

LSTM Neural Network Assisted Regex Development for Qualitative Coding*

Zhiqiang Cai^[0000-0002-2107-3378], Brendan Eagan^[0000-0002-9189-1747], Cody Marquart^[0000-0002-3387-6792], and David W. Shaffer^[0000-0001-9613-5740]

University of Wisconsin-Madison

{zhiqiang.cai, beagan, cody.marquart}@wisc.edu, dws@education.wisc.edu

Abstract. Regular expression (regex) based automated qualitative coding helps reduce researchers’ effort in manually coding text data, without sacrificing transparency of the coding process. However, researchers using regex based approaches struggle with low recall or high false negative rate during classifier development. Advanced natural language processing techniques, such as topic modeling, latent semantic analysis and neural network classification models help solve this problem in various ways. The latest advance in this direction is the discovery of the so called “negative reversion set (NRS)”, in which false negative items appear more frequently than in the negative set. This helps regex classifier developers more quickly identify missing items and thus improve classification recall. This paper simulates the use of NRS in real coding scenarios and compares the required manual coding items between NRS sampling and random sampling in the process of classifier refinement. The result using one data set with 50,818 items and six associated qualitative codes shows that, on average, using NRS sampling, the required manual coding size could be reduced by 50% to 63%, comparing with random sampling.

Keywords: Negative reversion · Qualitative coding · LSTM neural network.

1 Introduction

Coding in qualitative research is traditionally a complex manual process of theory discovery [10]. Qualitative researchers go through the data item by item to extract information they find interesting or pertinent to their study. They discover patterns or themes to form big-C *Codes* [12]. Once a code is well defined, they go back to the data and systematically identify where each code occurs through out the entire dataset. As Shaffer and Ruis [13] pointed out, coding is basically a process of “defining concepts and identifying where they occur” [13]. While a researcher may be able to find the patterns for code definition using a relatively small part of a data set, “identifying where they occur” requires the researcher to go through every item in the data set. Such manual coding is often

* Supported by Natural Science Foundation

very expensive and can be impossible for quantitative ethnographers, because they often deal with data that is too large to code by hand.

In contrast to manual coding, machine learning algorithms have been used to train classifiers for automatic coding. Such algorithms use part of the data that researchers manually coded to train a classification model, which is then used to classify the rest of the data [9, 1]. Researchers have identified challenges in using machine learning for automatic coding [4, 2]. The so called “black box” issue is a major challenge for ethnographers, qualitative researchers, or people working in the digital humanities interested in adopting automated qualitative coding methods. Namely, when an excerpt is coded as *positive* for a certain code, the *evidence* for the coding is often unclear from the output of the machine learning model, which makes it difficult for researchers to be successful in the *interpretation* stage in the analyses. Another challenge is that the amount of human coded data required for training machine learning algorithms is large, especially when the base rate (positive rate or frequency) of a code is small. For example, if a code occurs in about 1% of the items in the data set, to train a reliable neural network classifier often requires about 4,000 manually coded training items. In addition, machine learning algorithms tend to build classifiers based on high frequency expressions (words, phrases). Thus, the expressions used by minority groups or sub-populations could be under represented, resulting in biased and unfair coding.

Regex based coding has been proven to be a more efficient method for qualitative coding, which often requires manually coding only a few hundreds of excerpts in order to develop a reliable regex-based classifier. Once regexes are well developed, the coding is an automatic process and the matched text elements provide explicit coding evidence [13, 12]. However, regex based classifiers suffer from the *low recall* problem [3, 2]. That is, regex based classifiers can suffer from too many false negatives, where a researcher or domain expert would identify that a code is present, but the classifier does not. Various ways have been proposed to overcome the challenge of the low recall problem. For example, Latent Semantic Analysis [11] can be used to form a “snowball” method to help recursively expand keyword set for regex construction [3]. The latest discovery is the so called “*negative reversion set*” (NRS) in which false negative items are denser than in the negative set of regex classifiers [2]. This discovery makes it possible to greatly shorten the time in searching for missing patterns in regex development because researchers can use this set to focus on coding cases that are more likely to be false negatives and use the patterns shown in those examples to improve their classifiers.

Researchers developing new methodologies often conduct simulation studies to examine or compare the performance of different statistics or analytical approaches [7, 6]. This paper simulates the use of NRS in real world scenario and tackles the issue of required manual coding size. The main question we want to answer is, how much manual coding efforts can be reduced by using the NRS?

