

Evaluating NLP tools designed to assist instructors with formative assessment for large-enrollment STEM classes

Matthew Beckman
Penn State University

January 6, 2026

Slides: mdbeckman.github.io/JMM2026-WashingtonDC/

Two question pop quiz

- 1 Is your lucky (or favorite) number odd or even?
- 2 How would you describe the value of formative assessment?

Google Form



Figure 1: (QR Code) <https://forms.gle/hpW72fMYE1SsB19JA>

Same prompt, different task!

Value of Formative Assessment

mdb268@psu.edu [Switch account](#)

Not shared

Odd

How would you describe the value of formative assessment?

Your answer

[Back](#) [Submit](#) [Clear form](#)

Never submit passwords through Google Forms.

This form was created inside of psu.edu. [Report Abuse](#)

Google Forms

Value of Formative Assessment

mdb268@psu.edu [Switch account](#)

Not shared

Even

How would you describe the value of formative assessment?

- ☐ empowers students to monitor their own learning outcomes
- ☐ enables instructors to monitor learning outcomes of students on an individual or aggregate basis
- ☐ provides feedback that instructors can use to address misconceptions and/or adjust instruction
- ☐ amenable to low-stakes and high-frequency opportunities for students to engage with content
- ☐ I'm not sure

[Back](#) [Submit](#) [Clear form](#)

Motivation

- “Write-to-learn” tasks improve learning outcomes (Graham, et al., 2020)
- Critical for citizen-statisticians to communicate effectively (Gould, 2010)
- Frequent practice w/ communicating improves statistical literacy and promotes retention (Basu, et al., 2013)
- Formative assessment benefits both students & instructors (Black & William, 2009; GAISE, 2016; Pearl, et al., 2012)
- A majority of U.S. undergraduates at public institutions take at least one large-enrollment STEM course (Supiano, 2022)
- **Logistics of constructed response tasks jeopardize use in large-enrollment classes** (Garfield & Ben-Zvi, 2008; Woodard & McGowan, 2012)

Easy!



Erm...



Goal

Develop technology that can assist instructors for large (STEM) classes with providing targeted formative assessment feedback to students, such that instructor burden is similar to small class (~30 students)

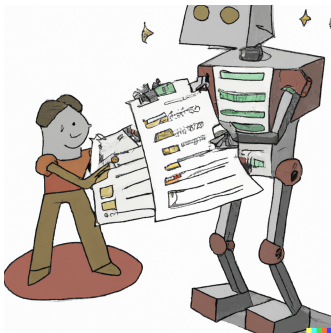
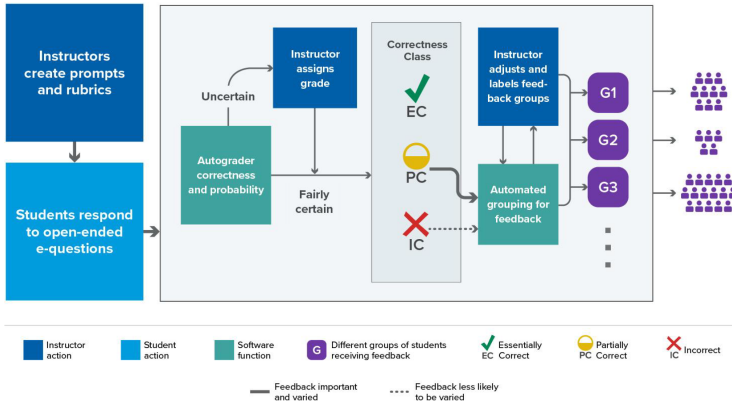


Figure 2: image created with assistance of DALL · E 2 by Open AI

Schematic



Goal: Computer-assisted formative assessment feedback for short-answer tasks in large-enrollment classes, such that instructor burden is similar to small class (~30 students)

Research Questions

- **RQ1:** What level of agreement is achieved among trained human raters labeling (i.e., scoring) short-answer tasks?
- **RQ2:** What level of agreement is achieved between human raters and an NLP algorithm?

