

Improving Tweet Representations using Temporal and User Context

Ganesh J¹, Manish Gupta^{1,2}, and Vasudeva Varma¹

¹ IIIT, Hyderabad, India

ganesh.j@research.iiit.ac.in, vv@iiit.ac.in

² Microsoft, India

gmanish@microsoft.com

Abstract. In this work we propose a novel representation learning model which computes semantic representations for tweets accurately. Our model systematically exploits the chronologically adjacent tweets ('context') from users' Twitter timelines for this task. Further, we make our model user-aware so that it can do well in modeling the target tweet by exploiting the rich knowledge about the user such as the way the user writes the post and also summarizing the topics on which the user writes. We empirically demonstrate that the proposed models outperform the state-of-the-art models in predicting the user profile attributes like spouse, education and job by 19.66%, 2.27% and 2.22% respectively.

1 Introduction

The short and noisy nature of tweets poses challenges in computing accurate latent tweet representations. We observe that Paragraph2Vec [1] which is good in computing document representation overfits when evaluated for tweets, mainly due to the short length of tweets. To overcome this problem we utilize additional context from Twitter itself. Specifically, we hypothesize that a principled usage of chronologically adjacent tweets from users' Twitter timelines can help in significantly improving the quality of the representation. The main challenge lies in assigning appropriate attention weights to context tweets such that semantically relevant tweets receive high weights compared to less relevant ones. Consider Fig 1³, where we want to learn the representation for the tweet $t(j)$. One can see that the target tweet $t(j)$ has less semantic interactions with the context tweet $t(j-2)$. To capture this, we propose an attention based model that assigns a variable weight to each context tweet that captures the semantic correspondence between the target tweet and the context tweet. We further augment the attention model to be user-aware so that it can do well in modeling the target tweet by exploiting the rich knowledge about the user such as the way the user writes the post, and also summarizing the topics on which the user writes. Our work is closest to [2] where documents are modeled based on their word context as well as document stream context. We differ from their work in two ways: (1) they naïvely assume that all the documents in a stream have equal amount of semantic interactions and, (2) they ignore the knowledge of user (or document author).

We summarize our main contributions below. In summary, our contributions are as follows. (1) Our work is the first to model the semantics of the tweet using the temporal

³ The tweets are borrowed from Barack Obama's Twitter timeline posted in Sep 2015.

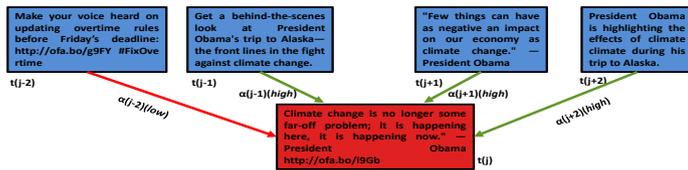


Fig. 1: $t(j-1)$, $t(j-2)$, $t(j+1)$ and $t(j+2)$ form the temporal context of $t(j)$. α 's denote the attention parameters of the proposed model.

context. (2) We introduce a novel attention based model that learns the weights for context tweets by back-propagating semantic loss. (3) We propose a novel way to learn user vector summarizing the content the user writes, which in turn helps in enriching the quality of the tweet embeddings. (4) We conduct quantitative analysis to showcase the application potential of the tweet representations learned from the model and also provide some interesting findings.

2 Related Work

Le et al. [1] adapt Word2Vec to learn document representations which are good in predicting the words present in the document. As seen in Section 5, for short documents like tweets, the model tends to learn poor document representations as the vector relies too much on the document content, resulting in overfitting. Djuric et al. [2] learn document representations using word context (same as [1]) along with document stream context in a hierarchical fashion. This work inspired us to learn tweet representations using user specific Twitter streams.

3 Problem Formulation

In this section we first introduce the notions of temporal context and attention, and then provide a formal problem statement.

Temporal context: Temporal context of a tweet $t(j)$ is the set of C_T tweets posted before and after $t(j)$ by the same user. The value C_T is a user specified parameter that defines the size of the temporal context to be considered to model a given tweet. For example, in Fig 1 we fix C_T as 2, the context tweets of $t(j)$ are $t(j-1)$, $t(j-2)$, $t(j+1)$ and $t(j+2)$.

Attention: An attention value is associated with a context tweet that defines the degree of semantic similarity between the context tweet and the target tweet. The more the latent semantic interactions between the tweets, the more is the attention. We denote the attention of context tweet $t(j-1)$ as $\alpha(j-1)$. For instance, in Fig 1, the attention value of context tweet $t(j-2)$ should be lower than that of context tweet $t(j-1)$ with respect to target tweet $t(j)$. In Fig 1, clearly $t(j-2)$ is not talking about the topic 'Climate Change' and so it makes sense to have a lower attention value.