2 Negative Reversion Set and Research Questions

2.1 Negative Reversion Set

The usefulness of Negative Reversion Set (NRS) was discovered by the *triangulation* of three “raters”: 1) a human; 2) a regex based classifier; and 3) a neural network model trained from the regex classifier. For a given data set D and a given code C , the human classification is denoted by $D = P + N$, where P is the set of human classified positive items and N the set of human classified negative items. The plus sign “+” denotes the union of the two sets. The classification by the regex classifier is denoted by $D = \tilde{P} + \tilde{N}$, where \tilde{P} and \tilde{N} are the regex classified positive and negative sets, respectively.

Taking human classification as ground truth, the four intersection sets of the two classifications by human and regex classifier are denoted by

- **true positive set** $P\tilde{P}$: both human and regex classify as positive;
- **false positive set** $N\tilde{P}$: human classifies as negative but regex classifies as positive;
- **false negative set** $P\tilde{N}$: human classifies as positive but regex classifies as negative; and
- **true negative set** $N\tilde{N}$: both human and regex classify as negative.

In the formulas above, “ XY ” denotes the intersection of two sets “ X ” and “ Y ”.

The error and accuracy of the regex classification are determined by the size of the above four sets. The following metrics are used in this paper to measure the error and accuracy of the regex classifier:

- proportion of true positives:

$$tp = \frac{|P\tilde{P}|}{|D|};$$

- proportion of false positives:

$$fp = \frac{|N\tilde{P}|}{|D|};$$

- proportion of false negatives:

$$fn = \frac{|P\tilde{N}|}{|D|};$$

- proportion of true negatives:

$$tn = \frac{|N\tilde{N}|}{|D|};$$

- precision:

$$precision = \frac{tp}{|\tilde{P}|} = \frac{tp}{tp + fp};$$

– recall:

$$recall = \frac{tp}{|P|} = \frac{tp}{tp + fn}; \text{ and}$$

– Cohen’s kappa:

$$\kappa = \frac{p_o - p_c}{1 - p_c}$$

where $p_o = tp + tn$ is the observed agreement and $p_c = (tp + fp)(tp + fn) + (tn + fp)(tn + fn)$ is the chance agreement.

A regex classifier provides input to a neural network model to train another classifier, which classifies the data as $D = \tilde{P} + \tilde{N}$. Replacing \tilde{P} and \tilde{N} by $\tilde{\tilde{P}}$ and $\tilde{\tilde{N}}$ in the above formulas, the metrics about the performance of the neural network model are computed the same way as the regex classifier described above.

Table 1 shows the intersections of the three classifications. The so called “Negative Reversion Set” is the set of items that regex classifies as negative but the neural network model classifies as positive, i.e., the intersection of the regex negative set \tilde{N} and the neural network positive set $\tilde{\tilde{P}}$, which can be written as the union of two sets — a “correct” negative reversion set $P\tilde{N}\tilde{\tilde{P}}$ and an “incorrect” negative reversion set $N\tilde{N}\tilde{\tilde{P}}$:

$$\begin{aligned} NRS &= \tilde{N}\tilde{\tilde{P}} \\ &= P\tilde{N}\tilde{\tilde{P}} + N\tilde{N}\tilde{\tilde{P}}. \end{aligned} \tag{1}$$

Table 1. Intersect sets of three classifications.

Human	Regex	Neural Network
P	$P\tilde{P}$	$P\tilde{\tilde{P}}$ $P\tilde{\tilde{N}}$
	$P\tilde{N}$	$P\tilde{N}\tilde{\tilde{P}}$ $P\tilde{N}\tilde{\tilde{N}}$
N	$N\tilde{P}$	$N\tilde{P}\tilde{\tilde{P}}$ $N\tilde{P}\tilde{\tilde{N}}$
	$N\tilde{N}$	$N\tilde{N}\tilde{\tilde{P}}$ $N\tilde{N}\tilde{\tilde{N}}$