Relevant Papers

- Lloyd, S. E., Beckman, M., Pearl, D., Passonneau, R., Li, Z., & Wang, Z. (2022). Foundations for AI-Assisted Formative Assessment Feedback for Short-Answer Tasks in Large-Enrollment Classes. In *Proceedings of the eleventh international conference on teaching statistics*. Rosario, Argentina.
- Li, Z., Lloyd, S., Beckman, M. D., & Passonneau, R. J. (2023). Answer-state Recurrent Relational Network (AsRRN) for Constructed Response Assessment and Feedback Grouping. *Findings of the Association for Computational Linguistics: EMNLP 2023*.
- Beckman, M., Burke, S., Fiochetta, J., Fry, B., Lloyd, S. E., Patterson, L., & Tang, E. (in review). Developing Consistency Among Undergraduate Graders Scoring Open-Ended Statistics Tasks. Preprint URL: <https://arxiv.org/abs/2410.18062>

Collaborators

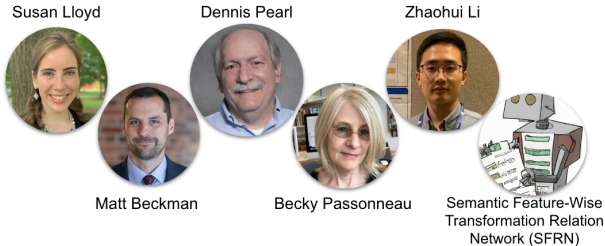


Figure 3: Lloyd et al., (2022); Li et al., (2023) Project Team

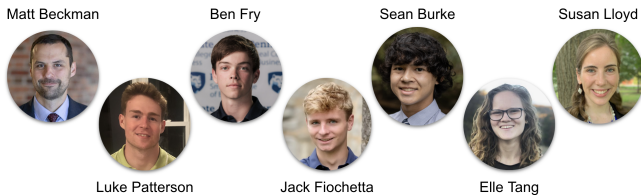



Figure 4: Beckman et al., (in review) Project Team

Methods (Short-answer task)

4. Walleye is a popular type of freshwater fish native to Canada and the Northern United States. Walleye fishing takes much more than luck; better fishermen consistently catch larger fish using knowledge about proper bait, water currents, geographic features, feeding patterns of the fish, and more. Mark and his brother Dan went on a two-week fishing trip together to determine who the better Walleye fisherman is. Each brother had his own boat and similar equipment so they could each fish in different locations and move freely throughout the area. They recorded the length of each fish that was caught during the trip, in order to find out which one of them catches larger Walleye on average.

a. Should statistical inference be used to determine whether Mark or Dan is a better Walleye fisherman? Explain why statistical inference should or should not be used in this scenario.



b. Next, explain how you would determine whether Mark or Dan is a better Walleye fisherman using the data from the fishing trip. *(Be sure to give enough detail that a classmate could easily understand your approach, and how he or she would interpret the result in the context of the problem.)*

Figure 5: Sample task including a stem and two short-answer prompts.

Methods (RQ1)

RQ1: What level of agreement is achieved among trained human raters labeling (i.e., scoring) short-answer tasks?

- Lloyd et al., (2022)
 - 3 raters typical of large-enrollment instruction team
 - (6 tasks) × (1,935 students) distributed among the team
 - sufficient intersection to assess inter-rater agreement
 - responses judged Correct / Partial / Incorrect against rubric
- Beckman et al., (in review)
 - 4 Undergraduate Teaching Assistants (UTAs) and 1 instructor
 - (4 tasks) × (63 students) scored by each UTA + Instructor
 - 5 sequential exercises associated with progression of scoring development

Results

- “short-answer” tasks are good for students, but hard to scale
- Can NLP tools help instructors with scale?
 - Evaluate & group student responses
 - Compare agreement between NLP & humans

Scoreboard¹

- (RQ1) Instructor agreement (QWK \approx 0.7 to 0.8+)
- (RQ1) UTA agreement (QWK \approx 0.6 to 0.7+)
- **What about... NLP algorithm & instructor agreement?**

¹Lloyd, et al. (2022); Beckman, et al. (in review)

Methods (RQ2)

RQ2: What level of agreement is achieved between human raters and an NLP algorithm?

