# Modeling Label Semantics for Predicting Emotional Reactions

**Radhika Gaonkar, Heeyoung Kwon, Mohaddeseh Bastan, Niranjan Balasubramanian**
Stony Brook University, Stony Brook, New York
`{rgaonkar, heekwon, mbastan, niranjan}@cs.stonybrook.edu`

**Nathanael Chambers**
US Naval Academy, Annapolis, MD
`nchamber@usna.edu`

## Abstract

Predicting how events induce emotions in the characters of a story is typically seen as a standard multi-label classification task, which usually treats labels as anonymous classes to predict. They ignore information that may be conveyed by the emotion labels themselves. We propose that the semantics of emotion labels can guide a model's attention when representing the input story. Further, we observe that the emotions evoked by an event are often related: an event that evokes joy is unlikely to also evoke sadness. In this work, we explicitly model label classes via label embeddings, and add mechanisms that track label-label correlations both during training and inference. We also introduce a new semi-supervision strategy that regularizes for the correlations on unlabeled data. Our empirical evaluations show that modeling label semantics yields consistent benefits, and we advance the state-of-the-art on an emotion inference task.

## 1 Introduction

Understanding how events in a story affect the characters involved is an integral part of narrative understanding. Rashkin et al. (2018) introduced an emotion inference task on a subset of the ROCStories dataset (Mostafazadeh et al., 2016), labeling entities with the emotions they experience from the short story contexts. Previous work on this and related tasks typically frame them as multi-label classification problems. The standard approach uses an encoder that produces a representation of the target event along with the surrounding story events, and then pushes it through a classification layer to predict the possible emotion labels (Rashkin et al., 2018; Wang et al., 2018).

This classification framework ignores the semantics of the emotions themselves. Each emotion label (e.g., joy) is just a binary prediction. However, consider the sentence, "Danielle was really short on money". The emotional reaction is FEAR of being short on money. First, if a model had lexical foreknowledge of "fear", we should expect an improved ability to decide if a target event evokes FEAR. Second, such a model might represent relationships between the emotions themselves. For example, an event that evokes FEAR is likely to evoke SADNESS and unlikely to evoke JOY. When previous models frame this as binary label prediction, they miss out on ways to leverage label semantics.

In this work, we show that explicitly modeling label semantics improves emotion inference. We describe three main contributions[1]. First, we show how to use embeddings as the label semantics representation. We then propose a label attention network that produces label-informed representations of the event and the story context to improve prediction accuracy. Second, we add mechanisms that can make use of label-label correlations as part of both training and inference. During training, the correlations are used to add a regularization loss. During inference, the prediction logits for each label are modified to incorporate the correlations, thus allowing the model's confidence on one label to influence its prediction of other labels. Third, we show that the label correlations can be used as a semi-supervised signal on the unlabeled portion of the ROCStories dataset.

Our empirical evaluations show that adding label semantics consistently improves prediction accuracy, and produces labelings that are more consistent than models without label semantics. Our best model outperforms previously reported results and achieves more than 4.9 points absolute improvement over the BERT classification model yielding a new state-of-the-art result for this task.

---

[1]https://github.com/StonyBrookNLP/emotion-label-semantics

## 2 Emotion Inference

The emotion inference task introduced by Rashkin et al. (2018) is defined over a subset of short stories from the ROCStories dataset (Mostafazadeh et al., 2016). It infers the reactions that each event evokes in the characters of the story, given the story context thus far. For each sentence (i.e. event) in a story, the training data includes annotations of eight emotions. Given a sentence $x_s$ denoting a single event in a story, the task is to label the possible emotional reactions that an event evokes in each character in the story. Since an event can evoke multiple reactions, the task is formulated as a multi-label classification problem.

The standard approach to this task has been as follows. For a given character $c$ and the target sentence $x_s$, collect all previous sentences $x_c$ in the story in which the character $c$ is mentioned as the character context. Encode the target sentence, and the character context to obtain a single representation, and use it as input to a multi-label classification layer for prediction. Rashkin et al. (2018) benchmark the performance of multiple encoders (see Section 5).