Problem Statement: Let the training tweets be given in the order in which they are posted. In particular, we assume that we are given a user set U of N_u tweet sequences, with each sequence $u(k) \in U$, containing N_t tweets, $u(k) = \{t(1), \dots, t(j), \dots, t(N_t)\}$ posted by user $u(k)$. Moreover, each tweet $t(j)$ is a sequence of N_w words, $t(j) = \{w(j, 1), \dots, w(j, i), \dots, w(j, N_w)\}$. The problem is to learn semantic low-dimensional representations for all the tweets in the sequences in set U .

4 Proposed Models

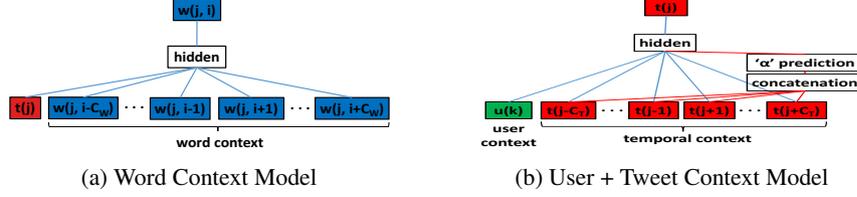


Fig. 2: Architecture diagram of our model.

Our model (Fig 2) learns tweet representations in a hierarchical fashion: learning from the words present in the tweet using word context model (Fig 2 (a)) along with the temporal tweets present in the user stream using tweet context model (Fig 2 (b)). Both the models will be discussed in detail below. Let $\mathbf{w}(\mathbf{j}, \mathbf{i})$, $\mathbf{t}(\mathbf{j})$ and $\mathbf{u}(\mathbf{k})$ denote the embedding for a word i from tweet j , tweet j and user $u(k)$ respectively, all of which have the size ‘ n ’. We will discuss details about both of these models in this section.

4.1 Word Context Model

The goal of the word context model is to learn tweet representations which are good at predicting the words present in the tweet. The model has three layers. The first layer contains the word embeddings, $\mathbf{w}(\mathbf{j}, \mathbf{i} - C_W), \dots, \mathbf{w}(\mathbf{j}, \mathbf{i} - 1), \mathbf{w}(\mathbf{j}, \mathbf{i} + 1), \dots, \mathbf{w}(\mathbf{j}, \mathbf{i} + C_W)$ near the i^{th} target word in tweet j , which denote the word context for the word i (i.e., $w(j, i)$) along with the tweet embedding $\mathbf{t}(\mathbf{j})$. Secondly, there is a hidden layer with size equal to the number of words in the vocabulary ($|V|$). The final layer is a softmax layer which gives a well-defined probability distribution over words in the vocabulary. The input to the word context model is all pairs of word context of word i and tweet $t(j)$ in the corpus. The objective is to maximize the likelihood of the word $w(j, i)$ occurring given its context, i.e., $\mathbb{P}(w(j, i) | w(j, i - C_W), \dots, w(j, i - 1), w(j, i + 1), \dots, w(j, i + C_W), t(j))$. Equation 1 represents the forward propagation step in our 1-hidden layer feed forward model, where W_{WC} and T_{WC} denote the additional parameters of the model.

$$\hat{y}_{|V| \times 1}(j) = \text{softmax}(W_{WC} \times \sum_{l \in \{i - C_W, i + C_W\} \setminus i} \mathbf{w}(\mathbf{j}, \mathbf{l}) + T_{WC} \times \mathbf{t}(\mathbf{j})) \quad (1)$$

4.2 User + Tweet Context Model

The goal of this model is to enrich the tweet representation learned from the word context, by modeling the current tweet conditioned on its temporal context and the proposed user context. The user context makes our model user-aware by exploiting the user characteristics such as the way the user writes the post and also summarizing the topics on which the user writes. These user vectors are learned automatically from the set of tweets posted by the user through this model. As a naïve solution, we can directly adopt Djuric et al. [2]’s approach and apply on the Twitter stream. As discussed in Section 3, this assumption is too strong for social media streams. Can we assign attention levels to the context tweets with respect to the tweet being modeled? To learn the optimal values of attention ($\alpha(j)$), we introduce the attention parameters as shown in Equation 2. The intuition is that semantic loss will be less if

the weights of each of the temporal context tweets are learned accurately. The values of $\alpha(j)$'s can be computed as shown in Equation 3. The objective of this model is to maximize the likelihood of the tweet j posted by user k given its temporal context $(\mathbf{t}(\mathbf{j} - \mathbf{C}_T), \dots, \mathbf{t}(\mathbf{j}-1), \mathbf{t}(\mathbf{j}+1), \dots, \mathbf{t}(\mathbf{j}+\mathbf{C}_T))$ and user context $(\mathbf{u}(\mathbf{k}))$, which is given by $\mathbb{P}(t(j)|t(j - C_T), \dots, t(j - 1), t(j + 1), \dots, t(j + C_T), u(k))$. Since the tweet space can be exponentially large, we use hierarchical softmax [3] instead of normal softmax to bring down the time complexity from $O(|T|)$ (or $O(|V|)$ for the previous model) to $O(\log|T|)$ (or $O(\log|V|)$).