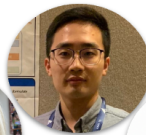
Susan Lloyd



Dennis Pearl



Zhaohui Li



Matt Beckman



Becky Passonneau



Semantic Feature-Wise
Transformation Relation
Network (SFRN)

Paper introducing SFRN

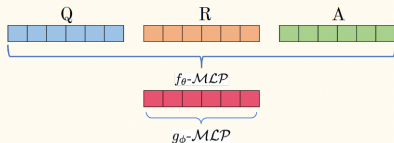
Li, Z., Tomar, Y., & Passonneau, R. J. (2021). A Semantic Feature-Wise Transformation Relation Network for Automatic Short Answer Grading. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6030–6040. Association for Computational Linguistics.
<https://aclanthology.org/2021.emnlp-main.487>

Meet the “machine”: NLP for Assessment

- Natural language processing (NLP) involves how computers can be programmed to analyze language elements
- NLP-assisted feedback for educational use:
 - automated short-answer grading (ASAG) from 2009
 - essays & long-answer tasks earlier
- Human-machine collaboration is a promising mechanism to assist rapid, individualized feedback at scale (Basu, 2013)
- Deep neural networks application since 2016
- Relational (neural) networks

Meet the “machine”: Relational Networks

Motivation for a Relation Network



- Many short-answer datasets have triples
 - Question prompt
 - Rubric OR Reference answers
 - Answer from student
- Transformers are less practical
 - Datasets are often relatively small
 - Learning a single vector can efficiently capture relational structure

Q: Susan has samples of 5 different foods. Using only the results of her experiment, how will Susan know which food contains the most sugar? (Gas volume is evaluated by tube)

R: Susan should compare the amount of gas in each bag. The bag with the most gas contains the food with the most sugar.

A: Susan will know how much sugar is in the foods by putting each bag in a volume tube. When her finder stops after pushing the top, the bottom of the part she pushes down will be on a number. That number is the milliliters of sugar in the food. Whichever number is the highest, that means that food has the most sugar.

Figure 6: Image credit: Becky Passonneau

- relation networks designed to learn generalizations that infer meaning in a data-efficient way
- much of the architecture inspired by work from computer vision

Meet the “machine”: SFRN Schematic

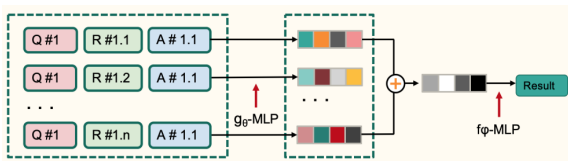


Figure 7: encoder (Left); fusion function (Middle); classifier (Right).

Semantic Feature-Wise Transformation Relation Network (SFRN):

- end-to-end model with three components:
 - (g_{θ} MLP) pretrained BERT encoder (LLM) » vector representations
 - (+) learned feature-wise transformation function fuses multiple representations, as needed (e.g., if multiple reference answers)
 - (f_{ϕ} MLP) is a classifier algorithm, i.e., neural network
- data augmentation during training step

Results

- “short-answer” tasks are good for students, but hard to scale
- Can NLP tools help instructors with scale?
 - Evaluate & group student responses
 - Compare agreement between NLP & humans

Scoreboard²

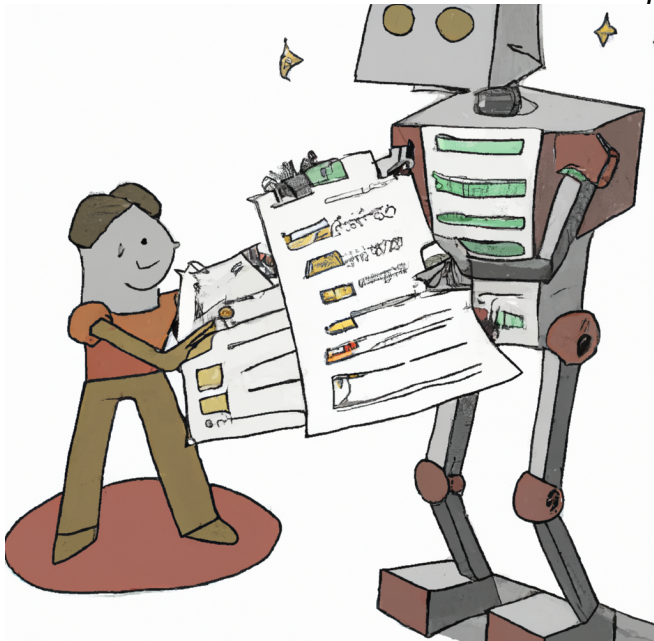
- (RQ1) Instructor agreement (QWK \approx 0.7 to 0.8+)
- (RQ1) UTA agreement (QWK \approx 0.6 to 0.7+)
- (RQ2) NLP algorithm & instructor agreement (QWK \approx 0.7+)
- **What if we combine the Human & Machine??**

