

Explaining the Hardest Errors of Contextual Embedding Based Classifiers

Anonymous ACL submission

Abstract

We seek to explain potential causes for incorrect classification of the most challenging documents, namely, documents that no classifier using state-of-the-art, very semantically-separable contextual embedding representations managed to accurately predict. To do so, we propose a misclassification taxonomy of incorrect predictions, which we used to perform qualitative human evaluation. We posed two (research) questions, achieving a high inter-evaluator agreement of 81.7%. We worked with three sentiment analysis datasets, two in the movie reviews domain and a third one containing product reviews. We quantified answers per category in our taxonomy across all datasets and computed their proportion. Differences were observed between the product and movie review domains, such as the prevalence of ambivalence in product reviews and sarcasm in movie reviews. Our analysis also revealed an unexpectedly high rate of human mislabeling in the datasets and a significant number of model errors that we cannot yet explain. To ensure reproducibility, our documentation, code, and datasets can be accessed on GitHub.¹

1 Introduction

In a scenario where the amount of user-generated content is growing exponentially, automatic text classification (ATC), one of the fundamental tasks of machine learning, plays a vital role in enabling the automatic categorization of texts into different semantic groups based on their distinctive characteristics (Li et al., 2022; Galke and Scherp, 2022).

The state-of-the-art in ATC is currently provided by Attention-Based Transformer methods (e.g., BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), BART (Lewis et al., 2019)), which produce contextual representations of words and documents. Indeed, in (de Andrade et al., 2023), the authors

show that these contextual representations are so (semantically) separable in the embeddings space that any classifier using them achieves similar effectiveness, no matter how simple (e.g., a Nearest-Centroid classifier) or complex it may be (e.g., a Gradient Boosted Decisions Tree (GBDT) or a Support Vector Machines). Some of the results obtained in that study are the highest (state-of-the-art) ever reported in the literature for effectiveness (e.g., Macro-F1) in several of the experimented datasets.

With such powerful text representations and results, sometimes achieving or even exceeding human parity (Hassan et al., 2018; Yan et al., 2023), a main question that arises is: *Are we reaching the limits of what can be automatically classified by a machine learning model?*

In this article, we delve deep into this question by analyzing the reasons for misclassification by classifiers using these powerful contextual representations. We go one step further to advance the literature and look into the **hardest cases**, i.e., documents that none of the strongest classifiers explored in the aforementioned study, using contextual embedding-based representations, were able to classify correctly.

A thorough review evidenced that such type of error or misclassification analysis is rarely performed in the literature, with a few exceptions (Martins et al., 2021). Misclassification analysis serves the purpose of revealing the how’s and why’s behind model (or human) failure. One of the main difficulties in performing such an analysis is the lack of standardized methodologies and methods for doing so. Accordingly, one of our contributions is the proposal of a **misclassification taxonomy** capable of categorizing incorrect classification *upon classifiers application*.

We propose and evaluate a *taxonomy of errors* using a sample of the documents for which none of the classifiers can achieve correct predictions. Due to their simplicity compared to more complex tasks,

¹Anonymous

081	such as topic categorization, we initially focus on	• Across all assessed datasets, the predominant category of reason for errors (> 50%) is “Sufficient information with model failure”, a result that can be potentially leveraged for model enhancement.	129
082	sentiment analysis (binary) tasks, adopting BERT		130
083	to generate the contextual representations for the		131
084	documents. We evaluate the proposed taxonomy		132
085	with a different sample of erroneous documents,		
086	using human evaluators with different backgrounds	• The evaluators found a significant amount of mislabeling in the datasets by humans – i.e., evaluators considered that the model provided a correct label but the human mislabeled the document – in around 33% of the documents in the product dataset and 16% in one of the movie datasets. As there was a high agreement among evaluators regarding those mislabelings (most cases with 4 evaluators having the same opinion), they are worth further investigation.	133
087	to assess how effective and useful the taxonomy		134
088	is to explain the errors.		135
089	Unlike previous work (Martins et al., 2021) –		136
090	which focuses on characterizing and assessing the		137
091	impact of “hard” instances in the effectiveness of		138
092	the polarity detection task using a single dataset		139
093	of movie reviews and unconcerned about textual		140
094	representation – we here focus on analyzing and		141
095	quantifying the reasons for the misclassification		142
096	of the hardest documents by all machine learning		
097	methods using some of the most separable representations in the literature. For this, we use datasets from two domains: movie and product reviews. We also contrast and compare the results in these two domains, gathering insights into the differences in the type of errors found in each of them.	• In movie reviews, sarcasm ² (> 23% of the cases) is a major reason for model error. We believe this is a particular characteristic of this domain.	143
098			144
099			145
100			
101		• In product reviews, ambivalence (40% of the cases) is the main reason for model error.	146
102			147
103	The main questions we seek to answer are:		
104	RQ1 <i>Is the proposed taxonomy for misclassification effective for misclassification analysis?</i> To answer RQ1, we analyze evaluators’ responses regarding their level of agreement – the higher the agreement, the more effective the taxonomy. We analyze inter-evaluator disagreement and correlate that with hardness in classifying.	• In at least 41% of the evaluations, evaluators reported sufficient information in the text, with the model failing in its prediction for reasons we cannot yet explain.	148
105			149
106			150
107			151
108		In sum, this paper’s main contributions include:	152
109			
110		1. A methodology for finding the hardest (most challenging documents) based on misclassifications by all classifiers using contextual embedding representations	153
111	RQ2 <i>Can the text misclassification taxonomy be used to reveal the main reasons for misclassifications? Are there significant differences in the results among different domains?</i> In RQ2, drawing on the consensus achieved, we quantify and analyse the main reasons and causes for the misclassifications, outlining potential differences between domains.		154
112			155
113			156
114		2. The development and evaluation of a taxonomy for categorizing the main causes of automatic classifiers’ misclassifications	157
115			158
116			159
117		3. A deep analysis of the results obtained upon applying the taxonomy to three different datasets in two different domains. The results may have interesting implications for the improvement of the next generation of textual classifiers and representations	160
118			161
119	Our experiments engaging eight human evaluators with two different backgrounds (Computer Science and Linguistics) and three datasets, two in the movie reviews domain and one in the product reviews domain, revealed that:		162
120			163
121			164
122			165
123		4. A release of a new dataset of challenging documents manually annotated by humans.	166
124	• The developed taxonomy proved effective, with an inter-evaluator agreement of over 81% for error category. This suggests that evaluators find it relatively easy to identify classification errors using the proposed taxonomy.		167
125			
126		This paper is organized as follows. Section 2 presents related work. Section 3 describes our methodology. Section 4 reports and discusses our results, followed by our conclusions in Section 5.	168
127			169
128			170
			171

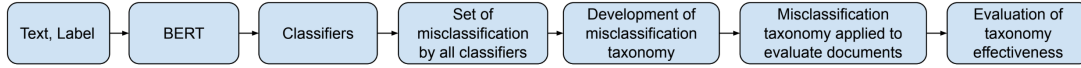


Figure 1: Flowchart of the Evaluation Methodology.