In Equation 1, the “correct” negative reversion set $P\tilde{N}\tilde{\tilde{P}}$ is the set of items that are falsely classified by regex as negative but correctly “reversed” back to positive by neural network. The “incorrect” negative reversion set $N\tilde{N}\tilde{\tilde{P}}$ is the set of items that regex correctly coded as negative but wrongly reversed by the neural network as positive. The proportion of correctly reversed items in

the negative reversion set is used to measure the false negative density in the negative reversion set. For any two sets $X \subseteq Y$, we define the density of X in Y by

$$d(X, Y) = \frac{|X|}{|Y|}.$$

Thus the density of correctly reversed false negative items in the negative reversion set can be written as

$$d(P\tilde{N}\tilde{\tilde{P}}, \tilde{N}\tilde{\tilde{P}}) = \frac{|P\tilde{N}\tilde{\tilde{P}}|}{|\tilde{N}\tilde{\tilde{P}}|},$$

and the false negative density in the regex negative set is given by

$$d(P\tilde{N}, \tilde{N}) = \frac{|P\tilde{N}|}{|\tilde{N}|}.$$

Cai et al. [2] found that the density of false negative items in NRS could be much higher than in the regex negative set. That is,

$$d(P\tilde{N}\tilde{\tilde{P}}, \tilde{N}\tilde{\tilde{P}}) \gg d(P\tilde{N}, \tilde{N}).$$

As a side note, we point out that similar notations can be used to define “Positive Reversion Set (PRS)” which could be written as the union of “Correct Positive Reversion” and “Incorrect Positive Reversion”:

$$PRS = \tilde{P}\tilde{N} = P\tilde{P}\tilde{N} + N\tilde{P}\tilde{N}.$$

Positive reversion set could be useful for improving precision of regex classifiers, which is an interesting issue but out of the scope of this paper.

2.2 NRS Assisted Regex Development

Based on their discovery, Cai et al. [2] proposed a regex development procedure (see Figure 1). In their procedure, the regex development starts from an initial list. This regex list is a set of “atomic regex”, each of which cannot be decomposed as two regexes combined by the symbol “|”. For example, “\bchair\btable” is not an atomic regex because it can be decomposed as “\bchair” and “\btable”. A researcher provides one or more keywords that can represent a part of the code. The keywords are then used to construct the atomic regexes.

Once the initial regex list is obtained, the data is coded by the regex list and a neural network model is then trained using a part of the data classified by the regex list. The reversion sets, including negative reversion set and positive reversion set, are then identified by intersecting the regex classification and neural network classification. Items from reversion sets are then selected and presented to the researcher to rate the occurrence of the code. A conflict occurs if the researcher’s rating is different from the regex classification. The researcher may

resolve the conflict by changing the rating or modifying the regex list. When the regex list is modified, a new neural network model is re-generated and new reversion sets are identified. New items are continuously selected and presented to the researcher, until the researcher’s rating and the regex classification reach an acceptably high level of agreement. In this procedure, reversion sets serve as providers of most likely conflict items.

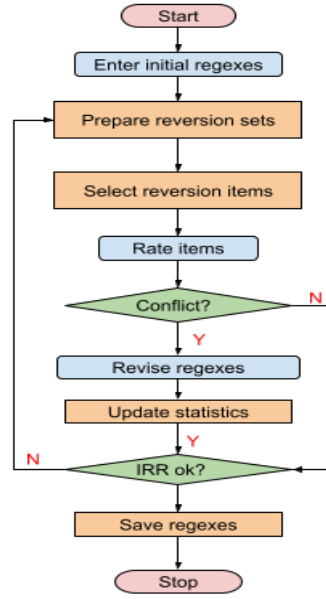


Fig. 1. Neural network assisted regex development for qualitative coding.

2.3 Research Questions

The regex development procedure suggested by Cai et al. [2] assumed that selecting items from a set with denser false negatives helps shorten the searching of new regexes. However, exactly how much such selection helps increase the efficiency of regex based classifier refinement remains an unanswered question. In this paper, we attempt to answer this question. Specifically, we ask:

RQ1 How does the false negative density change over multiple iterations? We know that the false negative density in the negative reversion set $\tilde{N}\tilde{P}$ could be much higher than in the regex negative set \tilde{N} . However, for each subsequent iterations, the number of false negative items reduces and thus the density could reduce quickly. This leads to the following questions, 1) How dense could it be at the beginning of the process? and 2) How quickly would it drop?