²Lloyd, et al. (2022); Beckman, et al. (in review)



Figure 8: Image credit: <https://www.slugmag.com/arts/film/film-reviews/terminator-genisys-time-is-not-on-my-side/>

Human-Machine Partnership



Human-Machine *Partnership*

Our approach to human-in-the-loop (HIL) did **not** make a recommendation (e.g., Left), it just shows examples to the human when it needs help (e.g., Right).

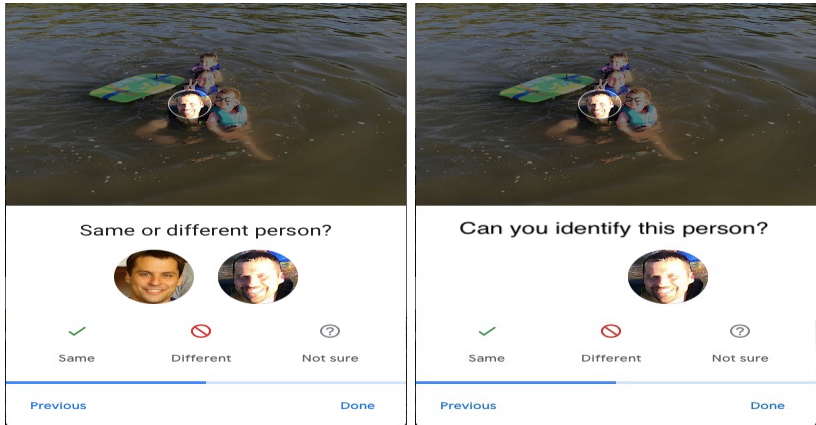


Figure 10: Illustration adapted from Google Photos

Human-Machine Partnership Method

Want to evaluate accuracy of marking algorithm when designed to “defer” to human judgment

- algorithm evaluates a probability for each label (EC, PC, IC)
 - if a label has high probability, use algorithm label
 - if no label has sufficiently high probability, defer to human
- interests
 - estimate how frequently the algorithm defers
 - estimate accuracy of the combined process

Human-Machine Partnership Results

Our work is first that we know of to implement controllable, selective prediction deferral policy for the classifier (i.e., scoring) step.

Threshold	Deferral Rate	Simulated HIL Accuracy
0.68	9.5%	0.855
0.75	13.2%	0.861
0.80	16.0%	0.871
0.85	20.2%	0.884
0.90	25.6%	0.899

Figure 11: Accuracy of Human-in-the-loop compared with expert label ground truth.

Results

- “short-answer” tasks are good for students, but hard to scale
- Can NLP tools help instructors with scale?
 - Evaluate & group student responses
 - Compare agreement between NLP & humans

Scoreboard³

- (RQ1) Instructor agreement (QWK ≈ 0.7 to $0.8+$)
- (RQ1) UTA agreement (QWK ≈ 0.6 to $0.7+$)
- (RQ2) NLP algorithm agreement with instructors (QWK $\approx 0.7+$)
- (RQ2) Human-Algorithm partnership may be even better? ($\approx 0.85+$)

³Li, et al., (2023)

Discussion

- **RQ1:** Substantial agreement achieved among trained human raters provides context for further comparisons
- **RQ2:** NLP algorithm produced agreement reasonably aligned to results achieved by pairs/groups of trained human raters
 - Human-in-the-Loop » Instructor / Algorithm partnership
- **What about feedback? A few avenues come to mind...**
 - Should AI just do it?
 - Clustering (or Classifier) Tools?
 - Topological Data Analysis?
 - Something completely different?

