
A Hierarchical Recurrent Encoder-Decoder Model for Supervised Analysis of Therapeutic Alliance in Online Therapy Conversations

Emily Tseng

Abstract

As online therapy delivered via asynchronous text-chat grows more and more popular, there emerges a need for language analysis systems that can automatically infer high-level attributes of therapy conversations. Rising tension, lessening interest, and other such attributes are easily discernable to the therapist and to any human observer, but the asynchronous nature of this type of therapy means human observers are not always on hand to make such judgments. If used wisely, automated systems might assist therapists and their supervisors in preemptively flagging conversations with the patients at highest risk for abandoning treatment, and ultimately enable them to provide better care—but to achieve such systems, we need models that can capture emerging attributes of conversations as they progress.

In this work, we provide a proof of concept for a conversational forecasting model that makes use of generative pre-training to develop a neural representation of high-level conversational dynamics. Our model uses this representation to infer the patient’s score for the working alliance between patient and therapist at a given point in a therapeutic conversation. Using a proprietary dataset of transcripts and patient outcomes from a major provider of online therapy, we show that this model outperforms baselines that do not use generative pretraining. Our work extends prior literature on conversational forecasting frameworks from social media exchanges to therapy conversations, which by nature have more complex long-term dependencies. We close by identifying several opportunities for next steps.

1. Introduction

In therapeutic and caregiving contexts, a patient’s felt sense of *alliance* with a care provider can make or break their treatment. This is particularly important in remote caregiving contexts like online therapy, where the vast majority of

the connection between the patient and the caregiver takes place entirely over text-chat. As the popularity of this form of therapy grows, the therapists and supervisors who work on these platforms become akin to call-center employees: shuffling quickly between communication channels with an increasing and neverending caseload of patients demanding their immediate attention. Equipping these experts with systems that analyze therapeutic conversations and forecast relative risk might enable them to more effectively direct their attention and resources, and ultimately provide better care to more people.

Such systems, however, are not straightforward to build. Our task belongs to the broader field of *conversational forecasting*, which has identified two core modeling challenges: (1) forecasting conversational dynamics requires capturing *high-level attributes* about a conversation that may not be discernible from individual utterances; and (2) because conversations have an *unknown horizon*, they require online processing methods and cannot rely on fixed-length inputs. Many combinations of methods have been applied to similar conversational forecasting tasks, ranging from handcrafted features developed from psycholinguistic intuition (Zhang et al., 2018) to risk scores computed via similarity metrics between utterances (Althoff et al., 2016). Most recently, Chang et al. 2019 showed that a model using generative pretraining to develop neural representations of high-level dynamics could then be fine-tuned to outperform these baselines on a supervised prediction task.

In this project, we provide a proof of concept for the use of neural representations of high-level conversational dynamics in forecasting alliance between patients and therapists. We describe a model for predicting alliance at any given point in a conversation, and examine its potential using a proprietary dataset of therapy transcripts and patient-provided alliance scores from a major provider of online therapy services. We compare the performance of our model against a traditional feature-engineered method, bag-of-words, and examine variants with and without attention. Our work provides a starting point for further work developing conversational forecasting methods suited to this important problem domain. To summarize, in this work we:

- develop a new model extending the current state-of-the-art in conversational forecasting to meet the challenges of online therapy conversations,
- show our model outperforms approaches that do not make use of generative pre-training prior to the supervised task, and
- lay a foundation for upcoming work examining more novel approaches using a proprietary real-world dataset of text-chat therapy transcripts.

2. Background

Online therapy data. Much of the current state-of-the-art in conversational forecasting has made use of datasets drawn from online social media conversations, for example publicly available Reddit threads (Zhang et al., 2018) or Wikipedia conversations (Chang & Danescu-Niculescu-Mizil, 2019). The language of online therapy is different, resembling something more akin to private SMS conversations.