We extend this previous work to integrate label semantics into the model by adding label embeddings (Section 3) and explicitly representing label-label correlations (Section 4).

## 3 Label Semantics using Embeddings

A simple strategy to model label semantics is to explicitly represent each with an embedding that captures the surface semantics of its label name. Since the emotion labels correspond to actual words (e.g., joy, fear, etc.), we can initialize them with their corresponding word embeddings (learned from a large corpus). We then use these label embeddings in two ways as detailed below.

### 3.1 Label Attention Network

The label embeddings can be used to guide an encoder network to extract emotion-related information from the sentences. We adopted the Label-Embedding Attentive Network (LEAM) architecture to produce label-focused representations (Wang et al., 2018). The main idea behind the LEAM model is to compute attention scores between the label and the representations of the tokens in the input that is to be classified[2]. This can

---

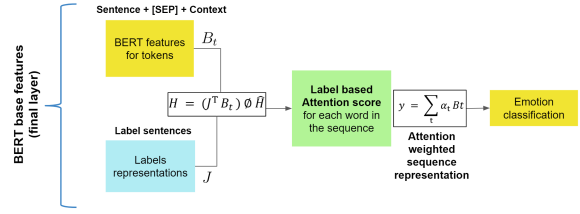[2]The original model used LEAM directly on top of Glove embeddings (Wang et al., 2018).



Figure 1: Label-Embedding Attentive Network using BERT Features. $y$ denotes the label attended story sentence and context representation, where $\alpha$ is the attention score.

then be used to appropriately weight the contributions of each token to the final representations. In this work, we use LEAM to compute an attention matrix computed over the hidden states produced by the encoder and the label embeddings. The encoder used is the BERT features for each token $B_t$ in the text and each of the label sentences $J$. The attention matrix is then used to produce a weighted combination of the contextual representations of the input, using the compatibility matrix $H$, as computed in (Wang et al., 2018). This gives emotion focused representations $y$ to use for classification:

$$H = (J^T B_t) \oslash \hat{H} \qquad (1)$$

Figure 1 illustrates the key steps in the model.

### 3.2 Labels as Additional Input

Rather than learning label embeddings from scratch, we also explore using contextual embeddings from transformer-based models like BERT. This allows us to use richer semantics derived from pre-training and also allows us to exploit the self-attention mechanism to introduce label semantics as part of the input itself. In addition to the target and context sentences, we also include emotion-label sentences, $L_s$, of the form "`[character] is [emotional state]`" as input to the classifier. For each instance, we add eight such sentences covering all emotional labels[3]. In this paper, we use the final layer of a pretrained Bert-base model to get representations for the input sentence and each of the emotion-label sentences. The self-attention mechanism will automatically learn to attend to these label sentences when constructing the representations for the input text.

---

[3]This is similar to how answer options are encoded in multiple choice question answering in transformer-based models.
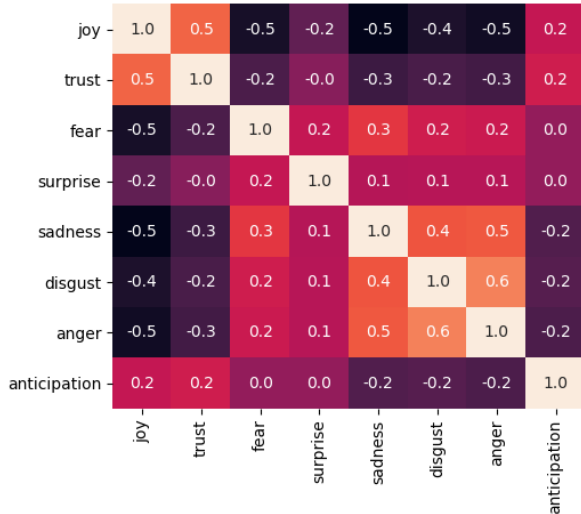
Figure 2: Emotion correlations as seen in the ground truth labels in the training data.