2 Related Work

Reasons for models’ misclassifications have been investigated for texts and images. In Lee et al. (2017), five categories for misclassification of objects in images are explored: 1) “similar labels”: the term³ representing the predicted object in the image is not in the ground truth (GT) but is semantically similar to the GT; 2) “not salient”: the predicted object exists in the image but is not present in the GT; 3) “challenging images”: the GT is challenging even for a human being; 4) “incorrect GT”: incorrect annotation by humans; and 5) “incorrect prediction class”: machine prediction is incorrect but with sufficient information in the image for humans to detect. Unlike Lee et al., we propose a taxonomy for text classification using textual data and associated labels as input. Furthermore, we focus on analyzing the hardest cases for the automatic classifiers.

In another study aimed to categorize prediction failures, (Meek, 2016) defines four error categories: 1) “mislabeling errors”: human labeling errors; 2) “representation errors”: limitations in the feature set used for evaluation; 3) “learner errors”: prediction errors when there is sufficient information for accurate classification; 4) “boundary errors”: correct predictions could be achieved by adding more examples, indicating absence of labeled examples for a specific class in the training set. Defining “boundary errors” is hard in our case due to our use of LLMs and the lack of details regarding their training.

A different approach is pursued by Pandey et al. (2022), who focus on evaluating human labeling by assessing the impact of (i) the number of documents; (ii) the time allocated to evaluators; and (iii) the order in the sequence of annotations in the labeling task. They also assess the effects of evaluator memory retention for the task. Unlike the above works, our study focuses on a set of (test) documents that have been misclassified by *all* classifiers using state-of-the-art, very separable (contextual) representations.

²In this work, we do not differentiate between irony and sarcasm as others works do (e.g. (Frenda et al., 2023a))

³The set of true labels (*Ground Truth (GT)*) is the set of terms that textually describe the objects in the image.

The closest work to ours is Martins et al. (2021), who analyze hard instances of text, i.e., instances assigned to a “neutral” – when polarity is not clearly defined – or “discrepant” category – when polarity differs from its associated labeling. Their study centers on characterizing and evaluating the influence of challenging cases on the classifier’s effectiveness when performing polarity detection using a single movie review dataset. In our study, we concentrate on analyzing and quantifying the factors contributing to misclassifications of the hardest documents of all, i.e., those misclassified by all classifiers using the current best (contextual) text representations. Compared to Martins et al.’s, our taxonomy has more fine-grained categories, such as sarcasm and ambivalence (contrasting opinions in the same text). We have additional goals as well, such as validating our taxonomy and contrasting the results in different datasets and domains, running qualitative experiments engaging human evaluators with different backgrounds.

3 Experimental Methodology

Our methodology, comprising seven steps, is summarized in Figure 1. The text and label for each document are used as input for fine-tuning a BERT model, resulting in an encoder that produces contextual embeddings vectors representing the documents using the CLS approach. We employ various classifiers with these embeddings as input, exploring different underlying techniques. From this set of classifiers, we select the set of documents for which none of the classifiers can produce correct predictions (according to the assigned labels in the datasets). Within this set, we sample documents for analysis to outline misclassification categories (“Development of the misclassification taxonomy” in Figure 1), which human evaluators will use to evaluate documents (“Application of misclassification taxonomy to evaluate documents” in Figure 1) in a second sample different from the first one. Upon the application of the taxonomy, we quantify the results and evaluate its efficiency. Details of the steps follow.

3.1 Datasets

Our study draws on three datasets developed for binary sentiment classification. Each dataset

was constructed with a text and an associated sentiment label. The first dataset comprises customers’ reviews after purchasing products on Amazon’s website (Keung et al., 2020). Products are assigned a rating from 1 to 5 stars by customers. We collected reviews containing ratings of 1 and 2 stars and labeled them as negative, while reviews containing ratings of 4 and 5 stars were labeled as positive. We discarded reviews with 3 stars (deemed neutral). The second (PangMovie (Pang and Lee, 2004)) and the third (VaderMovie (Ribeiro et al., 2016)) datasets are made up of movie reviews, comprising a text and a sentiment label (positive or negative). Table 1 presents some statistics of the datasets. As we can see, class distribution into positive and negative instances is basically balanced in the three datasets.

Dataset	Documents	Avg words	Positive	Negative	Classifier Macro-F1
Amazon	168000	33	84000	84000	94.2
PangMovie	10662	19	5331	5331	86.8
VaderMovie	10568	19	5242	5326	89.1

Table 1: Datasets Statistics

3.2 Data Representation

We fine-tuned BERT, adapting this LLM to the specific domain of sentiment classification using the texts and labels in our datasets. The aim is to improve the representation and enhance the model’s effectiveness for sentiment classification. The model’s fine-tuning produces an encoder, which generates CLS-based 768-dimensional embedding vectors to represent the documents. As discussed in (de Andrade et al., 2023), this fine-tuning process is fundamental to ensure the quality of the representation and the separability (into semantic classes) of the generated embedding space.

To perform fine-tuning, we used the suggested hyper-parameterization, fixing the learning rate with the value 2×10^{-5} , the batch size with 64 documents, adjusted the model to five epochs and set the maximum size of each document to 256 tokens.

In our experiments, we employ a five-fold stratified cross-validation procedure. This means that all procedures of fine-tuning, training, and optimizing the classifiers’ parameters with the validation sets are repeated five times. Reported results correspond to the average of the five test folds.

We used BERT in our study, but other Transformers can be easily applied within our methodology. Indeed, experiments in (de Andrade et al., 2023) show that the contextual representations produced by different transformers (e.g., Roberta, BART)

are quite similar regarding class separability, the main aspect driving our evaluations.