Two key attributes of this type of data make it a fit for the neural approach to conversational forecasting. First, therapeutic conversations by definition involve therapists deploying high-level conversational strategies not visible in individual utterances: therapists are trained to engage with patients in a way that nudges them towards self-driven behavior change. These conversational strategies manifest differently in different therapists, according to their trainings, personal conversational styles, and levels of comfort with the patient. Thus, detection of these high-level attributes cannot be distilled to keyword-based searches or other handcrafted features.

Second, therapeutic conversations exhibit additional levels of variation in their progression over time. Unlike therapeutic conversations in hotline contexts (Althoff et al., 2016) or in traditional, in-person therapy, there are no enforced time bounds in online therapy: a patient and a therapist converse asynchronously in a chatroom that is always available to both parties. Our data additionally show variation in *how* people use these online spaces: some patients send many short messages to their therapists within compressed timeframes ('bursts'), while others send longer messages at longer intervals; therapists similarly exhibit variable messaging behaviors. (For more dataset statistics, refer to section 4.) Thus, there is no way to systematize a predetermined breakpoint at which it may be appropriate to make a forecast, and online processing methods are required for the task.

Therapeutic alliance and WAI. The psychology literature has had a longstanding interest in the role of the relationship between the therapist and the patient in successful psychotherapy. Early works posited that a beneficial

attachment between patient and therapist enables the patient to trust the therapist enough to use his or her interpretations to bring about positive change (for a review, see Ardito & Rabellino 2011). Bordin 1979 defined the concept of a *working alliance* as the relational underpinning of many approaches to psychotherapy, and outlined that different types of alliance emerge from different approaches. Regardless of approach, Bordin argues, it is the *strength* of the alliance that best predicts positive outcomes.

As a way to measure alliance in a given patient-therapist relationship, Horvath & Greenberg 1989 proposed the Working Alliance Inventory (WAI), a set of self-reported scales that measure the quality of the alliance along the three dimensions defined in Bordin's theoretical framework: (1) the *bond* between patient and therapist, (2) the agreement on the *goals* of the therapy, and (3) the agreement on the *tasks* required to achieve those goals. Available in variants from the perspective of the patient, the therapist, and a third-party observer, the WAI has been validated as a reliable metric for alliance in several contexts (for a meta-analysis, see Martin et al. 2000).

3. Related Work

Our work adapts and extends the Conversational Recurrent Architecture for Forecasting (CRAFT), a framework integrating generative pre-training with a supervised fine-tuning model to achieve improved predictive ability on conversation-level attributes, e.g., whether an online conversation will derail into personal attacks (Chang & Danescu-Niculescu-Mizil, 2019). CRAFT itself is an improvement on previous models applied to conversational forecasting that extracted hand-crafted features from fixed-length sliding windows within a conversation (Zhang et al., 2018). By making use of generative pre-training, CRAFT learns a neural representation of the high-level attributes of a conversation, and uses that representation to make downstream predictions. Whereas CRAFT focused on prediction of a binary outcome (derailment vs. non-derailment), our work focuses on classification of progressive segments of a conversation into one of multiple outcomes (buckets of WAI scores, as described in section 4). Thus we describe the models reported in this work as variants of CRAFT-Multilabel.

We additionally extend CRAFT to consider a data type with significantly different modeling considerations. CRAFT was developed to predict the risk of derailment in social media conversations, including threads on Reddit and Wikipedia. In contrast, we study therapeutic conversations, which consist of longer utterances and contain more long-range dependencies, thus presenting a more challenging modeling task. Prior work on similar data most notably includes Althoff et al. 2016, which attempted to quantify

higher-level conversational attributes via vector similarities between utterances using a dataset of transcripts from an SMS-based crisis hotline. A related body of prior work has applied existing measures in discourse analysis, for example synchrony, to therapeutic conversations, showing that effective therapists tend to stylistically mimic their patients (Doré & Morris, 2018).

4. Dataset

We consider a dataset of therapy transcripts and associated patient outcomes from Talkspace, an online therapy platform. Due to the highly sensitive nature of therapy, we put significant effort into respecting patients’ privacy and autonomy. All represented patients gave informed consent for the use of their data in research, and transcripts were anonymized by Talkspace before they were handed to our research team. Study data was stored on a secure remote server maintained by research IT staff at our institution and accessible only when connected to our institutional VPN; at no point was study data downloaded to another local or cloud environment. Our study protocol was approved by our institutional IRB.