## 4 Label Semantics using Correlations

When more than one emotion is evoked by an event, they aren't independent. Indeed, as shown in Figure 2, there are strong (positive and negative) correlations between the emotion labels in the ground truth. For instance, there is a negative correlation ($\rho = -0.5$) between JOY and SAD labels and a positive correlation between JOY and TRUST ($\rho = 0.5$). We propose two ways to incorporate these label correlations to improve prediction.

### 4.1 Correlations on Labeled Data

In a multi-label setting, a good model should respect the label correlations. If it is confident about a particular label, then it should also be confident about other positively correlated labels, and conversely less confident about labels that are negatively correlated.

Following Zhao et al. (2019), we add (i) a loss function that penalizes the model for making incongruous predictions, i.e. those that are not compatible with the label correlations, and (ii) a component that multiplies the classification logit vector $z$ with the learned label relations encoded as a learned correlation matrix $G$. This component transforms the raw prediction score of each label to a weighted sum of the prediction scores of the other labels. For each label, these weights are given by its learned correlation with all the other labels. Therefore, the prediction score of each label is affected by the prediction score of the other labels, based on the correlation between label pairs. The final prediction scores are then calculated as shown in the equation:

$$e = \sigma(z \cdot G) \qquad (2)$$

The overall loss then comprises of two loss functions - the prediction loss ($\mathcal{L}_{BCE}$), and the correlation-loss ($\mathcal{L}_{corr}$):

$$\mathcal{L}(\theta) = \mathcal{L}_{BCE}(e, y) + \mathcal{L}_{corr}(e, y') \qquad (3)$$

Where $\mathcal{L}_{corr}$ computes BCE Loss with continuous representation of the true labels $y$, using the learned label correlation $G$:

$$y' = y \cdot G \qquad (4)$$

### 4.2 Semi-supervision on Unlabeled Data

We also introduce a new semi-supervision idea to exploit label correlations as a regularization signal on unlabeled data. The multi-label annotations used in this work (Rashkin et al., 2018) only comprises a small fraction of the original ROCStories data. There are ~40k character-line pairs that have open text descriptions of emotional reactions, but these aren't annotated with multi-label emotions, and therefore were not used in the above supervised emotion prediction tasks. We propose a new semi-supervised method over BERT representations that augments the soft-training objective used in Section 4.1 with a label correlation incompatibility loss defined over the unlabeled portion of the ROCStories dataset.

We use two loss functions: the loss computed in Equation 3, and the regularization loss on the unlabeled training data (Equation 5). For the semi-supervised training, we use an iterative batch-wise training. In the first step, all weights of the model are minimized by minimizing the loss in Equation 3. In the next step, the learned label correlations are updated using:

$$\mathcal{L}_{reg} = \sum_{i,j} G_{ij} \cdot d(e_i, e_j) \qquad (5)$$

$$d(e_i, e_j) = \begin{cases} \|e_i - e_j\| & \text{for } G_{ij} \geq 0, \\ \|e_i - e_j\| - 1 & \text{otherwise.} \end{cases}$$

This loss helps the model to produce consistent predictions based on the correlations by forcing positively correlated labels to have similar scores and negatively correlated ones to have dissimilar scores.

| | Model | Precision | Recall | F1 |
|---|---|---|---|---|
| Baselines | Rashkin et al. (2018) | | | |
| | BiLSTM | 25.31 | 33.44 | 28.81 |
| | CNN | 24.47 | 38.87 | 30.04 |
| | REN | 25.30 | 37.30 | 30.15 |
| | NPN | 24.33 | 40.10 | 30.29 |
| | Paul and Frank (2019)[*] | 59.66 | 51.33 | 55.18 |
| | BERT | 65.63 | 56.91 | 60.96 |
| Adding Label Semantics | **Label Embeddings** | | | |
| | LEAM w/ GloVe | 59.81 | 54.46 | 57.03 |
| | LEAM w/ BERT Features | 67.29 | 54.48 | 60.22 |
| | BERT + Labels as Input | 63.05 | 61.70 | 62.36 |
| | **Label Correlation** | | | |
| | Learned Correlations | 56.50 | 71.47 | 63.11 |
| | Semi-supervision | 57.94 | 76.35 | 65.88 |

Table 1: Comparison Results on ROCStories with Plutchik emotion labels

## 5 Experimental Setup

We compare our proposed models with the models presented in Rashkin et al. (2018), the LEAM architecture of Wang et al. (2018), and fine-tuned BERT models (Devlin et al., 2019) for multi-label classification without label semantics. For all the models we report the micro-averaged Precision, Recall and F1 score of the emotion prediction task.