Discussion

- **RQ1:** Substantial agreement achieved among trained human raters provides context for further comparisons
- **RQ2:** NLP algorithm produced agreement reasonably aligned to results achieved by pairs/groups of trained human raters
 - Human-in-the-Loop » Instructor / Algorithm partnership
- **What about feedback? A few avenues come to mind...**
 - ~~Should AI just do it?~~
 - Clustering (or Classifier) Tools?
 - Topological Data Analysis?
 - Something completely different?

Current Events & Next Steps

- challenge system with diverse tasks, institutions, student populations;
 - partnering with ISU, MSU, PSU, UCSB, UF, UTEP, & UoA
 - approx 45,000 student responses; 27 task prompts
 - *influence of linguistic diversity?*
- accumulated data to be shared with broader NLP community
 - this will be among the largest **open** data sources of it's kind
 - addresses barriers imposed by proprietary data sources on NLP research
- algorithm development
 - contrastive loss function
 - flexibility for task structure
 - studying influence of rubric features
 - LLM performance benchmarking (Wei et al, in review)

Acknowledgments

- US National Science Foundation (NSF DUE-2236150: Project CLASSIFIES)
- Penn State Center for Socially Responsible Artificial Intelligence
- Strategic partnership between University of Auckland and Penn State University
- Thanks to students and faculty at partner institutions that have assisted us with data collection.

References (1/3)

- 1 Basu, S., Jacobs, C., & Vanderwende, L. (2013). Powergrading: a Clustering Approach to Amplify Human Effort for Short Answer Grading. *Transactions of the Association for Computational Linguistics*, 1, 391–402. https://doi.org/10.1162/tacl_a_00236
- 2 Beckman, M. (2015). Assessment Of Cognitive Transfer Outcomes For Students Of Introductory Statistics. <http://conservancy.umn.edu/handle/11299/175709>
- 3 Beckman, M., Burke, S., Fiochetta, J., Fry, B., Lloyd, S. E., Patterson, L., & Tang, E. (in review). Developing Consistency Among Undergraduate Graders Scoring Open-Ended Statistics Tasks. Preprint URL: <https://arxiv.org/abs/2410.18062>
- 4 Black, P., & Wiliam, D. (2009). Developing the theory of formative assessment. *Educational Assessment, Evaluation and Accountability*, 21, pp 5-31. <https://doi.org/10.1007/s11092-008-9068-5>
- 5 GAISE College Report ASA Revision Committee (2016). Guidelines for Assessment and Instruction in Statistics Education College Report 2016. URL: <http://www.amstat.org/education/gaise>
- 6 Gould, R. (2010). Statistics and the Modern Student. *International Statistical Review / Revue Internationale de Statistique*, 78(2), 297–315. <https://www.jstor.org/stable/27919839>

References (2/3)

- 7 Guo, W., Diab, M. (2012) Modeling Sentences in the Latent Space. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*, pages 864–872. Association for Computational Linguistics.
- 8 Graham, S., Kiuahara, S. A., & MacKay, M. (2020). The Effects of Writing on Learning in Science, Social Studies, and Mathematics: A Meta-Analysis. *Review of Educational Research*, 90(2), 179–226.
- 9 Gwet, K. (2014). Handbook of Inter-Rater Reliability: The Definitive Guide to Measuring the Extent of Agreement Among Raters
- 10 Li, Z., Tomar, Y., & Passonneau, R. J. (2021). A Semantic Feature-Wise Transformation Relation Network for Automatic Short Answer Grading. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 6030–6040. Association for Computational Linguistics.
<https://aclanthology.org/2021.emnlp-main.487>
- 11 Li, Z., Lloyd, S., Beckman, M. D., & Passonneau, R. J. (2023). Answer-state Recurrent Relational Network (AsRRN) for Constructed Response Assessment and Feedback Grouping. *Findings of the Association for Computational Linguistics: EMNLP 2023*. <https://doi.org/10.18653/v1/2023.findings-emnlp.254>
- 12 Lloyd, S. E., Beckman, M., Pearl, D., Passonneau, R., Li, Z., & Wang, Z. (2022). Foundations for AI-Assisted Formative Assessment Feedback for Short-Answer Tasks in Large-Enrollment Classes. In *Proceedings of the eleventh international conference on teaching statistics*. Rosario, Argentina.