In total, our dataset consists of 5.7M messages exchanged between patients and therapists, representing 11,233 patients’ full courses of treatment. 1,906 therapists are represented, with an average of 9 patients per therapist. As discussed in section 2, the language of online therapy is highly variable: utterances in our dataset averaged 78.0 words (median=32, std=139.66, range=0-3363).

Outcome annotations are provided in the form of patients’ responses to surveys issued approximately every 3 weeks. Patients vary in the number of utterances they provide before each survey, and in the number of surveys completed overall. In total, 13,742 WAI scores were provided by 6,702 patients. Patients provided an average of 2.1 WAI scores (range 1-24, stdev 2.2). As depicted in Figure 1, The overall distribution of scores skewed strongly towards the positive end of the spectrum (more strongly allied).

Given this, in this project we formulated our predictive task as multiclass classification between three buckets of scores: 0-14, 15-19, and 20. These can be interpreted, respectively, the patient perceiving a less-than-ideal alliance with their therapist, the patient perceiving a good alliance with their therapist, and the patient perceiving the strongest possible alliance with their therapist. We confirmed this interpretation with our domain-expert partner at Talkspace.

5. Model

Our model is an adaptation of the Conversational Recurrent Architecture for Forecasting (CRAFT) (Chang &

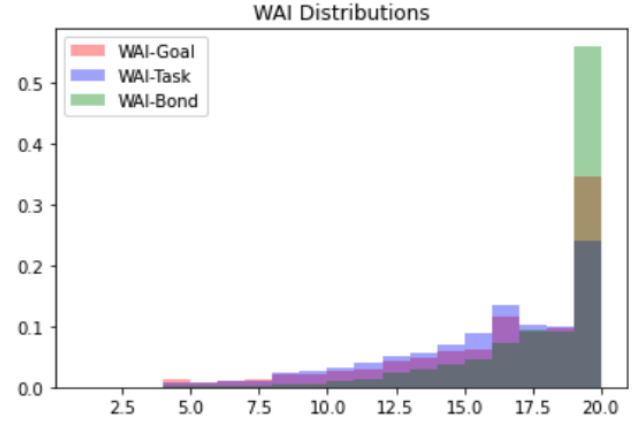


Figure 1. Distribution of WAI scores across subdimensions.

Danescu-Niculescu-Mizil, 2019), which integrates a generative dialogue model and a supervised fine-tuning component to produce predictions utterance-by-utterance about some high-level conversational state.

Problem definition. We define a conversation C as a variable-length sequence of n utterances, $C = \{u_1, \dots, u_n\}$. Utterances are variable-length sequences of tokens w , and thus $u_n = \{w_1, \dots, w_{M_n}\}$, where M_n is the length in tokens of utterance n .

Given a therapy exchange $C = \{u_1, \dots, u_n\}$, we generate h_n^{con} , a neural representation of high-level conversational state up to utterance u_n . We then use h_n^{con} as the input to a classifier that predicts y_n , the label attached to the conversation up to utterance n . We define the label y_n based on the WAI score provided by the patient at utterance u_n : 0 if the score is in the bucket 0-14, 1 if 15-19, and 2 if 20.

Generative component. Following Chang et al. 2019, we adopted for our generative component the hierarchical recurrent encoder-decoder (HRED) architecture proposed in Sordoni et al. 2015 and Serban et al. 2016. Built to model high-level conversational context, including temporal structure and dependencies between consecutive sequential inputs, HREDs are uniquely suited for conversational forecasting tasks.

HREDs are comprised of three component recurrent neural networks (RNNs): an utterance encoder, a conversation encoder, and a decoder. First, the *utterance encoder* generates for each utterance a semantic vector representation via its hidden state $h_m^{enc} \in \mathbb{R}_{enc}^d$, where d_{enc} is the desired dimension. For each token w_m in utterance n of length M , the encoder updates its h_m^{enc} like so:

$$h_m^{enc} \leftarrow f_{enc}^{RNN}(w_m, h_{m-1}^{enc}) \quad (1)$$

The utterance encoder’s hidden state at the last step, h_M^{enc} , in theory represents an embedding for the entire utterance.