RQ2 Is using the negative reversion set more efficient than random sampling for the purpose of regex based classifier refinement, if so how much more? We hypothesize that sampling from a set with dense false negatives could help a researcher more quickly identify missing regexes. The question is, for a given level of required coding accuracy, how many fewer items would a researcher need to examine when using the negative reversion set compared to random sampling?

3 Method

3.1 Data

The data we used in this study were collected by previous researchers from an engineering virtual internship called *Nephrotex* [5, 8]. Participants worked as interns at a fictitious company that designs and manufactures ultrafiltration membranes for hemodialysis machinery used to treat end-stage renal failure. The work was divided into two phases. In the first phase, participants were grouped into teams of five. Each team worked on a specific task. In the second phase, participants were regrouped, or jigsawed, into new teams of five members from different teams in the first phase. The task in the second phase was to reflect the work each member did in the first phase to collaborate with their new team with varied expertise. The utterances in all online team chat conversations were collected, resulting in a data set with 50,888 utterances. After filtering out 70 longest utterances, 50,818 utterances were used in this study, each of which had a length not more than 100 words.

3.2 Codes

Previous researchers created six codes, including *Tech Constraints*, *Performance*, *Collaboration*, *Design*, *Data* and *Requests*. They developed and validated regexes for each code. Table 2 shows the name, description, example regex and IRR (Inter-Rater Reliability) of the six codes. The three numbers in the IRR column are the kappa values between three pairs of raters: human rater 1 : human rater 2, regex : human rater 1, and regex : human rater 2. Since the regexes were well developed and validated, in this study we consider the regex lists “complete” and use the full regex classification as human classification.

3.3 Three Classifiers

This study involves three types of classification processes. The first type named “Full Regex”, which uses the complete regex list developed and validated by previous researchers. This classification is used as “ground truth” and is considered equivalent to “human” classification. The second type is called “Partial Regex”, which uses a subset of the a full regex list to classify the data. This type is used to simulate incomplete regex classifiers that are under development. The third

Table 2. Code definition, example regexes and IRR.(All codes were validated at a kappa threshold of 0.65 and a rho threshold of 0.05)

Code	Description	Example	IRR
TECH CONSTRAINTS	Referring to inputs: material, processing method, surfactant, and CNT.	\bPESPVP, \bdry-jet, \bnegative charge, \bsurfactant, ...	0.96 1.00 0.96
PERFORMANCE	Referring to attributes: flus, blood cell reactivity, marketability, cost, or reliability.	\bafforda, \bBCR, \bflux, \bexpensive, \bmarketa, ...	0.88 0.93 0.84
COLLABORATION	Facilitating a joint meeting or the production of team design products.	\bmeeting, \bwe all, \bdiscussion, \bwhat should,...	0.76 0.87 0.76
DESIGN	Referring to design and development prioritization, tradeoffs, and design decisions.	\bfinal decision, \bdecision, \bwent with, \bbased each design,...	0.89 0.86 0.84
DATA	Referring to or justifying decisions based on numerical values, results tables, graphs, research papers, or relative quantities.	\bchart, \bequal value, \bresults, ...	0.94 0.90 0.89
REQUESTS	Referring to or justifying decisions based on internal consultant's requests or patient's health or comfort.	\buser, \bDuChamp, \bPadma, \bsafety, \bhospital, ...	0.88 0.94 0.94

type is the LSTM neural network (see Figure 2), which uses a small part of the data classified by a partial regex classifier to train a predictive model. Details about this neural network model can be found in Cai et al. [2].

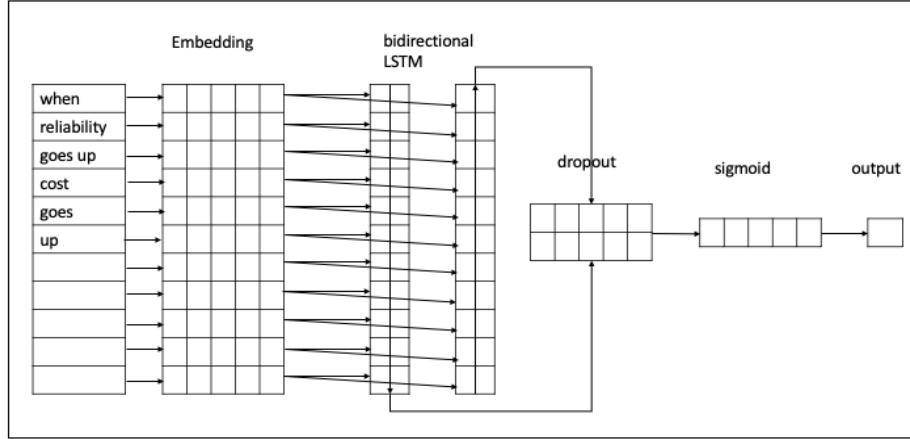


Fig. 2. LSTM neural network for negative reversion.

