

Temporal Recurrent Activation Networks

Giuseppe Manco, Giuseppe Pirrò, and Ettore Ritacco

ICAR - CNR, via Pietro Bucci 7/11C, 87036 Arcavacata di Rende (CS), ITALY,
{name.surname@icar.cnr.it}

Abstract. We tackle the problem of predicting whether a target user (or group of users) will be active within an event stream before a time horizon. Our solution, called **PATH**, leverages recurrent neural networks to learn an embedding of the past events. The embedding allows to capture influence and susceptibility between users and places closer (the representation of) users that frequently get active in different event streams within a small time interval. We conduct an experimental evaluation on real world data and compare our approach with related work.

1 Introduction

There is an increasing amount of streaming data in the form of sequences of events characterized by the time in which they occur and their mark. This general model has instantiations in many contexts, from sequences of tweets characterized by a (re)tweet time and identity of the (re)tweeter and/or the topic of the tweet, to sequences of locations characterized by the time and location of each check-in. We focus on influence-based activation networks, that is, event sequences where the occurrence of an event can boost or prevent the occurrence of another event. Understanding the structural properties of these networks can provide insights on the complex patterns that govern the underlying evolution process and help to forecast future events. The problem of inferring the topical, temporal and network properties characterizing an observed set of events is complicated by the fact that, typically, the factors governing the influence of activations and their dependency from times are hidden. Indeed, we only observe activation times and related marks, (e.g. retweet time) while, activations can depend on several factors including the stimulus provided by the ego-network of a user or his attention/propensity towards specific themes. The goal of this paper is to introduce **PATH** (Predict User Activation from a Horizon), which focuses on scenarios where there is the need to *predict whether a target user (or group of users) will be active before a time horizon T_h* . **PATH** can be used, for instance, in market campaigns where target users are the potential influencers that if active, before T_h , can contribute to further spread an advertisement and trigger the activation of influencees that can consider a given product/service. **PATH** learns

an embedding of the past event history via Recurrent Neural Networks that also cater for the diffusion memory. The embedding allows to capture influence and susceptibility between users and places closer (the representation of) users that frequently get active in different streams within a small time interval.

Related Work. We conceptually separate related research into: (i) approaches like DeepCas [9] and DeepHawkes [2] that tackle the problem of predicting the length that a cascade will reach within a timeframe or its incremental popularity; (ii) approaches like Du et al. [4] and Neural Hawkes Process (NHP) [10] model and predict time event markers and time; (iii) approaches based on Survival Factorization (SF) [1] that leverage influence and susceptibility for time and event predictions; (iv) other approaches that do not use neural networks (e.g., [5, 3, 12]). PATH adopts a different departure point from these approaches: it focuses on predicting the activation of (groups of) users before a time horizon. Differently from (i) PATH considers time and uses an embedding to capture both influence and susceptibility between users and predict future activations. Moreover, (i) focuses on the prediction of cumulative values only (e.g., cascade size). Differently from (ii), we do not assume that time and event are independent. Besides, (ii) focuses on predicting event types (e.g., popular users), which is not enough in the scenarios targeted by PATH (e.g., targeted market campaigns) where one is interested in predicting the behavior of specific users and not their types. As for (iii), it fails in capturing the cumulative effect of history while PATH captures by using an embedding. As for (iv), the main difference is that PATH can automatically learn (via neural networks) an embedding representing influence/susceptibility.

The contributions of the paper are as follows: **(i)** PATH, a classification-based approach based on recurrent neural networks allowing to model the likelihood of observing an event as a combined result of the influence of other events; **(ii)** an experimental evaluation and a comparison with related work.

The remainder of the paper is organized as follows. We introduce the problem in Section 2. We present PATH in Section 3. We compare our approach with related research in Section 4. We conclude and sketch future work in Section 5.

2 Problem Definition

We focus on network of individuals who react to solicitations along a timeline. An activation network can be viewed as an instance of a marked point processes on the timeline, defined as a set $X = \{(t_h, s_h)\}_{1 \leq h \leq m}$. Here, $t_h \in \mathbb{R}_+$ denotes the events of the point process, and $s_h \in \mathbb{M}$ denote the marks in the measurable space \mathbb{M} . Relative to activation networks, the specification of s_h occurs by means of the realizations u_h , c_h and \mathbf{x}_h , where $u_h \in \mathcal{V}$ (with $|\mathcal{V}| = N$) represent individuals, $c_h \in \mathcal{I}$ (with $|\mathcal{I}| = M$) represent solicitations and \mathbf{x}_h is side information which characterizes of the reaction of the entity, described as an instance relative to a feature space of interest. For example, \mathcal{V} can represent users who are engaged in online discussions \mathcal{I} , and the tuple $(t_h, (u_h, c_h, \mathbf{x}_h))$ represents the contribution of u_h to discussion c_h with the post \mathbf{x}_h . It is convenient to view the process as a set of *cascades*: that is, for each $c \in \mathcal{I}$ we can consider the subset $\mathcal{H}^c =$