Following Serban et al. 2016, h_0^{enc} is initialized as the zero vector $\mathbf{0}$, and following Chang et al. 2019, we use the GRU (Cho et al., 2014) as our nonlinear gating function f^{RNN} .

Next, the *conversation encoder* uses the hidden states from each consecutive comment in a sequence of length N to produce an embedding h_n^{con} for the conversation up to the utterance at that point (u_N):

$$h_n^{con} \leftarrow f_{con}^{\text{RNN}}(h_{M_n}^{enc}, h_{n-1}^{con}) \quad (2)$$

The conversation encoder also initializes its hidden state h_0^{con} with the zero vector $\mathbf{0}$, and also uses the GRU as its nonlinearity. We denote the dimension of h_n^{con} as d_{con} .

The *decoder* uses the embedded conversational context h_n^{con} to generate a response to utterance n . Following Sordani et al. 2015, it does this by first initializing its own hidden state $h^{dec} \in \mathbb{R}^{d_{dec}}$ using a nonlinear activation of h_n^{con} :

$$h_0^{dec} = \tanh(Dh_n^{con} + b_0) \quad (3)$$

Where $D \in \mathbb{R}^{d_{dec} \times d_{con}}$ projects the context embedding into decoder space, and $b_0 \in \mathbb{R}^{d_{dec}}$. The decoder then updates its own hidden state for each response token using the following recurrence:

$$h_t^{dec} \leftarrow f_{dec}^{\text{RNN}}(w_{t-1}, h_{t-1}^{dec}) \quad (4)$$

The decoder then produces the next token in its response by producing a probability distribution over words from h_t^{dec} :

$$w_t = f^{out}(h_t^{dec}, w_{t-1}) \quad (5)$$

In our implementation, following CRAFT, we supplement f^{out} with attention (Luong et al., 2015), the intuition being that certain dimensions of the context encoder states c_t – output from f_{con}^{RNN} alongside h_n^{con} – may be more informative for decoding than others. Our final f^{out} is as follows:

$$f^{out} = \text{softmax}(\tanh(W_c[h_t^{dec}|c_t])) \quad (6)$$

Per Luong et al. 2015, we utilize the *concat*-based scoring function in our global attention to compute c_t .

Predictive component. Our predictive component uses the conversational embedding up to utterance u_n to generate a prediction for the WAI score at that utterance. We operationalize this as a multi-layer perceptron (MLP) that takes in the conversational state h_n^{con} and produces a distribution $p(Y_n|h_n^{con})$ over possible labels $Y = \{0, 1, 2\}$. Adapting from Chang et al. (2019), our MLP uses three fully-connected layers and leaky ReLU activations between each layer; but for our task of multilabel classification, we use a softmax activation. The result is a model that creates a probability distribution in which each score can be interpreted as the likelihood of the given label.

We further implemented two variants of the predictive component: classifiers with and without attention (Luong et al., 2015). The intuition here was that h_n^{con} may encode many different high-level conversational attributes, not just those that have bearing on our particular outcome variable. By learning to weight some dimensions of h_n^{con} above others, we may achieve better performance on targeted tasks.

Parameters. The parameters of our model include those used within each RNN of the generative component, $\theta_{enc}, \theta_{con}, \theta_{dec}$, as well as those used within the predictive component, θ_{clf} .

Each RNN within the generative component contains a set of parameters for its constituent GRU (see Sordani et al. 2015 for a longer explanation). θ_{dec} also includes parameters D and b_0 transforming h_n^{con} into h_0^{dec} , as well as all the parameters for *concat*-scored global attention per Luong et al. 2015.

As for the predictive component, across both with- and without-attention variants, θ_{clf} is operationalized as a three-layer feedforward MLP. The attention variant additionally contains a layer of attention weights.