References (3/3)

- 13 Pearl, D. K., Garfield, J. B., delMas, R., Groth, R. E., Kaplan, J. J., McGowan, H., & Lee, H. S. (2012). Connecting Research to Practice in a Culture of Assessment for Introductory College-level Statistics. URL: http://www.causeweb.org/research/guidelines/ResearchReport_Dec_2012.pdf
- 14 U.S. Department of Education, Office of Educational Technology (2023). Artificial Intelligence and Future of Teaching and Learning: Insights and Recommendations, Washington, DC.
- 15 Viera, A. J., & Garrett, J. M. (2005). Understanding interobserver agreement: the kappa statistic. *Family Medicine*, 37(5), 360–363.
- 16 Woodard, R., & McGowan, H. (2012). Redesigning a large introductory course to incorporate the GAISE guidelines. *Journal of Statistics Education*, 20(3).
- 17 Wei, Y., Pearl, D., Beckman, M., Passonneau, R. (2025). Concept-based Rubrics Improve LLM Formative Assessment and Data Synthesis. Preprint URL: <https://arxiv.org/pdf/2504.03877>

Thank You

Evaluating NLP tools designed to assist
instructors with formative assessment for
large-enrollment STEM classes

Matthew Beckman
Penn State University

January 6, 2026

Slides: mdbeckman.github.io/JMM2026-WashingtonDC/

Supporting Slides

Results: Instructors as Graders

RQ1: What level of agreement is achieved among trained human raters labeling (i.e., scoring) short-answer tasks?

Comparison	Reliability
Rater A & Rater C	QWK = 0.83
Rater A & Rater D	QWK = 0.80
Rater C & Rater D	QWK = 0.79
Rater A 2015 & 2021	QWK = 0.88

Figure 12: Interrater agreement among three instructors; intra-rater agreement for Rater A with several years delay

Reliability intuition⁴: moderate < 0.6 < substantial < 0.8 < near perfect < 1.0

⁴Viera & Garret (2005)

Results: Instructor and UTA Graders

RQ1: What level of agreement is achieved among trained human raters labeling (i.e., scoring) short-answer tasks?

Raters	Day 1	Day 5	Week 10
A & E	0.46 (0.35, 0.58)	0.57 (0.47, 0.67)	0.58 (0.49, 0.67)
A & F	0.61 (0.50, 0.71)	0.72 (0.64, 0.79)	0.78 (0.71, 0.85)
A & G	0.63 (0.55, 0.72)	0.73 (0.66, 0.80)	0.73 (0.66, 0.81)
A & H	0.72 (0.65, 0.80)	0.71 (0.63, 0.78)	0.68 (0.59, 0.78)

Figure 13: Pairwise agreement between UTAs and an instructor (Rater A)

Reliability intuition: moderate < 0.6 < substantial < 0.8 < near perfect < 1.0

Results: Instructor & UTA (cont'd)

Raters	QWK	95% CI
A	0.82	(0.76, 0.88)
E	0.57	(0.46, 0.68)
F	0.74	(0.67, 0.82)
G	0.66	(0.56, 0.76)
H	0.74	(0.67, 0.81)

Figure 14: Intra-rater agreement (self-consistency) for each participant as measured with Quadratic Weighted Kappa (QWK) while scoring the same set of student responses on two occasions approximately 10 weeks apart.

Date (Exercise)	Rubric Description	AC_2	95% CI
Day 1 (Ex 1)	Solution with Verbal Instructions	0.688	(0.63, 0.74)
Day 5 (Ex 4)	Expert Rubric, Part 1	0.784	(0.75, 0.82)
Week 10 (Ex 5)	Expert Rubric, Part 2	0.778	(0.74, 0.81)

Figure 15: Group agreement among four undergraduate TAs and one instructor, as measured with Gwet's (2014) AC_2 ; 95% confidence intervals accompany each estimate.

