths
- Sample diverse elective courses to discover potential specialization areas
- Set specific, measurable goals for grade improvement in one subject area
- Develop consistent study skills and time management strategies

Career Pathway Guidance:

- Keep multiple pathways open—specialization decisions are not yet urgent
- Explore career clusters through informational interviews and job shadowing
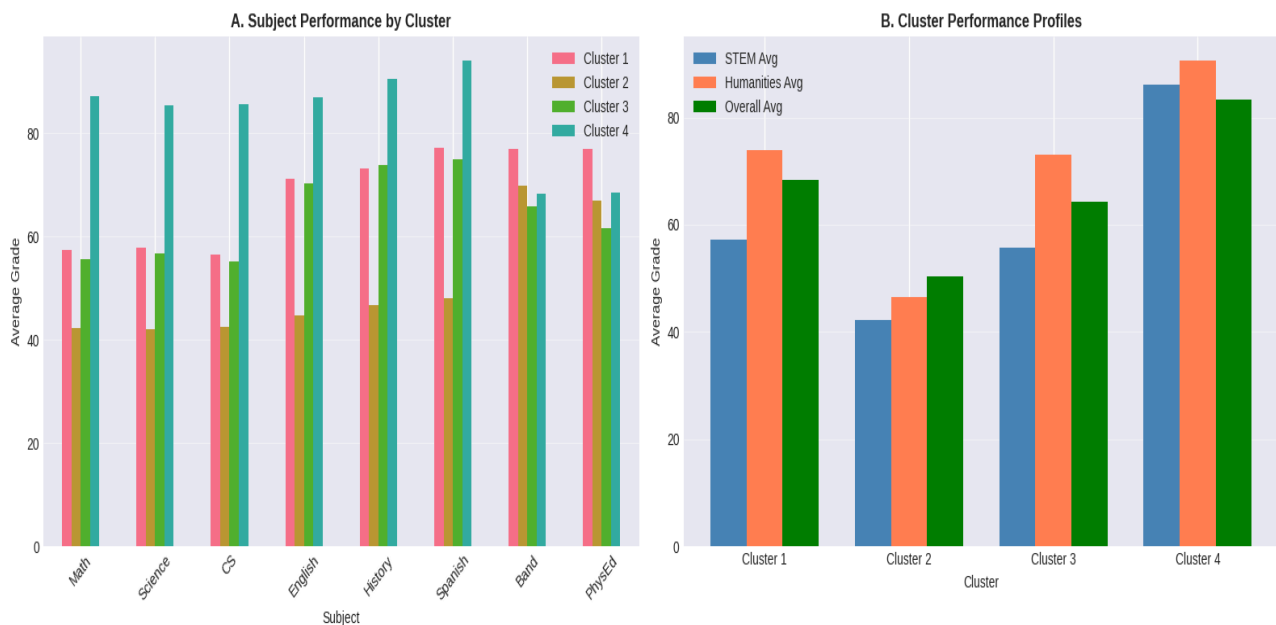- Connect academic subjects to real-world applications to increase engagement

Cluster 4: Low Performers (22 students, 22%)
*Profile: Below-average performance, particularly weak in Mathematics*

Immediate Recommendations:

- Conduct comprehensive diagnostic assessment to identify specific learning barriers
- Establish daily tutoring in Mathematics with consistent, supportive provider
- Reduce course load if necessary to focus on foundational skills
- Ensure continued enrollment in Band/Physical Education where success is demonstrated



PCA Cluster Profile Analysis

Q15. Identify three specific, actionable recommendations for the school administration derived from this analysis.

Recommendation 1: Implement Differentiated Mathematics Pathways

Problem Statement:
Mathematics exhibits the lowest average performance (51.6), left-skewed distribution, and serves as a gateway subject for STEM success. The one-size-fits-all Mathematics curriculum is failing 60% of students while potentially under-challenging the 28% of students in Cluster 1 who demonstrate advanced quantitative ability.

Proposed Solution:
Restructure Mathematics course offerings into three parallel pathways rather than a single sequential track:

Pathway A: Accelerated/Enrichment (25% of students)

- Faster pacing through standard curriculum
- Additional enrichment topics (cryptography, data science, mathematical modeling)
- Preparation for AP Calculus by Grade 11
- Target population: Cluster 1 students, STEM-focused high performers

Pathway B: Standard/Applied (50% of students)

- Traditional pacing aligned with grade-level standards
- Increased emphasis on applied problems and real-world connections
- Integration with Science and Computer Science curricula
- Target population: Cluster 3 students, balanced medium performers, Humanities-focused high performers needing Mathematics maintenance

Pathway C: Foundational/Supported (25% of students)

- Reduced pacing with additional instructional time
- Diagnostic-driven remediation of prerequisite skill gaps
- Co-requisite support model (enroll in grade-level content WITH concurrent foundational skills course)
- Emphasis on quantitative literacy rather than calculus preparation
- Target population: Cluster 4 students, transitional at-risk group, students with persistent Mathematics difficulty

Implementation Steps:

1. Administer universal diagnostic assessment to all students currently enrolled in Mathematics
2. Assign students to pathways based on demonstrated skill levels, not prior grades alone
3. Provide professional development for Mathematics faculty on differentiated instruction
4. Schedule pathway courses concurrently to allow student movement between tiers
5. Evaluate effectiveness through pre-post assessment and student progression data

Expected Outcomes:

- Reduction in Mathematics failure rates within one academic year
- Improved engagement among both struggling and advanced students
- More appropriate challenge levels for all students
- Increased enrollment in advanced STEM coursework

Resource Requirements:

- Additional Mathematics sections (splitting one large course into three smaller pathways)
- Instructional materials for applied/enrichment content
- Diagnostic assessment system
- Professional development investment

Recommendation 2: Develop Early Identification System for At-Risk Students

Problem Statement:
t-SNE analysis revealed a transitional "at-risk" group of 12 students positioned between low and medium performers. These students currently have overall averages of 55-60 and are not yet identified as failing, but exhibit clear early warning signs—particularly declining Mathematics performance and emerging STEM-Humanities gaps. Current identification systems activate only after students have already failed, representing lost opportunities for preventive intervention.