3.3 Text Classifiers

For document classification, we used the textual representations generated by the Transformer as input to four of the strongest classifiers used in (de Andrade et al., 2023), namely: KNN, Random Forests (RFs), Support Vector Machines (SVMs) and Logistic regression (LR). Indeed, despite using different rules and heuristics, the effectiveness of these classifiers (and of all other classifiers tested in (de Andrade et al., 2023)) is basically the same in all tested datasets when using the contextual embedding representations. This is due to the fact that these representations are already so semantically (by class) separated in the embeddings space that the employed classifier has little work. The Macro-F1 of these classifiers (remember that all classifiers are basically tied when using contextual embedding representations according to (de Andrade et al., 2023)) is shown in Table 1.

Our decision was to explore classifiers based on different approaches – multiple decision rules (RFs), local neighborhoods (kNN), global maximum margins (SVMs and LR) – so that if all of them misclassify the same document, this can be ascribed to the misclassified document being hard to classify. *And we do want to understand the reasons why!*

Hence, we selected the set of documents that all classifiers misclassified in the three datasets. A sub-sample from this set was used as a basis for devising our taxonomy and a different (disjoint) sub-sample for actual evaluation, as described next. Table 2 shows the number of misclassified instances by all classifiers. As we can see, there is no significant skewness in the distribution of positive and negative misclassified documents in the three datasets. We took a random sample of 60 misclassified documents from each dataset for evaluation, and the results are presented in Section 4.

Dataset	Misclassification	Positive	Negative
Amazon	216	115	101
PangMovie	120	54	66
VaderMovie	85	37	48

Table 2: Set of misclassifications by all classifiers.

3.4 Taxonomy Development

We conducted a preliminary round of assessment using a set of 15 randomly selected documents from PangMovie and Amazon. During this round,

	A	B	C	D	E	F
1	ID	Text	Label assigned by human	Label assigned by model	Who misclassified the text?	Based on your answer to question 1, why do you think the model (or the human) misclassified the text?
2	3669	Rope completely broke off after a couple of months	Positive	Negative		

Figure 2: Screenshot of form provided to evaluators.

E

Who misclassified the text?

Human

Model

Don't know

F

Based on your answer to question 1, why do you think the model (or the human) misclassified the text?

Sarcasm

Ambivalence

Insufficient Information

Sufficient information with model failure

Sufficient information with human failure

None of the above

(a) Question 1

(b) Question 2

Figure 3: Screenshot of form with questions and options for selection.

Category	Description
Sarcasm	Text contains ironic expressions (words that are the opposite of what one means), humorous expressions, figurative language (metaphors, idioms)
Ambivalence	Text contains both positive and negative opinions
Insufficient Information	Text is very brief and lacks information to assign the predominant sentiment
Sufficient information with model failure	Text contains sufficient information but model was unable to correctly assign the predominant sentiment
Sufficient information with human failure	Text contains sufficient information but human was unable to correctly assign the predominant sentiment
None of the above	None of the above categories can be said to account for the misclassification

Table 3: Misclassification categories and their description

Text	Category
A film of precious increments artfully camouflaged as everyday activities	Sufficient information with model failure
Rope completely broke off after a couple of months	Sufficient information with human failure
Final verdict: you've seen it all before.	Sarcasm
Expensive but won't oxidize metal. Maybe better than soap	Ambivalence

Table 4: Texts illustrating misclassification categories included in evaluators' guidelines

we convened to discuss potential sources of misclassification, aiming to better comprehend the reasons behind incorrect predictions. Through this process, we agreed upon a set of potential reasons, which represent the bulk of the categories in our taxonomy of errors. We carried out a subsequent evaluation with another set of 15 documents from each dataset, refining definitions, instructions, and the evaluation process. Upon concluding this iteration, we excluded all documents used in the preliminary stage and proceeded with a new evaluation. We randomly selected 60 samples from each dataset for manual human evaluation.

3.5 Distribution of documents

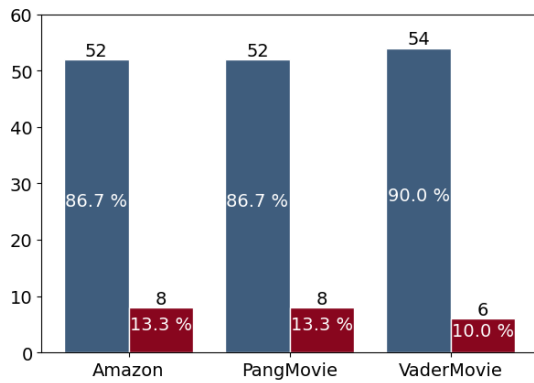
To evaluate the selected texts, we recruited eight participants with expertise in Computer Science and Linguistics, all with prior experience in annotation tasks for NLP. The participants comprised two professors holding a Ph.D. in Computer Science, one with a Ph.D. in Linguistics, and five students pursuing their bachelor's or master's degrees. All evaluators had advanced

proficiency in English, which attests to their ability to perform the evaluation task.

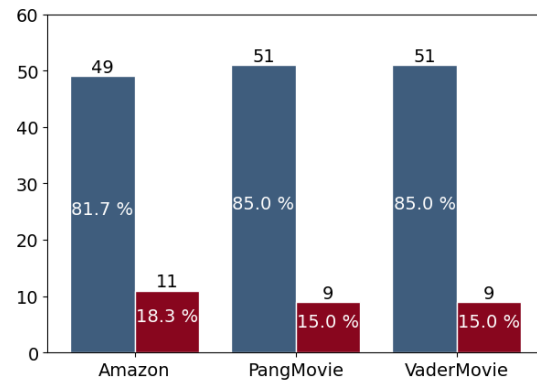
Each participant was assigned 30 out of the 60 documents in each of the three datasets, totaling 90 documents per evaluator. Each document was assigned to be evaluated by four participants, two having a computer science background and the other two having a linguistics background. The decision to assess each document by two evaluators from each field was meant to enable quantification of (dis)agreement within the same background groups and between the two groups with different backgrounds. Section 4 presents our results considering all four evaluators. Results regarding impact of evaluators' background are presented in our Appendix.

3.6 Evaluation Form

Individual forms were created for each evaluator and shared on a web cloud provider, ensuring evaluators could not access each other's forms. Our evaluation form comprised four tabs, the first containing instructions on how to evaluate the



(a) Q1: Consensus and No Consensus



(b) Q2: Consensus and No Consensus

Figure 4: Total Consensus and No Consensus for Questions 1 and 2.

documents and the remaining ones having one sample of documents per tab, each line containing a document and the categories to be assigned to it.