3.4 Simulation Procedure

To answer our research questions, we simulated coding process under two conditions, one was called “n” condition and the other was called “random” condition. In “n” condition, the neural network modeling was included in the loop to identify negative reversion sets. The “n” condition procedure went as follows.

- 1 Select one of the six codes;
- 2 Select a single regex with base rate not less than 0.001 from the full regex list as lead regex (there are a total of 119 such regexes in the six classifiers);
- 3 Use the selected lead regex as partial regex list;
- 4 Code the data with the partial regex list;
- 5 Compute kappa between full regex coding and partial regex coding;
- 6 Sample 400 items containing at least 10 positive items and train a neural network model;
- 7 Code the data using the neural network model;
- 8 Find the negative reversion set;
- 9 Sample 30 items from the negative reversion set;
- 10 Compute the false negative density in the 30 items;
- 11 If there are no false negative items, go to step 14 ;
- 12 Find the regexes in the full regex list that match any of the false negative item in the 30 items and add the matched regexes into the partial regex list;

- 13 Go to step 4;
- 14 If there are more regexes that have not been selected as lead regex, go to 2;
- 15 If there are more codes, go to 1;
- 16 Stop.

The procedure for the random condition was the same as the “nn” condition, except that the 30 items in each iteration were randomly selected from the regex negative set, instead of NRS.

3.5 Data Aggregation

The false negative density and the achieved kappa in each iteration were averaged over the 119 leading regexes. The 95% confidence intervals were also computed for these means.

4 Results

4.1 False Negative Density

Figure 3 shows the average false negative density in each 30-item iteration under the two conditions (nn and random). In the first iteration, the 30 items from NRS (nn condition) had about 36% false negatives on average, while the random sampling from the regex negative set (random condition) had only around 11% on average. The density in both conditions dropped quickly in the earlier iterations and became about the same from the 6th iteration on. This indicates that the use of NRS is more effective at the beginning stages of regex list development.

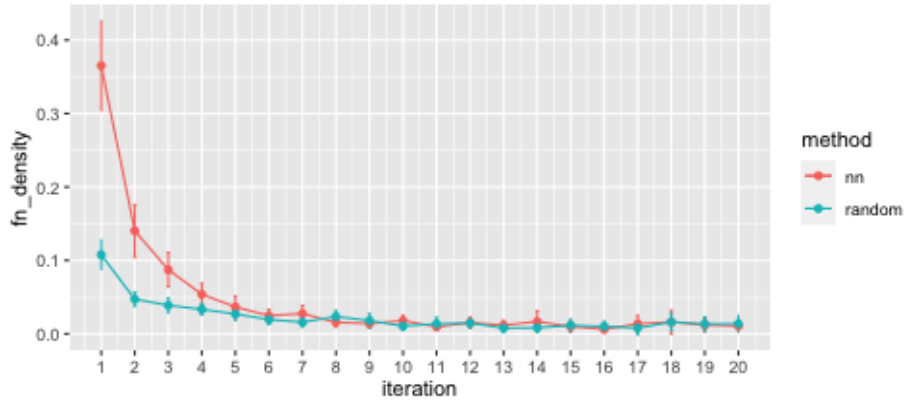


Fig. 3. False negative density in neural network identified negative reversion set (NRS) and in random sample from partial regex negative set

4.2 Required Manual Coding

Figure 4 shows the average kappas with 95% confidence intervals for each iteration, each condition. At iteration 1, for a given code and given lead regex, the partial regex list for both conditions were the same. And that is because both conditions started with the same lead regex. Therefore, the kappas between full regex list and partial list for both conditions were the same. From iteration 2 on, the items in the partial regex list for the nn condition were sampled from negative reversion set, while unsurprisingly, items were randomly sampled from the negative regex set for the random method condition, which is why we see the divergence in kappas from iteration 2. The kappa values for nn condition immediately became higher than the random condition, which implies that, the full regex list could be more quickly identified using the nn method. For example, to get a kappa $\kappa = 0.8$, the nn condition needed to iterate up to the 4th iteration, which required 90 items; while the random condition needed to iterate up to 9th iteration, which required 240 items. Table 3 shows the estimated average number of required items for the two conditions. Overall, the nn to random ratio is less than 0.5. For higher kappa, the ratio is even smaller.