6. Training

For the generative component, training involves maximizing the log-likelihood of the provided context-reply pairs. For a given conversation C comprised of N utterances:

$$L(\theta) = \sum_{n=1}^N \log(P(u_n|u_{1:n-1})) \quad (7)$$

$$= \sum_{n=1}^N \sum_{m=1}^{M_n} \log(P(w_{n,m}|w_{n,1:m-1}, u_{1:n-1})) \quad (8)$$

Optimization is done by applying the back-propagation through time (BPTT) algorithm standard to RNN training (Rumelhart et al., 1986).

Similarly, training for the predictive component back-propagates the cross-entropy loss between the model output and the label. Note that when training the predictive component, we follow the principles of *fine-tuning* and back-propagate through the entire model, all the way back to the encoder (Howard & Ruder, 2018). This nudges the embeddings learned during pretraining towards greater applicability to the predictive task.

7. Methods

Generative pre-training. We began by training our generative model following the structure outlined in section 6. From our dataset of 5.7M messages we randomly subsampled 250k pairs of contexts (sequences of utterances) and

replies. Randomization allowed our model to see a variety of contexts from a variety of conversations, as opposed to repeatedly seeing subsets of the same conversations. For example, in context-reply pairs $(\{u_1, \dots, u_{n-1}\}, u_n)$ and $(\{u_1, \dots, u_n\}, u_{n+1})$ from the same conversation, the subsequence of utterances $\{u_1, \dots, u_{n-1}\}$ repeat across both contexts.

Our generative model was implemented using Pytorch (Paszke et al., 2019), and optimized using that framework’s built-in Adam optimizer. Encoder and conversation encoder learning rates were set at 0.0001, and the decoder learning rate was set at 0.0005. Training took place on our remotely administered server per the bounds of our IRB.

Task-specific fine-tuning. All experiments with the predictive component of CRAFT-Multilabel first initialized the encoder, conversation encoder, and decoder parameters from the final iteration of the generative model. For this proof-of-concept, we subsampled 320 context-label pairs from our dataset and constrained analysis to one outcome dimension, WAI-Goal. As mentioned in 4, we adapted the data to bucket the provided raw outcome measures into three classes: 0-14, 15-19, and 20. Our core evaluation metric was overall accuracy.

Using a train-test split of 70-30, we trained for 10 epochs with a batch size of 16 and a constant learning rate of 0.001. All models used a hidden size of 100. All training took place on our remotely administered CPU per the bounds of our IRB. Training and evaluation runs on this machine took an average of 2.5 hours per run.

CBoW comparison. Core to the proposed applicability of the CRAFT-Multilabel framework is its ability to represent the high-level conversational dynamics that emerge from utterances in sequence (e.g., derailment, the level of working alliance present), and make predictions about those high-level conversational dynamics by using those representations as features in a downstream task.

To test whether our model was indeed representing higher-level semantics, we compared CRAFT-Multilabel against a cumulative bag-of-words (CBoW) model. Given a vocabulary V and an input context $C_n = \{u_1, \dots, u_n\}$, CBoW represents C_n as a $|V|$ -length vector of the counts of each word in V in the context ($f(C_n)$). In our implementation, we normalized the raw counts of each word in V according to their relative importance within their contexts, quantified by their term frequency-inverse document frequency (tf-idf). CBoW is inherently suited for online learning of progressive parts of the conversation: When faced with the next utterance in the sequence, $C_{n+1} = \{u_1, \dots, u_{n+1}\}$, CBoW recomputes $f(C_{n+1})$ to reflect the *cumulative* word counts across the updated context.

The $|V|$ -length feature vectors $f(C)$ are used as input to

a simple feedforward MLP that mimics the structure of the predictive component of CRAFT-Multilabel: that is, it produces a distribution $p(Y|f(C))$ over possible labels $Y = \{0, 1, 2\}$. Note that the CBoW by definition only captures relative word occurrences; it does not capture word order, or even boundaries between the utterances in our contexts, to say nothing of higher-level conversational dynamics. Thus, comparison of CBoW against CRAFT-Multilabel should demonstrate the effect of incorporating such conversation-level attributes into this predictive task.