Figure 2 shows the form provided to evaluators, with columns for text ID, text to be evaluated, label assigned by a human, and label assigned by the machine model. Two additional columns were assigned to be filled in by evaluators with their answer to two questions: (i) “Who misclassified the text?”, for which one out of three options could be chosen: “Model”, “Human”, and “I don’t know” (see Figure 3a); and (ii) “Based on your answer to question 1, why do you think the model (or human) misclassified the text?”, for which one out of six options could be chosen, as shown in Figure 3b. Table 3 provides a description of the available options.

3.7 Categories To Evaluate Misclassification

The second question in our evaluation form required the evaluator to choose a category that could account for either the model’s or human misclassification. The instructions tab provided evaluators with examples of each category, some of which are presented in Table 4. The first row shows an example of a text misclassified due to the model’s failure in spite of sufficient information to make the correct decision. In this case, the model assigned a negative sentiment, though the text contains a positive opinion “precious increments artfully....”. The second row shows an example of misclassification due to human failure with sufficient information. The wording “completely broke off” indicates a negative opinion, but it is labeled as positive. The third row is a misclassification ascribed to sarcasm, where “seen it before” is a negative opinion ironically expressed. The fourth row exemplifies a misclassification

due to *ambivalence*, where despite the negative word “expensive”, there are two favorable opinions (“won’t oxidize” and “better than soap”).

4 Results

Documents were assessed by four evaluators. Question 1 required selecting one out of three alternatives, whereas Question 2 had six alternatives. Consensus was defined as one of the alternatives having the *majority* of votes – 4, 3, or 2 votes.⁴ If there was no majority of votes for a document, we classified it as “No consensus”.

4.1 Taxonomy effectiveness

To answer our first research question: “Is the proposed taxonomy for misclassification effective to be used for misclassification analysis?”, we analyzed the responses from questions 1 and 2 provided by the evaluators. We consider a taxonomy effective if there is high consensus among evaluators upon the defined categories and if there is low consensus in a category that has no definition, in our case, “Don’t know” for Question 1 and “None of the above” for Question 2.

Figure 4 shows the consensus percentages obtained for Questions 1 and 2 in the three evaluated datasets. For Question 1, out of 60 documents, 54 attained high inter-evaluator agreement in VaderMovie, and 52 in Amazon and PangMovie. In other words, in at least 86.7% of the cases (52/60), consensus was achieved in some category defined for Question 1 in the three evaluated datasets, implying low difficulty for evaluators to define a type of misclassification.

⁴In the case of two votes, provided that the remaining two alternatives have one vote each.

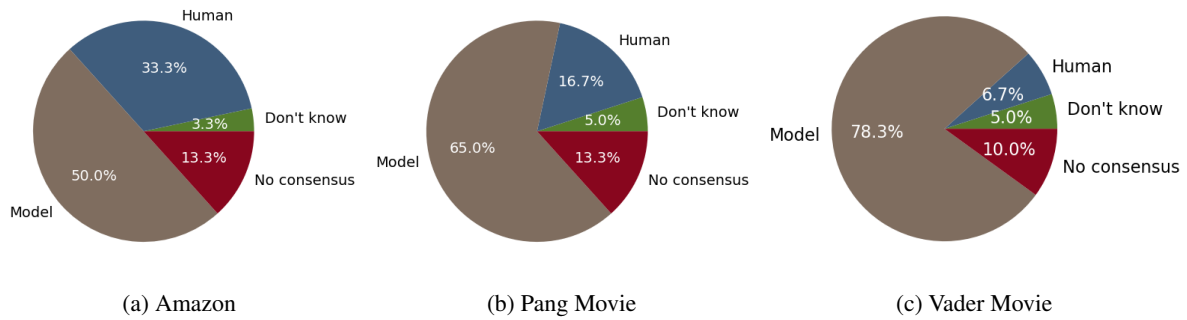


Figure 5: Percentages for answers to Question 1 in the three datasets.

We break down those numbers in the Appendix to show the consensus distribution per document and per evaluator background. As shown there, the vast majority of the documents had the same categorization assigned by 4 or 3 evaluators, emphasizing the high agreement and taxonomy effectiveness.

Figure 4b shows the consensus percentages for Question 2. It is important to bear in mind that in Question 2, six options were available, likely leading to a higher difficulty in achieving agreement. Nonetheless, we can observe a high consensus in all datasets for this question, with the lowest value being obtained in the Amazon dataset, 49 out of 60 documents reaching at least 81.7% consensus. As also shown in Figure 5, “No Consensus” was below 14% for Question 1 and below 19% for Question 2. In the Appendix, we show examples (in Table 5) of documents that posed difficulties for evaluators.

4.2 Response Analysis

Upon concluding our consensus analysis, we explored the responses given by evaluators. Figure 5 shows that 50% of the misclassifications in the Amazon dataset were ascribed to the model. This is even higher in the movie datasets, emerging as the main misclassification reason in 65% of the cases in PangMovie and almost 80% in VaderMovie. Percentages for the option “Dont know” are very low in all datasets. Together with the option “No consensus”, they achieve at most 18.3% in PangMovie (and 16.3% and 15% in Amazon and VaderMovie, respectively) of all analyzed documents in all datasets.

Though lower than Model’s errors, the percentage of errors ascribed to the “Human” category is significant, mainly in the Amazon dataset (33%). This means that in 33% of the misclassifications, 3 or 4 evaluators (in the majority of the cases) considered that the model classified the document correctly and the human provided a wrong label.

Though lower in the movie domain, human mislabeling is not negligible in the two datasets – 16.7% in PangMovie and 6.7% in VaderMovie. This relatively high percentage of human mislabeling merits further investigation in future studies, though manual labeling has been acknowledged as a complex task, very much prone to errors (Zhu et al., 2023).

Figure 6 presents the results for Question 2. Consensus cases show clear differences between the two domains. The main reason for misclassifications in Amazon was “Ambivalence”, with 30% of the cases, whereas “Sarcasm” is almost non-existent. In the product review domain, texts tend to be more focused on features of a product, so called *aspects*, there being less irony or sarcasm in the reviews. Moreover, most misclassifications occurred when the text concomitantly expressed both positive and negative opinions about product aspects. These are a challenge both for the model and the human to predict the “correct polarity” for the document. This raises the question as to whether there is indeed a single correct polarity label for these documents or whether different aspects of the product should be given different polarity labels (Brauwerters and Frasincar, 2022).