Rashkin et al. (2018) modeled character context and pre-trained on free response data to predict the mental states of characters using different encoder-decoder setups, including BiLSTMs, CNNs, the recurrent entity network (REN) (Henaff et al., 2016), and neural process networks (NPN) (Bosselut et al., 2017). Additionally, we compare with the self-attention architecture proposed in (Paul and Frank, 2019), without the knowledge from ConceptNet (Speer and Havasi, 2012) and ELMo embeddings (Peters et al., 2018).

To compare against LEAM, we compare it against our proposal of the LEAM+BERT model, where our label attention is computed from BERT representations of each of the label sentences, and words in the input sentence. We also encode the sentence and context separately in a BiLSTM layer as done in Rashkin et al. (2018).

We also fine-tuned a BERT-base-uncased model for emotion classification, using $x_s$, $x_c$ and $L_s$ as inputs. This beats the other baselines by a significant margin, and is thus a strong new baseline. All our models are evaluated on the emotion reaction prediction task over the eight emotion labels (Plutchik categories) annotated in the Rashkin et al. (2018) dataset. We follow their evaluation setup, and report the final results on the test set. We use pre-trained GloVe embeddings (100d) and BERT-base-uncased representations with the LEAM model. The final classifier used in all models is a feed-forward layer, followed by a sigmoid.

## 6 Results

Table 1 compares the performance of the baselines with our models that use label semantics. Among the baselines, the fine-tuned BERT base model obtains the best results. Adding label embeddings (section 3.1) to the basic BiLSTM via LEAM model provides substantial increase, more than 27 absolute points in F1. We swapped in BERT features instead of GloVe and found a further 3 point improvement. The BERT baseline beat both of these, but appending label sentences as additional input to fine-tuned BERT increased its performance by 1.4 F1 points.

A further increase of 2 points in F1 is achieved by tracking label-label correlations through training loss and inference logits. In addition, adding semi-supervision yields the best gain of more than 4.9 points in F1 over basic BERT, providing a significant advance in state-of-the-art results for emotion inference in this dataset. We also checked the statistical significance of the Semi-supervision model (Table 1) against the Learned Correlations, BERT+Labels as Input, LEAM w/ BERT Features and the BERT model using the Randomization Test (Smucker et al., 2007). This involved comparing the outputs of the Semi-supervision model with the

| Sentence | Ground Truth | LS | NoLS |
|---|---|---|---|
| And nobody could give him any direction | Sad, Disgust Surprise | Sad, Disgust Anger | Sad |
| She said Mark can come for free | Joy, Trust Anticipation | Joy ,Trust Anticipation | Joy, Anticipation |
| He is relieved that it was not harmed | Joy, Surprise Anticipation | Joy, Surprise Anticipation | Fear, Surprise |
| The marshmallows were totally smooshed | Anger, Sad | Anger, Sad | Joy, Anticipation |

Table 2: Prediction of labels with label semantics (LS) versus without label semantics (NoLS). Including label semantics helps the model predict semantically labels (high correlations), with high probability.

above mentioned models after creating 100,000 random permutations. The Semi-supervision model achieved statistically significant improvement over all the baselines.

We did further qualitative analysis of the results on the dev set to better understand the performance of the Semi-supervised Label Semantics model. Compared to base BERT, this model predicts more emotion classes per instance (8839 vs 5024). The wrong predictions of this model have lower probabilities than the correct labels suggesting that classification could be further improved with proper threshold identification. This model is also better at capturing the semantic relations between labels during prediction. This is highlighted through some examples in Table 2.