Our CBoW was an MLP trained via scikit-learn’s Adam optimizer (Pedregosa et al., 2011) for 200 iterations with a constant learning rate of 0.001. Similar to the predictive component of CRAFT-Multilabel, the MLP consisted of 2 hidden layers of size 100, each of which used a ReLU activation. We trained and evaluated this model using the same train-test splits as the CRAFT-Multilabel predictive component described above.

8. Results

Table 1 depicts the performance of our models on our held-out test data. As shown, the CRAFT-Multilabel variant with attention (CRAFT-MA) outperforms both the CBoW baseline and the variant without attention (CRAFT-M).

To confirm that our models were indeed training, we examined losses over time for both the generative pre-training and the two CRAFT variants (Figure 2). Note that in this proof-of-concept we set constant learning rates and did not anneal. Our plots show that our generative model did learn and improve on loss values through the early part of training before reaching a noisy equilibrium, presumably because the learning rate was then too large. Both CRAFT variants also show they were able to learn over time, as reflected in their decreasing loss values.

9. Discussion

Examination of the confusion matrices from each testing run (Figure 3), as well as the precision, recall and F1 values (Table 1) provides a window into the model’s performance. First, all models appeared to be uniquely good at distinguishing which samples earned a score of 20, as reflected in the fact that all models’ overall highest F1 was for samples in the “perfect alliance” bucket (20), and the density of correctly classified samples in that bucket on the confusion matrices. This is a surprising result not explained by sampling bias: despite the skew in the overall dataset (Figure 1), we balanced samples from all three classes in training and testing. We had hypothesized that the cutoff between “good” alliance (15-19) and “perfect” alliance (20) would not be quite as clear-cut; evidently, the models are capable of learning the difference.

Model	Test-set Acc.	Class 0 (WAI-Goal 0-14)			Class 1 (WAI-Goal 15-19)			Class 2 (WAI-Goal 20)		
		Precision	Recall	F1	Precision	Recall	F1	Precision	Recall	F1
CBoW	0.5368	0.62	0.47	0.53	0.37	0.47	0.41	0.68	0.68	0.68
CRAFT-M	0.5579	0.55	0.42	0.48	0.41	0.55	0.47	0.72	0.65	0.68
CRAFT-MA	0.6105	0.77	0.38	0.51	0.5	0.59	0.54	0.65	0.78	0.7

Table 1. Test-set performance for the CBoW baseline and CRAFT-Multilabel variants. CRAFT-Multilabel with attention is denoted as CRAFT-MA, and without attention as CRAFT-M. All trials performed with a subset of 320 context/label pairs balanced between the 3 classes and split 70-30 between training and testing. For each model, accuracy figures reflect the top-performing version.

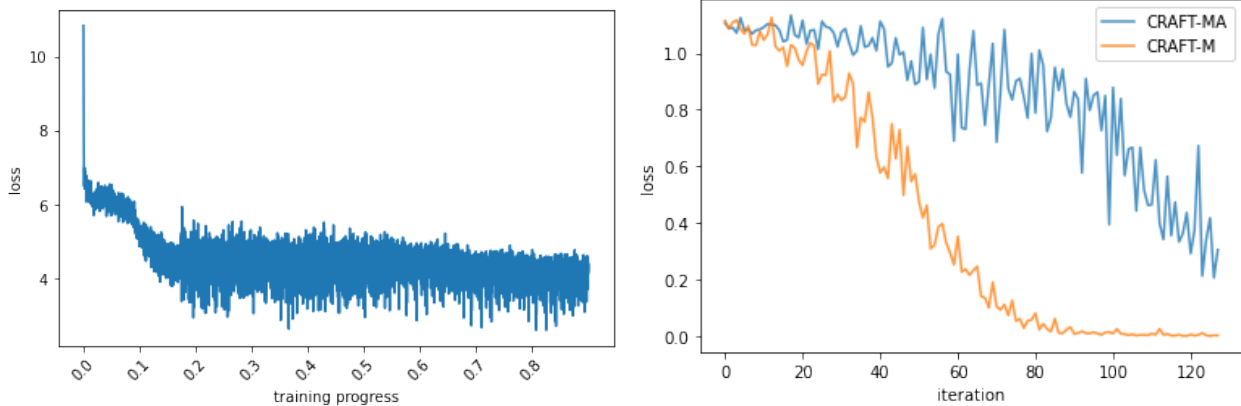


Figure 2. Losses over training iterations during generative pretraining (L) and finetuning of both CRAFT-M variants (R).