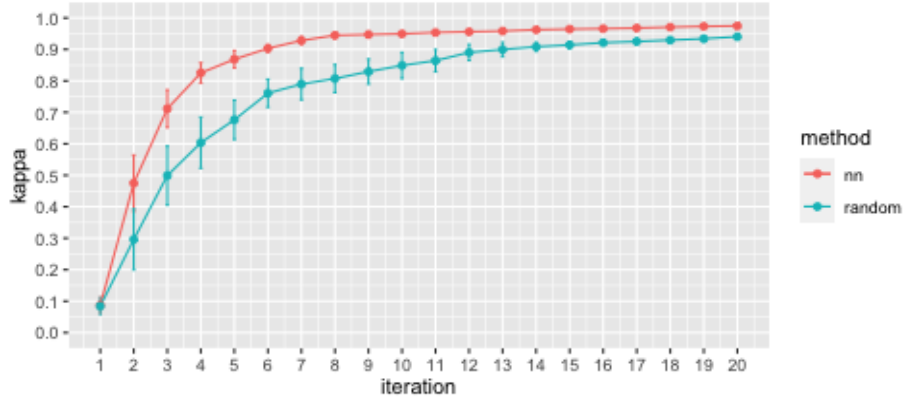


Fig. 4. Kappa between full regex coded data and partial regex coded data at each iteration for two methods

Table 3. Estimated number of manually coded items for developing regex with (nn)/without (random) NRS.

kappa	0.5	0.6	0.7	0.8	0.9
nn	45	60	70	90	180
random	90	120	150	240	420

5 Conclusions and Limitations

In this paper, we explored the use of the negative reversion technique in regex classifier development through a simulated coding procedure. We obtained two major results. One result is about the change of the density of the false negative items in the iterative procedure. We found that the false negative items were dense in early regex development but quickly dropped as the procedure iterated. This result has multiple implications. One is that the use of negative reversion set is more effective at the beginning of the iterations. When the regex list gets close to complete, the benefit from using negative reversion set is limited. Another implication is that, the dropping rate of the false negative density may be used as an indicator of the convergence of the iterative procedure. Namely, when the density and the dropping rate become small, further iteration may not be needed and the quality of the regex list under development may be high enough. However, if we consider the quality of the regex list as a function of density change, the code base rate needs to be taken into account, because for codes with smaller base rate, density is smaller and the change rate may be different. Further exploration on this issue is needed but is beyond the scope of this paper.

In addition, we found that sampling from the negative reversion set could be two or three times more efficient than random sampling from the regex negative set. In other words, using negative reversion set could largely reduce the number of items human needs to manually code when developing regex based classifiers.

Our results are averaged over 119 lead regexes from six classifiers with fairly tight confidence intervals providing evidence for the robust of our conclusions and the use of the NRS more broadly. However, when this technique is applied to other data sets, several factors may impact the results. The first factor is the code base rate. The six codes used in this paper have base rates ranging from 7% to 16%. When this technique is applied to codes with base rates far away from this range, the density drop rate, the kappa growth rate, and, more importantly, the number of required manual coding items could be different. The second factor is the number of items selected in each iteration. In this paper, we chose 30 items in each iteration just because it was easier for demonstrating performance. However, this number may affect the results. A smaller number may further reduce required manual coding. Meanwhile, a smaller number also implies more frequent updating of the neural network models, which takes computing time. The R code we ran on a mac laptop needs about 1 minute to update an LSTM neural network model with 400 items in the training data. When the computing time is an issue, less frequent neural network updating could be considered. However, if the neural network updating can be further optimized and becomes fast enough, fewer items, or even single item per iteration could be considered.