## 7 Related Work

One of the most widely-used work in narrative understanding introduced ROCStories, a dataset for evaluating story understanding (Mostafazadeh et al., 2016). On a subset of these stories (Rashkin et al., 2018) added annotations for causal links between events in stories and mental states of characters. They model entity state to predict emotional reactions and motivations for causing events occurring in ROCStories. Additionally, they also introduce a new dataset annotation that tracks emotional reactions and motivations of characters in stories. Other work looked at encoding external knowledge sources to augment motivation inference (Paul and Frank, 2019) on the same dataset. Both treat labels as anonymous classes, whereas this work explores modeling the semantics of the emotion labels explicitly. Recent work in multi-label emotion classification has shown that using the relation information between labels can improve performance. (Kurata et al., 2016) use the label co-occurrence

information in the final layer of the neural network to improve multi-label classification. Correlation-based label representations have also been used for music classification styles (Zhao et al., 2019). Our work builds on these and adds a similar result showing that label correlations can have significant impact for emotion label inference.

## 8 Conclusions

We present new results for the multi-label emotion classification task of Rashkin et al. (2018), extending previous reported results by 10.7 F1 points (55.1 to 65.8). The multi-label nature of emotion prediction lends itself naturally to use the correlations between the labels themselves. Further, we showed that modeling the class labels as semantic embeddings helped to learn better representations with more meaningful predictions. As with many tasks, BERT provided additional context, but our integration of these label semantics showed significant improvements. We believe these models can improve many other NLP tasks where the class labels carry inherent semantic meaning in their names.

## Acknowledgments

# References

Antoine Bosselut, Omer Levy, Ari Holtzman, Corin Ennis, Dieter Fox, and Yejin Choi. 2017. Simulating action dynamics with neural process networks. *arXiv preprint arXiv:1711.05313*.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.

Mikael Henaff, Jason Weston, Arthur Szlam, Antoine Bordes, and Yann LeCun. 2016. Tracking the world state with recurrent entity networks. *arXiv preprint arXiv:1612.03969*.

Gakuto Kurata, Bing Xiang, and Bowen Zhou. 2016. Improved neural network-based multi-label classification with better initialization leveraging label co-occurrence. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 521–526.

Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong He, Devi Parikh, Dhruv Batra, Lucy Vanderwende, Pushmeet Kohli, and James Allen. 2016. A corpus and cloze evaluation for deeper understanding of commonsense stories. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 839–849, San Diego, California. Association for Computational Linguistics.

Debjit Paul and Anette Frank. 2019. Ranking and Selecting Multi-Hop Knowledge Paths to Better Predict Human Needs. In *Proceedings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics*, volume 1, pages 3671–3681, Minneapolis, Minnesota, USA.

Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.

Hannah Rashkin, Antoine Bosselut, Maarten Sap, Kevin Knight, and Yejin Choi. 2018. Modeling naive psychology of characters in simple commonsense stories. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2289–2299, Melbourne, Australia. Association for Computational Linguistics.

Mark D Smucker, James Allan, and Ben Carterette. 2007. A comparison of statistical significance tests for information retrieval evaluation. In *Proceedings of the sixteenth ACM conference on Conference on information and knowledge management*, pages 623–632.

Robyn Speer and Catherine Havasi. 2012. Representing general relational knowledge in Concept-Net 5. In *Proceedings of the Eighth International Conference on Language Resources and Evaluation (LREC'12)*, pages 3679–3686, Istanbul, Turkey. European Language Resources Association (ELRA).

Guoyin Wang, Chunyuan Li, Wenlin Wang, Yizhe Zhang, Dinghan Shen, Xinyuan Zhang, Ricardo Henao, and Lawrence Carin. 2018. Joint embedding of words and labels for text classification. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2321–2331, Melbourne, Australia. Association for Computational Linguistics.

Guangxiang Zhao, Jingjing Xu, Qi Zeng, Xuancheng Ren, and Xu Sun. 2019. Review-driven multi-label music style classification by exploiting style correlations. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2884–2891, Minneapolis, Minnesota. Association for Computational Linguistics.