In the movie domain, we see a different result, with “Sarcasm” pointed out as the main reason for misclassification in VaderMovie and the second main one in PangMovie, almost tied with “Ambivalence“. We believe sarcasm is a particular characteristic of the movie review domain, possibly due to the fact that reviewers assess artistic productions and feel the need to use figurative language to express their opinion about them. As in the Amazon dataset, “Ambivalence” is a major reason for misclassifications, especially in PangMovie. This suggests that in the movie domain, reviewers also tend to point out both positive and negative aspects, bringing a challenge both for models and humans

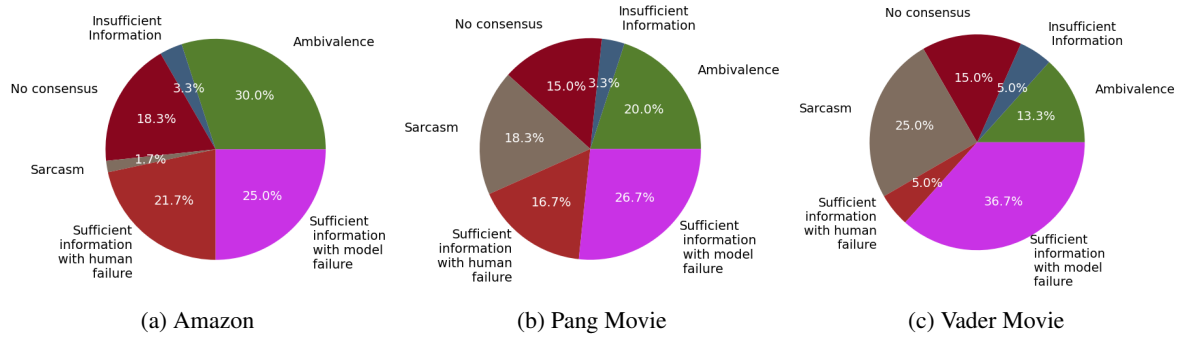


Figure 6: Percentages for answers to Question 2 in the three datasets.

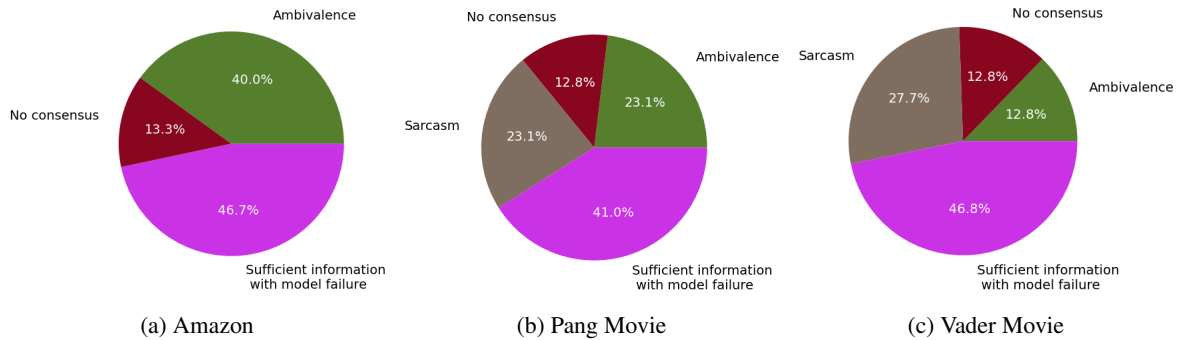


Figure 7: Response analysis for Question 2 in cases where "Model Failure" was selected for Question 1.

to ascribe polarity to the texts. In this sense, *sarcasm detection* (Verma et al., 2021) and *aspect analysis* (Brauwert and Frasincar, 2022) are both interesting lines of investigation worth pursuing.

Finally, "Sufficient Information with Model Failure" is the major reason for errors in both movie datasets (36.7% in Vader and 26.7% in PangMovie) and the second most frequent in the product dataset (25% of the cases) for Question 2. Indeed, if we delve into the results by exploring the answers for which evaluators selected Model failure in Question 1, we can see (Figure 7) that almost half of the error cases lie in this category for the three datasets. This means that evaluators considered there was sufficient information in the text for the model to produce the correct classification, but for some unknown reason (neither "Ambivalence" nor "Sarcasm"), which they could not point out (or which our taxonomy failed to elicit), the model got it wrong. This could be due to lack of enough training data, which is hard to assess when using LLMs, or to borderline cases between two classes, not yet explored. We will address these challenging cases in the future.

5 Conclusion

We addressed the hard task of trying to unveil the reasons why models misclassify the hardest docu-

ments, those which no classifier could correctly classify, despite using the best, most separable state-of-the-art text representations. For this, we devised a taxonomy of errors and ran a qualitative experiment which requested eight evaluators with two distinct backgrounds to use our taxonomy to qualify the errors. The high consensus among the evaluators emerged as a strong finding in our analysis.

Regarding our experimental results, we have found significant differences regarding reasons for misclassifications in the product and movie reviews domains. Sarcasm is very pronounced in movie reviews, while Ambivalence, though occurring in both domains, is more prevalent in product reviews. Further interesting findings were the high proportion of human labeling errors, mainly in the product reviews, and a noteworthy number of Model errors, which we cannot yet account for. We also believe that our methodology serves to identify challenging documents based on evaluator disagreement.

In the future, we will focus on "unexplained" cases and seek a better understanding for the high human mislabelings rate. We also intend to investigate sarcasm detection and aspect analysis as pre-processing steps as a way to deal with the hardest cases. Finally, we intend to apply multi-perspective analysis to our study (Frenda et al., 2023b).

596 Limitations

597 Despite all our contributions, our study has some
598 limitations. Our evaluation targets two domains,
599 three datasets and the task of sentiment analysis.
600 Increasing the number of dataset domains and
601 expanding our analysis to the task of Topic
602 Classification will provide new valuable insights.
603 The size of our evaluation group is relatively small,
604 although this is common in qualitative studies (Sil-
605 verman, 2004). We will increase the number
606 of evaluators in future studies. Our work uses
607 BERT’s contextual representations. Although (de
608 Andrade et al., 2023) shows BERT produces rep-
609 resentations that are as separable (semantically) in
610 the embedding space as representations produced
611 by other Transformers (e.g., RoBERTa, BART),
612 we intend to test our methodology with different
613 Transformers in the future.