Results (RQ2)

RQ2: What level of agreement is achieved between instructors and the machine (an NLP algorithm)?

Comparison	Reliability
Rater A & SFRN	QWK = 0.79
Rater C & SFRN	QWK = 0.82
Rater D & SFRN	QWK = 0.74

Figure 16: Pairwise agreement with SFRN algorithm

Reliability intuition: moderate < 0.6 $<$ substantial < 0.8 $<$ near perfect < 1.0

Methods (RQ3): Humans

How similar is feedback provided by two instructors for some group of students?

- Two instructors independently evaluated 100 “partial credit” responses
- Each instructor provided free-text feedback to each student
- Verbatim feedback captured for each instructor and cross-tabulated for analysis.
- *Results:*
 - The two instructors gave substantially equivalent feedback to 66 of 100 responses
 - Evidence of two large “clusters” (and quite a few singletons)

Methods (RQ3): Machines

- Experiment #1
 - retrain k-means & k-medoids clustering & evaluate stability
 - compare representations with higher & lower dimensionality
 - **Results:**
 - SFRN ($D = 512$): cluster stability 0.62
 - Highest stability among competing algorithms was 0.88, achieved using a matrix factorization method that produces static representations ($D = 50$; WTMF; Guo & Diab, 2011)
 - *cursed*
- Experiment #2:
 - ~~clustering~~ \Rightarrow FB Classifier?
 - Both Humans & Machines attempt
 - **Results:**
 - NLP Algorithm was more consistent with instructor A on one task and instructor B on the other task tested.
 - *meh*

Results (RQ3 humans)

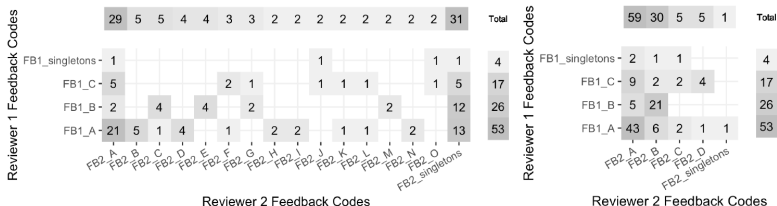


Figure 17: Cross-tabulation of feedback distribution for the two reviewers for the initial feedback (left) compared with the same analysis for the portion of feedback related to the statistical concept at issue (right).

- Reviewer 1 favored feedback on statistical concepts (only).
- Reviewer 2 provided same, plus a quote from the student
- Reviewer 2 parsed feedback to compare remarks related to the statistical concepts (only) with that of Reviewer 1.

Results (RQ3 humans)

Feedback Code	Feedback verbatim text suggested by the Reviewer
FB1_A (Reviewer 1)	What can we do to evaluate whether [the] result is better than we would expect for someone that is strictly guessing?
FB2_A (Reviewer 2)	Think about what inferential statistical method might we use to evaluate the percentage of correctly identified notes.
FB1_B (Reviewer 1)	Good idea to have a threshold for comparison, but it's very important that it be established carefully. For example, how might you establish a threshold that...
FB2_B (Reviewer 2)	Why this threshold? What inferential statistical method might we use to evaluate the percentage of correctly identified notes?

Figure 18: Verbatim feedback most frequently provided by each reviewer for responses to task 2B.

Results (RQ3 machines)

RQ3: What sort of NLP representation leads to good clustering performance, and how does that interact with the classification algorithm?

- SFRN ($D = 512$) produced reasonably consistent clusters when retrained (0.62)
- Highest consistency (0.88; $D = 50$) was achieved using a matrix factorization method that produces static representations (WTMF; Guo & Diab, 2011)
- AsRRN compared to humans (A & B) grouping students by pre-determined feedback categories:

Task	Sample Size	A & B	A & AsRRN	B & AsRRN
1	90	0.71	0.53	0.69
2	100	0.45	0.70	0.41