Our results beg further inquiry into what, precisely, the neural representations of higher-level conversational dynamics h^{con} are actually learning. We reserve an in-depth evaluation of those representations for next steps.

10. Conclusion

In this work, we implemented, trained and tested a proof of concept for a model that can infer working alliance from the texts of a talk-therapy conversation. For a given segment of conversation, our model utilizes conversation-level embeddings learned via generative pre-training to forecast alliance scores across three buckets: 0-14, 15-19, and 20. Tested on a held-out sample of 96 context/label pairs, our model significantly outperforms both a baseline bag-of-words approach and a variant without attention (Table 1).

This proof of concept opens many avenues for next steps. One obvious avenue is to expand analysis to a larger subset of our data. Now that this project has laid the groundwork for training and testing on this protected dataset, expansion to larger training sets that require longer training times is readily achievable. Another obvious avenue for further exploration is to incorporate additional information into the initial utterance encodings: namely, the timestamps associated with each. Such information might help the model learn whether response times and the ‘burstiness’ of mes-

saging patterns has bearing on a patient’s felt sense of alliance. In this work we also adapted the data to forecast working alliance across three buckets, 0-14, 15-19, and 20, in order to map to the distribution in our dataset. Future work might consider regressing to an exact value instead of framing the problem as multiclass classification.

Looking beyond optimizations to model to alternative approaches to the task, we are compelled by the possibility of using publicly available pretrained models for the English language such as BERT or GPT-2 to develop the conversation embeddings used as inputs to our predictive task. Trained on massive datasets, these models have recently shown success in a variety of natural language problems. Finally, we are also compelled by the possibility of variational approaches to the task, although explorations of a fundamentally different modeling framework may be best suited to a dataset with fewer constraints.

Lastly, there remains significant work in developing frameworks for interpreting neural representations of high-level conversational attributes, in order to better understand what components of the source conversations lead to scores in one direction or another.