614 References

615 Gianni Brauwerters and Flavius Frasincar. 2022. [A survey](#)
616 [on aspect-based sentiment classification](#). 55(4).

617 Claudio M.V. de Andrade, Fabiano M. Belém, Washing-
618 ton Cunha, Celso França, Felipe Viegas, Leonardo
619 Rocha, and Marcos André Gonçalves. 2023. [On the](#)
620 [class separability of contextual embeddings repre-](#)
621 [sentations – or “the classifier does not matter when](#)
622 [the \(text\) representation is so good!”](#). *Information*
623 *Processing & Management*, 60(4):103336.

624 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
625 Kristina Toutanova. 2019. [BERT: Pre-training of](#)
626 [deep bidirectional transformers for language under-](#)
627 [standing](#). pages 4171–4186.

628 Simona Frenda, Viviana Patti, and Paolo Rosso. 2023a. [When sarcasm hurts: Irony-aware models for abu-](#)
629 [sive language detection](#). In *Experimental IR Meets*
630 *Multilinguality, Multimodality, and Interaction - 14th*
631 *International Conference of the CLEF Association, CLEF 2023, Thessaloniki, Greece, September 18-21,*
632 *2023, Proceedings*, volume 14163 of *Lecture Notes*
633 *in Computer Science*, pages 34–47. Springer.

634 Simona Frenda, Alessandro Pedrani, Valerio Basile,
635 Soda Marem Lo, Alessandra Teresa Cignarella, Raf-
636 faella Panizzon, Cristina Marco, Bianca Scarlini, Vi-
637 viana Patti, Cristina Bosco, and Davide Bernardi.
638 2023b. [EPIC: multi-perspective annotation of a cor-](#)
639 [pus of irony](#). In *Proceedings of the 61st Annual*
640 *Meeting of the Association for Computational Lin-*
641 *guistics (Volume 1: Long Papers), ACL 2023, Toronto,*
642 *Canada, July 9-14, 2023*, pages 13844–13857. Asso-
643 ciation for Computational Linguistics.

644 Lukas Galke and Ansgar Scherp. 2022. [Bag-of-words](#)
645 [vs. graph vs. sequence in text classification: Ques-](#)

646 [tioning the necessity of text-graphs and the surpris-](#)
647 [ing strength of a wide MLP](#). In *Proceedings of the*
648 *60th Annual Meeting of the Association for Compu-*
649 *tational Linguistics (Volume 1: Long Papers)*, pages
650 4038–4051, Dublin, Ireland. Association for Compu-
651 tational Linguistics. 652 653

Hany Hassan, Anthony Aue, Chang Chen, Vishal
654 Chowdhary, Jonathan Clark, Christian Feder-
655 mann, Xuedong Huang, Marcin Junczys-Dowmunt,
656 William D. Lewis, Mu Li, Shujie Liu, Tie-Yan Liu,
657 Renqian Luo, Arul Menezes, Tao Qin, Frank Seide,
658 Xu Tan, Fei Tian, Lijun Wu, Shuangzhi Wu, Yingce
659 Xia, Dongdong Zhang, Zhirui Zhang, and Ming Zhou.
660 2018. [Achieving human parity on automatic chinese](#)
661 [to english news translation](#). *ArXiv*, abs/1803.05567. 662

Phillip Keung, Yichao Lu, György Szarvas, and Noah A.
663 Smith. 2020. [The multilingual amazon reviews cor-](#)
664 [pus](#). In *Proceedings of the 2020 Conference on Em-*
665 *pirical Methods in Natural Language Processing*. 666

Han S. Lee, Alex A. Agarwal, and Junmo Kim. 2017. [Why do deep neural networks still not recognize these](#)
667 [images?: A qualitative analysis on failure cases of](#)
668 [imagenet classification](#). *CoRR*, abs/1709.03439. 669 670

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan
671 Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
672 Ves Stoyanov, and Luke Zettlemoyer. 2019. [Bart: De-](#)
673 [noising sequence-to-sequence pre-training for natural](#)
674 [language generation, translation, and comprehension](#).
675 *arXiv preprint arXiv:1910.13461*. 676

Qian Li, Hao Peng, Jianxin Li, Congying Xia, Renyu
677 Yang, Lichao Sun, Philip S. Yu, and Lifang He. 2022. [A survey on text classification: From traditional to](#)
678 [deep learning](#). 13(2). 679 680

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-
681 dar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
682 Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized BERT pretraining](#)
683 [approach](#). *CoRR*, abs/1907.11692. 684 685

Karen Martins, Pedro O.S Vaz-de Melo, and Rodrygo
686 Santos. 2021. [Why do document-level polarity clas-](#)
687 [sifiers fail?](#) In *Proceedings of the 2021 Conference*
688 *of the North American Chapter of the Association for*
689 *Computational Linguistics: Human Language Tech-*
690 *nologies*, pages 1782–1794, Online. Association for
691 Computational Linguistics. 692

Christopher Meek. 2016. [A characterization of predic-](#)
693 [tion errors](#). *CoRR*, abs/1611.05955. 694

Rahul Pandey, Hemant Purohit, Carlos Castillo, and
695 Valerie L. Shalin. 2022. [Modeling and mitigating](#)
696 [human annotation errors to design efficient stream](#)
697 [processing systems with human-in-the-loop machine](#)
698 [learning](#). *International Journal of Human-Computer*
699 *Studies*, 160:102772. 700

Bo Pang and Lillian Lee. 2004. [A sentimental education:](#)
701 [Sentiment analysis using subjectivity summarization](#)
702 [based on minimum cuts](#). page 271–es. 703

704 Filipe N Ribeiro, Matheus Araújo, Pollyanna Gonçalves,
705 Marcos André Gonçalves, and Fabrício Benevenuto.
706 2016. Sentibench-a benchmark comparison of state-
707 of-the-practice sentiment analysis methods. *EPJ*
708 *Data Science*, 5:1–29.

709 D. Silverman. 2004. *Qualitative Research: Theory,*
710 *Method and Practice*. SAGE Publications.

711 Palak Verma, Neha Shukla, and A.P. Shukla. 2021.
712 *Techniques of sarcasm detection: A review*. In *2021*
713 *International Conference on Advance Computing and*
714 *Innovative Technologies in Engineering (ICACITE)*,
715 pages 968–972.