$$\hat{y}_{|T| \times 1}(j) = \text{softmax}(T_{TC} \times \sum_{l \in \{j - C_T, j + C_T\} \setminus j} \alpha(l) \times \mathbf{t}(l)) \quad (2)$$

$$(\alpha(j - C_T) \cdots \alpha(j - 1) \alpha(j + 1) \cdots \alpha(j + C_T)) = \text{softmax}(A[\mathbf{t}(\mathbf{j} - \mathbf{C}_T); \cdots; \mathbf{t}(\mathbf{j}-1); \mathbf{t}(\mathbf{j}+1); \cdots; \mathbf{t}(\mathbf{j}+\mathbf{C}_T);]\}) \quad (3)$$

where the parenthesis inside the softmax function represents concatenation of all context representations $((2 \times C_T \times n) \times 1$ in size). A is the additional weight matrix (of size $(2 \times C_T) \times (2 \times C_T \times n)$) added as parameters to the model. In practice, we observe that multiple passes ('epochs') on the training set are required to fine tune these attention values. The overall objective function intertwining both the models in a hierarchical fashion to be maximized can be summarized as shown in Equation 4. We use the cross-entropy as the cost function between the predicted distribution $\hat{y}(j)$ and target distributions $t(j)$ and $w(j, i)$, for modeling using the temporal and word context respectively. We train the model using back-propagation [4] and Adam [7] optimizer.

$$\mathcal{L}(\theta) = \sum_{u(k) \in U} \left[\sum_{t(j) \in u(k)} \sum_{w(j, i) \in t(j)} \log \mathbb{P}(w(j, i) | w(j, i - C_W), \dots, w(j, i - 1), w(j, i + 1), \dots, w(j, i + C_W), t(j)) + \log \mathbb{P}(t(j) | w(j, 1), \dots, w(j, N_w)) + \log \mathbb{P}(t(j) | t(j - C_T), \dots, t(j - 1), t(j + 1), \dots, t(j + C_T), u(k)) \right] + \log \mathbb{P}(u(k) | t(1), \dots, t(N_T)) \quad (4)$$

5 Experimental Evaluation

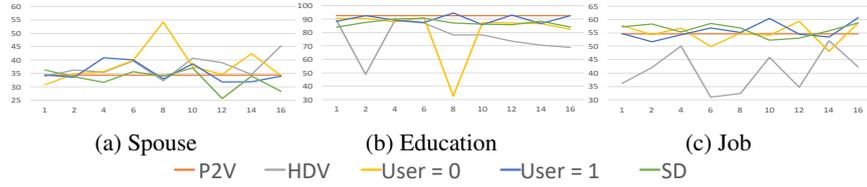
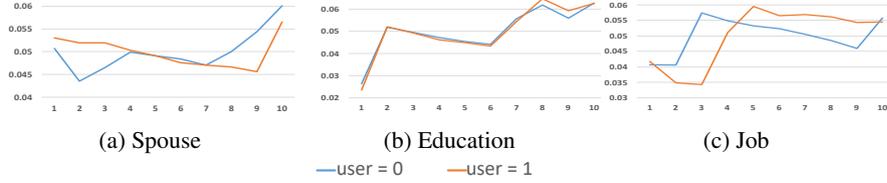
In this section we discuss details of our dataset, experiment, and then present quantitative analysis of the proposed models.

Algorithm	Spouse	Education	Job
Paragraph2Vec [1]	0.3435	0.9259	0.5465
Simple Distance model (SD)	0.3704	0.9068	0.5872
HDV [2]	0.4526	0.8901	0.521
Ours (User = 0)	0.5416	0.9098	0.5935
Ours (User = 1)	0.4082	0.9274	0.6067

Table 1: User profile attribute classification - F1 Score

5.1 Dataset Description

We use the publicly available dataset described in Li et al. [5] for all the experiments. It contains tweets pertaining to three profile attributes (spouse, education and job) of a user. Specifically, it has a set of tweets from users' Twitter timelines, that talk about the attribute ('positive' tweets) and those that do not ('negative' tweets). We randomly sample 1600 users from the dataset and use 70-10-20 ratio to construct train, validation and test splits. Tweet embeddings are randomly initialized while the word embeddings are initialized with the pre-trained word vectors from Pennington et al. [6].

Fig. 3: Model performance w.r.t. temporal context size C_T .Fig. 4: Mean attention w.r.t. distance from the center C_T .