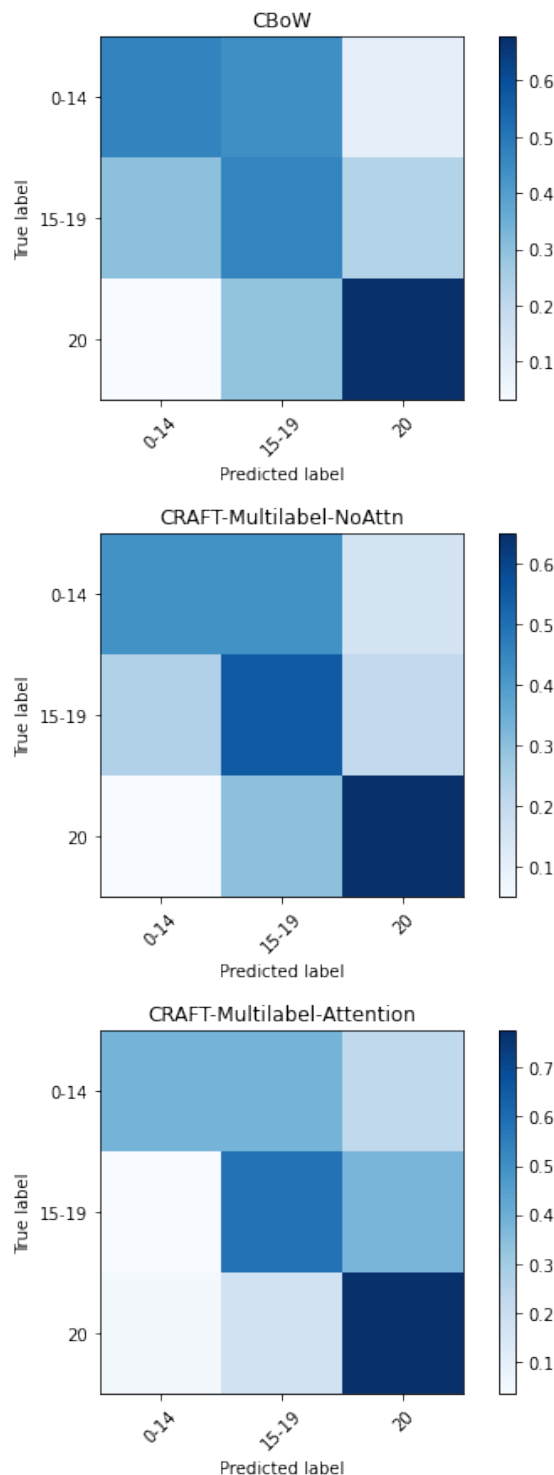


Figure 3. Confusion matrices for the CBoW and CRAFT variants.

References

- Althoff, Tim, Clark, Kevin, and Leskovec, Jure. Large-scale analysis of counseling conversations: An application of natural language processing to mental health. *Transactions of the Association for Computational Linguistics*, 4:463–476, 2016.
- Ardito, Rita B and Rabellino, Daniela. Therapeutic alliance and outcome of psychotherapy: historical excursus, measurements, and prospects for research. *Frontiers in psychology*, 2:270, 2011.
- Bordin, Edward S. The generalizability of the psychoanalytic concept of the working alliance. *Psychotherapy: Theory, research & practice*, 16(3):252, 1979.
- Chang, Jonathan P and Danescu-Niculescu-Mizil, Cristian. Trouble on the horizon: Forecasting the derailment of online conversations as they develop. In *Proceedings of EMNLP*, 2019.
- Cho, Kyunghyun, Van Merriënboer, Bart, Gulcehre, Caglar, Bahdanau, Dzmitry, Bougares, Fethi, Schwenk, Holger, and Bengio, Yoshua. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.
- Doré, Bruce P and Morris, Robert R. Linguistic synchrony predicts the immediate and lasting impact of text-based emotional support. *Psychological science*, 29(10):1716–1723, 2018.
- Horvath, Adam O and Greenberg, Leslie S. Development and validation of the working alliance inventory. *Journal of counseling psychology*, 36(2):223, 1989.
- Howard, Jeremy and Ruder, Sebastian. Universal language model fine-tuning for text classification. *arXiv preprint arXiv:1801.06146*, 2018.
- Luong, Minh-Thang, Pham, Hieu, and Manning, Christopher D. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015.
- Martin, Daniel J, Garske, John P, and Davis, M Katherine. Relation of the therapeutic alliance with outcome and other variables: a meta-analytic review. *Journal of consulting and clinical psychology*, 68(3):438, 2000.
- Paszke, Adam, Gross, Sam, Massa, Francisco, Lerer, Adam, Bradbury, James, Chanan, Gregory, Killeen, Trevor, Lin, Zeming, Gimelshein, Natalia, Antiga, Luca, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, pp. 8024–8035, 2019.

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

Rumelhart, David E, Hinton, Geoffrey E, and Williams, Ronald J. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.

Serban, Iulian V, Sordoni, Alessandro, Bengio, Yoshua, Courville, Aaron, and Pineau, Joelle. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.

Sordoni, Alessandro, Bengio, Yoshua, Vahabi, Hossein, Lioma, Christina, Grue Simonsen, Jakob, and Nie, Jian-Yun. A hierarchical recurrent encoder-decoder for generative context-aware query suggestion. In *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*, pp. 553–562, 2015.

Zhang, Justine, Chang, Jonathan P, Danescu-Niculescu-Mizil, Cristian, Dixon, Lucas, Hua, Yiqing, Thain, Nithum, and Taraborelli, Dario. Conversations gone awry: Detecting early signs of conversational failure. *arXiv preprint arXiv:1805.05345*, 2018.