716 Ming Yan, Haiyang Xu, Chenliang Li, Junfeng Tian,
717 Bin Bi, Wei Wang, Xianzhe Xu, Ji Zhang, Songfang
718 Huang, Fei Huang, et al. 2023. Achieving human par-
719 ity on visual question answering. *ACM Transactions*
720 *on Information Systems*, 41(3):1–40.

721 Yu Zhu, Yingchun Ye, Mengyang Li, Ji Zhang, and
722 Ou Wu. 2023. *Investigating annotation noise for*
723 *named entity recognition*. *Neural Comput. Appl.*,
724 35(1):993–1007.

Appendix

Consensus Distribution

This subsection presents evaluator consensus distribution for Questions 1 and 2, analyzed in Section 4.1. Regarding Question 1, as can be seen in Figure 8a, out of the 52 documents that achieved evaluator consensus in the Amazon dataset, 33 reached full agreement among all four evaluators, 16 documents reached full agreement among three, and 3 documents reached full agreement between two evaluators. This points to documents with full agreement among three or four evaluators representing a significant portion of the total number of documents with consensus, in turn demonstrating robustness of our final results. Similar results were obtained for VaderMovie and PangMovie regarding the joint proportion (i.e., sum of the proportions) of evaluations with 4 and 3 agreements.

Regarding Question 2, results show less consensus among the evaluators, which may be due to the number of categories they had to choose from. This is reflected in the graphs in Figure 9. The Amazon dataset showed higher consensus among a higher number of evaluators, possibly accounted for by the type of review - product review. As movie reviews assess artistic productions and implicate more sarcasm and figurative language, full consensus is harder to achieve, though still attainable.

Regarding documents for which there was no consensus among the evaluators (Figure 4a), there are 8 for the Amazon dataset, 8 for the PangMovie dataset and 6 for the VaderMovie dataset. As for question 2 (Figure 4b), there are 11, 9, and 9 documents without consensus for Amazon, Pang Movie, and Vader Movie datasets, respectively. To exemplify challenging documents, we provide three examples from each dataset in the “No consensus” category for Question 2, as shown in Table 5.

The first row in Table 5 shows an Amazon product review where the text begins positively but then brings in an issue with the product. Row 4 shows a movie review from the Pang Movie dataset, where the reviewer uses the words “distended” and “dragging”, creating uncertainty for categorization. Row 6 shows a series of references to other movies and directors, which requires previous knowledge of those movies and their evaluations. Therefore, we believe that the methodology of this study serves to identify challenging documents based on evaluator disagreement.

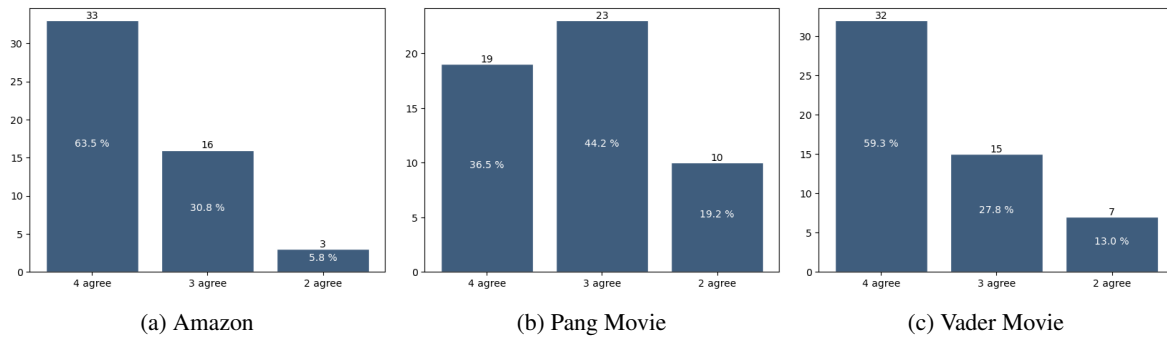


Figure 8: Consensus for Question 1.

Inter-evaluator agreement for Question 2 in cases of “Human Mislabeling”

In Figure 10, similar to Figure 7, we have the quantification of Question 2, but now restricted to the documents that were evaluated as human mislabelings in Question 1. In other words, documents the evaluator considered to have been correctly classified by the Model but which had been incorrectly labelled by the human (positive or negative). We can observe that, in general, the number is lower; for instance, in the Amazon dataset, we have 20 documents evaluated as mislabeled by humans.

Additionally, we can observe a high prevalence of the category human mislabeling with sufficient information in the text, which corresponds to 65% in the Amazon and 70% in the Pang Movie datasets. This means that the evaluator considered the document to have been mislabeled by the human, despite there being sufficient information in the text for the human to choose the ‘right’ label according to the evaluator’s assessment.

Regarding the VaderMovie dataset, numbers are low, which may bias some proportions – there are only four mislabeled documents evaluated as human mislabelings, and only 1 sample was considered to have been mislabeled by the human, despite there being sufficient information in the text for the human to choose the ‘right’ label.

Differences in Evaluation carried out by Computer Scientists and Linguists

We carried out an additional analysis focusing on evaluators’ background. Since each document was rated by two evaluators having a Linguistics background and two a Computer Science one, we examined our data to investigate differences ascribable to evaluators’ backgrounds. Figure 11 represents the quantification of the responses to question 1 by evaluators having a Computer

Science background (11a, 11b, 11c) and a Linguistics one (11d, 11e, 11f), in which there was inter-evaluator agreement of the two evaluators. We can notice that evaluators’ backgrounds had little impact on the results for all datasets.

813
814
815
816
817

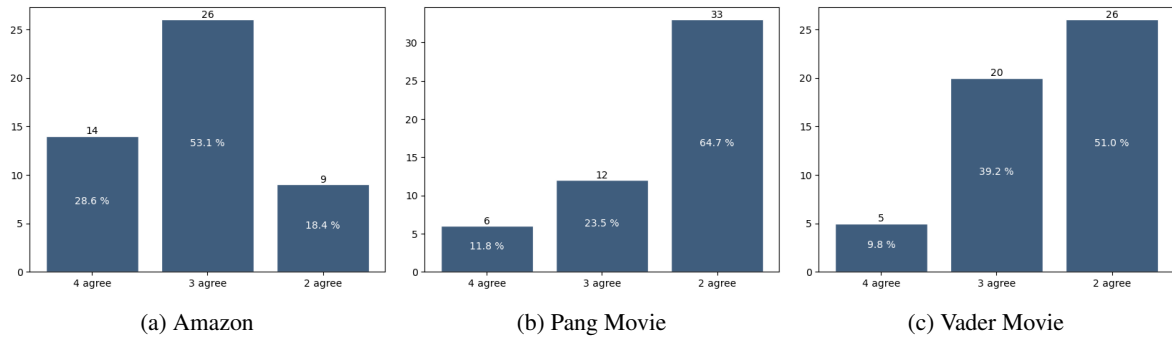


Figure 9: Consensus for Question 2.