5.2 Experimental Protocol

We consider the binary task of predicting whether a given entity mention corresponds to particular users’ profile attribute or not. We build our model to get the tweet vector and the entity vector by computing an average of all the tweet vectors for the entity. We tune the penalty parameter of a linear Support Vector Machine (SVM) on the validation set. Note that we use a linear classifier so as to minimize the effect of variance of non-linear methods on the classification performance and subsequently help in interpreting the results. We compare our model with three baselines: (1) Paragraph2Vec [1], (2) Simple Distance model (SD): A model that assigns attention weight to the context tweet which is inversely proportional to the distance of the tweet from the target tweet, (3) HDV [2], (4) Ours (User = 0): Our model when the user context is excluded from the temporal context, (5) Ours (User = 1): Our model when the user context is included in the temporal context. We empirically set n and C_W to 200 and 10 respectively for all the models. In case of SD, HDV and our models, we try values in $\{1, 2, 4, 6, 8, 10, 12, 14, 16\}$ to fix the temporal context size parameter (i.e., C_T) which is crucial in improving the semantics of the tweet.

5.3 Comparative analysis

From Table 1, we see that Paragraph2Vec overfits the validation set, resulting in poor accuracy during testing. HDV’s assumption of giving equal attention value to the temporal context also results in lower accuracy compared with our models. SD model outperforms HDV in two tasks, which substantiates our claim against HDV’s naïve assumption for social media. Our model with user vector outperforming the baselines for Education and Job attribute classification, shows the need to consider the user characteristics while modelling his/her tweets. The poor results for Spouse task suggest that this dataset has too many topic shifts and that the user vector turned out to be less accurate. Fig 3 displays the F1 results for different values of C_T , which is a vital parameter controlling the influence of temporal context. We observe that in some cases HDV outperforms the SD model, mainly due to the inability of the SD model to utilize the context information

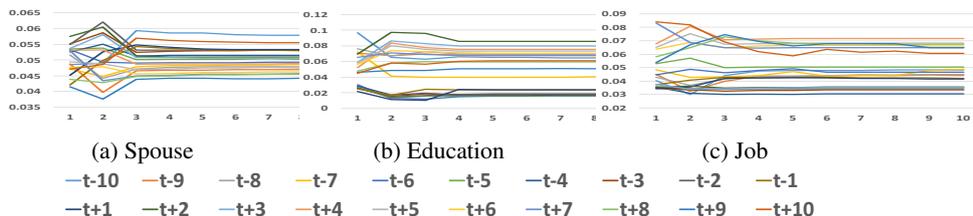


Fig. 5: Mean attention w.r.t. epoch for our model when the user context is included in the temporal context.

from farther tweets which are relevant with respect to the target tweet. Our models are 19.66%, 2.27% and 2.22% better compared to the baselines for the spouse, education and job attributes respectively.

5.4 Impact of Variable Attention

We plot the attention mean across each position of the context tweet with respect to the epoch number. From Fig 5, we see that mean attention at each context position are approximately in the ballpark. Mean attention weights vary for each context position, exhibiting no relation with respect to the increase in distance (as seen in Fig 4). These findings indicate the complexity of giving attention to tweets in the temporal context. Initially, we see that the mean attention weights are changing drastically indicating their sub-optimality. It is interesting to see the convergence of these weights to the optimal solution is fast (in terms of no. of epochs) in the model which uses user context when compared to the model that does not use it.

6 Conclusions

We proposed a model to learn generic tweet representations which have a wide range of applications in NLP and IR field. We discovered that the principled usage of the tweets in the temporal context is an important direction in enriching the representations. We also explored learning a novel user context vector to make our model user-aware while predicting the adjacent tweets. Through experimental analysis, we identified the cases when modeling the user characteristics help enhance the embedding quality. In future, we plan to understand the application potential of the user vector learned through our approach.

References

1. Quoc V. Le and Tomas Mikolov. 2014. Distributed Representations of Sentences and Documents. *Proceedings of the 31th International Conference on Machine Learning, ICML*, 1798–1828.
2. Nemanja Djuric, Hao Wu, Vladan Radosavljevic, Mihajlo Grbovic, and Narayan Bhamidipati. 2015. Hierarchical Neural Language Models for Joint Representation of Streaming Documents and their Content. *Proceedings of the 24th International Conference on World Wide Web, WWW*, 248–255.
3. Frederic Morin and Yoshua Bengio. 2005. Hierarchical Probabilistic Neural Network Language Model. *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics, AISTATS*.
4. David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. 1988. Learning Representations by Back-propagating Errors. *Neurocomputing: Foundations of Research*, 696–699.
5. Jiwei Li, Alan Ritter, and Eduard H. Hovy. 2014. Weakly Supervised User Profile Extraction from Twitter. *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL*, 165–174.
6. Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global Vectors for Word Representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP*, 1532–1543.
7. Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.