Training size for neural network models is another factor that may impact the results. In this paper, we chose training size as 400 items and required that there are at least 10 positive items among these 400 items. How much the results rely on this choice is unclear. Future work in finding ways to optimize this choice is under consideration. One note on this is that, we don't believe using a very large

training size will give us better results. The reason is that, larger the training size is, closer the neural network model will be to the regex classifier. That is, with a large training size, the training data could "force" the neural network to have a high level of agreement with the regex classifier and thus reduce the density of false negatives in the negative reversion set. We also point out that, any text classification models, not necessarily neural network models, could be used to construct the negative reversion set. It wouldn't be a surprise to us that there are models work better than the one we used in this paper.

As a final note, while this paper only considered false negatives, similar logic can be applied to construct "positive reversion set". In practical use, we believe both positive reversion set and negative reversion are useful. However, the negative reversion set could be more helpful when a code has a low base rate.

Acknowledgements This work was funded in part by the National Science Foundation (DRL-1661036, DRL-1713110, DRL-2100320, LDI-1934745), the Wisconsin Alumni Research Foundation, and the Office of the Vice Chancellor for Research and Graduate Education at the University of Wisconsin-Madison. The opinions, findings, and conclusions do not reflect the views of the funding agencies, cooperating institutions, or other individuals.

References

1. Bai, X.: Text classification based on lstm and attention. In: Thirteenth International Conference on Digital Information Management (ICDIM). pp. 29–32 (2018)
2. Cai, Z., Marquart, C., Shaffer, D.: Neural recall network: A neural network solution to low recall problem in regex-based qualitative coding. In: Mitrovic, A., Bosch, N. (eds.) Proceedings of the 15th International Conference on Educational Data Mining. pp. 228–238. International Educational Data Mining Society, Durham, United Kingdom (July 2022). <https://doi.org/10.5281/zenodo.6853047>
3. Cai, Z., Siebert-Evenstone, A., Eagan, B., Shaffer, D.W., Hu, X., C., G.A.: ncoder+: A semantic tool for improving recall of ncoder coding. In: Advances in Quantitative Ethnography: ICQE Conference Proceedings. pp. 52–65 (October 2019)
4. Chen, N.C., Drouhard, M., Kocielnik, R., Suh, J., Aragon, C.R.: Using machine learning to support qualitative coding in social science: Shifting the focus to ambiguity. *ACM Trans. Interact. Intell. Syst.* **8**(2), 9:1–9:20 (June 2018). <https://doi.org/10.1145/3185515>, <https://doi.org/10.1145/3185515>
5. Chesler, N., Ruis, A., Collier, W., Swiecki, Z., Arastoopour, G., Shaffer, D.: A novel paradigm for engineering education: virtual internships with individualized mentoring and assessment of engineering thinking. *Journal of Biomechanical Engineering* **137**(2), 1–8 (2015)
6. Eagan, B., Brohinsky, J., Wang, J., Shaffer, D.: Testing the reliability of inter-rater reliability. In: Proceedings of the Tenth International Conference on Learning Analytics and Knowledge. pp. 454–461 (2020)
7. Eagan, B., Swiecki, Z., Farrell, C., Shaffer, D.: The binary replicate test: Determining the sensitivity of cscl models to coding error. In: Proceedings of the 13th International Conference on Computer Supported Collaborative Learning (CSCL). pp. 328–335 (2019)

8. Gautam, D., Swiecki, Z., Shaffer, D.W., Graesser, A.C., Rus, V.: Modeling classifiers for virtual internships without participant data. In: Proceedings of the 10th International Conference on Educational Data Mining. pp. 278–283 (2017)
9. Georgieva-Trifonova, T., Duraku, M.: Research on n-grams feature selection methods for text classification. In: IOP Conference Series: Materials Science and Engineering. vol. 1031, p. 012048. IOP Publishing (2021)
10. Glaser, B., Strauss, A.: The discovery of grounded theory: Strategies for qualitative research. Aldine, Chicago (1967)
11. Landauer, T.K., Foltz, P.W., Laham, D.: An introduction to latent semantic analysis. *Discourse processes* **25**(2-3), 259–284 (1998)
12. Shaffer, D.: *Quantitative Ethnography*. Cathcart Press, Madison, WI (2017)
13. Shaffer, D.W., Ruis, A.R.: How we code. In: *Advances in Quantitative Ethnography: ICQE Conference Proceedings*. pp. 62–77 (February 2021)