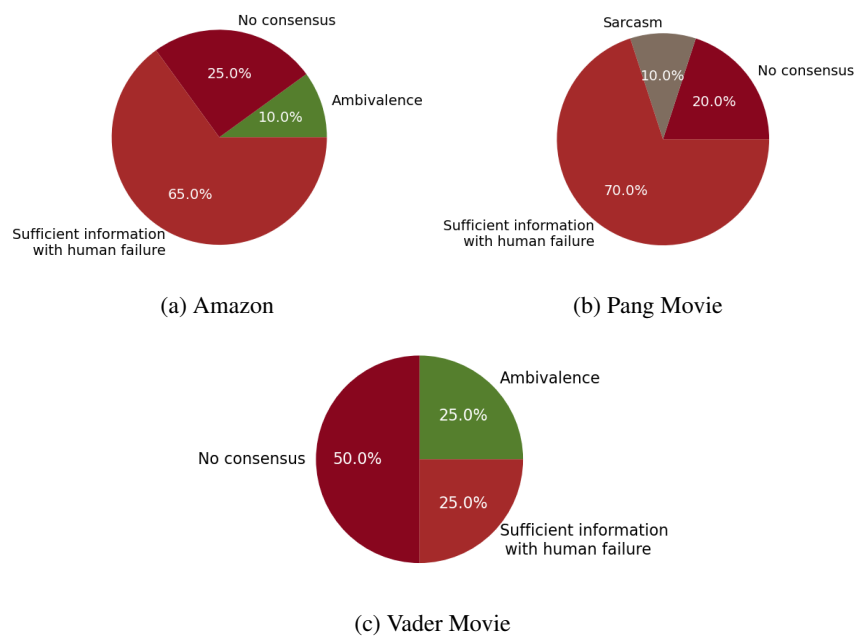


Figure 10: Results for Question 2 in cases where Response to Question 1 was "Human' Failure'.

Text	Dataset
They are ok except. The fitted pops off.	Amazon
I tried a few LED harnesses and none were bright enough to see my black dog at night running through the woods. This vest, as long as not directly in front of head/tail, is super visable.	Amazon
Very short on the sides. Overall, good fit but I do not like to show my belly. Too bad lad got that. Fabric very soft.	Amazon
The script kicks in, and mr. hartley's distended pace and foot-dragging rhythms follow.	Pang Movie
Eastwood winces, clutches his chest and gasps for breath. it's a spectacular performance - ahem, we hope it's only acting.	Pang Movie
Parts seem like they were lifted from terry gilliam's subconscious , pressed through kafka's meat grinder and into buñuel's casings	Pang Movie
The recording session is the only part of the film that is enlightening and how appreciative you are of this depends on your level of fandom.	Vader Movie
It shows that some studios firmly believe that people have lost the ability to think and will forgive any shoddy product as long as there's a little girl on girl action.	Vader Movie
A light, engaging comedy that fumbles away almost all of its accumulated enjoyment with a crucial third act miscalculation.	Vader Movie

Table 5: Texts illustrating the “No consensus” category

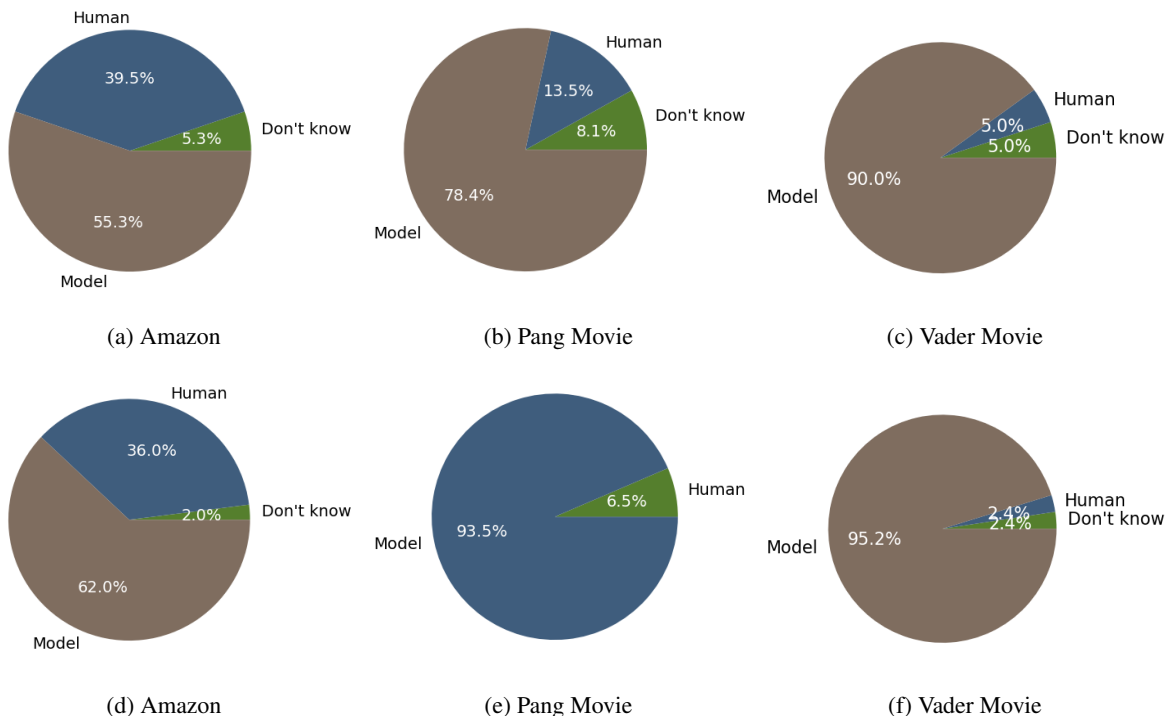


Figure 11: Percentages for answers to Question 1 by evaluators with a Computer Science background (a, b and c) and a Linguistics background (d, e and f).