Proposed Solution:
Implement a data-driven early warning system that identifies at-risk students based on patterns revealed in PCA and t-SNE analysis:

Early Warning Indicators:

1. Mathematics grade < 65 AND declining trend over two consecutive grading periods
2. Mathematics-Humanities gap > 15 points (in either direction)
3. STEM average < 60 despite passing Humanities grades

4. Overall average decline of > 5 points from previous academic year
5. Unexcused absences in Mathematics or Science courses

Tiered Intervention Protocol:

Tier 1 (Universal Screening - All Students):

- Quarterly review of all students on early warning indicators
- Automated dashboard for counselors and administrators

Tier 2 (Targeted Support - At-Risk Identified):

- Counselor check-in within 5 school days of identification
- Diagnostic assessment in identified weak subject area
- Teacher-team meeting to develop classroom-level accommodations
- Parent notification and engagement

Tier 3 (Intensive Support - Persistent Risk):

- Daily tutoring in identified subject area
- Reduced course load if appropriate
- Student success team meeting with all stakeholders
- Individualized learning plan development

Implementation Steps:

1. Configure student information system to generate early warning reports
2. Train counselors and administrators on early warning indicators and protocols
3. Establish weekly student success team meetings to review identified students
4. Document interventions and track student outcomes
5. Refine indicators based on predictive validity analysis

Expected Outcomes:

- Reduction in students progressing from "at-risk" to "low performer" status
- Earlier identification of learning barriers
- More efficient allocation of intervention resources
- Improved student outcomes with less intensive intervention required

Resource Requirements:

- Data system configuration (minimal cost, primarily staff time)
- Counselor and administrator training
- Tutoring resource allocation

- Early intervention is significantly less expensive than remediation after failure

Recommendation 3: Preserve and Protect Non-Academic Success Pathways

Problem Statement:
Band demonstrates the highest average performance (72.1) and lowest variability (std=10.6) of any subject, with near-zero correlation to academic performance. Physical Education shows a bimodal distribution but provides above-average performance for students in Cluster 4 and the transitional at-risk group. Despite this demonstrated success, non-academic subjects are frequently deprioritized or eliminated for struggling students under the mistaken belief that removing "non-essential" courses allows more time for academic remediation.

Proposed Solution:
Formalize Band and Physical Education as protected success pathways, particularly for academically struggling students, and investigate pedagogical practices from these subjects for potential cross-curricular application.

Policy Recommendations:

1. Protect Enrollment:
Establish policy that students identified as academically struggling will NOT be removed from Band or Physical Education for additional remediation time. Exceptions require student success team approval and documented rationale.

2. Investigate Transferable Practices:
Conduct structured observation and interview study of Band and successful Physical Education sections to identify pedagogical practices that may transfer to academic subjects:

- Band-specific practices to investigate:
  - Mastery-based progression (students advance when ready, not on fixed schedule)
  - Ensemble/collaborative learning model
  - Public performance as authentic assessment
  - Deliberate practice with immediate feedback
  - Clear skill progression with visible milestones
- Physical Education-specific practices to investigate:
  - Multiple assessment modalities (skill demonstration, fitness, knowledge, participation)
  - Self-competition rather than peer competition (personal improvement focus)
  - Choice within structured parameters

- ○ Low-stakes skill development before high-stakes assessment

3. Expand Success Definition:
Revise school accountability metrics to value student success in non-academic domains, not solely academic proficiency. Recognize Band and Physical Education achievement in honor roll criteria, award ceremonies, and college recommendation letters.

4. Pilot Cross-Curricular Integration:
Develop pilot units where Mathematics or Science concepts are taught through music or physical activity contexts. Evaluate engagement and learning outcomes compared to traditional instruction.

Implementation Steps:

1. Conduct administrative review of current policies regarding course enrollment for struggling students
2. Document cases where students were removed from arts/PE for remediation and track subsequent outcomes
3. Develop formal policy protecting arts/PE enrollment with clear exception criteria
4. Initiate pedagogical study with Band and PE department chairs
5. Report findings to faculty and identify potential instructional applications

Expected Outcomes:

- Preservation of student motivation and engagement through continued success experiences
- Reduced learned helplessness among academically struggling students
- Potential identification of instructional practices that could improve academic outcomes
- More balanced, holistic view of student achievement

Resource Requirements:

- Administrative time for policy review and development
- Released time for faculty collaboration and observation
- No additional budgetary resources required

Conclusion:
These three recommendations address the most significant findings from the analysis: the systemic Mathematics challenge, the missed opportunity for early intervention, and the underutilized value of non-academic success pathways. Implementation requires modest resources relative to the potential impact on student outcomes.