$\{(t, u, \mathbf{x}) | (t, (u, c, \mathbf{x})) \in X\}$ of elements marked by c , with $m_c = |\mathcal{H}^c|$. Also, \mathbf{t}^c and \mathcal{U}^c represent the projections on the first and second column of \mathcal{H}^c . We also denote by $\mathcal{H}_{<t}^c$ (resp. $\mathcal{H}_{\leq t}^c$) the set of events $e_i \in \mathcal{H}^c$ such that $t_i < t$ (resp., $t_i \leq t$). The terms $\mathbf{t}_{<t}^c$ and $\mathcal{U}_{<t}^c$ can be defined accordingly. The relationship $u \prec_c v$ denotes that both u and v are active in \mathcal{H}^c and there are some events relative to u and v such that u precedes v in some events. Finally, $\mathcal{C} = \{\mathcal{H}^1, \dots, \mathcal{H}^M\}$ denote a collection of M cascades over \mathcal{V} and \mathcal{I} .

Modeling diffusion. We start from the observation that what is likely to happen in the future (viz. which user will be active and when) depends on what happened in the past (viz. the chain of previously active users). One important point to take into account is the susceptibility of users, that is, the extent to which they are influenced by specific previously activated users. Our model should be flexible enough to reflect both *exciting* and *inhibitory* effects. While the former boosts the likelihood of observing u active in c , the latter actually could prevent it to do so. Given a cascade \mathcal{H}^c , a timestamp $t \geq 0$ and a user $u \notin \mathcal{U}_{<t}^c$, the goal is to obtain an estimate of the density function $f(t, u | \mathcal{H}_{<t}^c)$, which can be used to model the following evolution scenario: *given a time horizon T_h^c ; how likely is it that u will become active in c within T_h^c ?*

The challenge, at this point, is how to concretely formulate the density f . We can decouple its specification as follows:

$$f(t, u | \mathcal{H}_{<t}^c) = g(t | u, \mathcal{H}_{<t}^c) \cdot h(u | \mathcal{H}_{<t}^c), \quad (2.1)$$

where the first component represents the likelihood that u becomes active within t , given the current history, and the second component represents the likelihood that u activates (independent of the time) as a reaction to the current history.

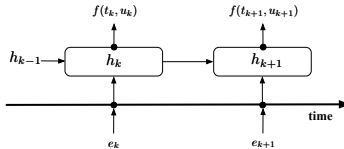
As for $\mathcal{H}_{<t}^c$, explicit information includes features like the sequence of user activations, their activation times, the relative activation speed, and possibly the topic of the cascade. Nevertheless, our assumption is that $\mathcal{H}_{<t}^c$ can also encode latent information including susceptibility and influence between users that can be derived, for instance, from neighborhood information in a network (e.g., follower/followee relations in Twitter) or user behaviors (e.g., users that retweet after a certain set of other influential user (re)tweet). This is exactly what we want to unveil in our modeling.

Embedding history. We want to learn an embedding of users in a latent K -dimensional space such that users in the same cascade are closer in the embedding, and users within different cascades are distant. We make usage of two matrices $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_N]$, $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_N] \in \mathbb{R}^{N \times K}$ that represent the susceptibility and influence, respectively. Matrices are computed by relying on the standard network architecture borrowed from the `word2vec` paradigm [11]:

$$\mathbf{a}_v = \mathbf{W}_e \mathbf{v} \quad \mathbf{s}_u = \mathbf{V}_e \mathbf{u}$$

Here, \mathbf{u} , \mathbf{v} represents the one-hot encodings of u and v . The matrices \mathbf{W}_e , \mathbf{V}_e represent the embeddings, obtained by minimizing an adapted form of contrastive loss [6] that penalizes the distance of users within the same cascades and the closeness of users in different cascades.

Capturing the diffusion memory. To encode temporal relationship within $\mathcal{H}_{<t}^c$ we use recurrent neural networks (RNNs). An RNN is a recursive structure that, at the current step, gets as input the previous network state (the outputs form the hidden units) along with the current input to compute a new state. The following picture provides an overview of a simple RNN cast to our context. At each step k , we feed into the network an event $e_k \in \mathcal{H}^c$ that encodes the current user (u_k) and its activation time (t_k). The learned hidden state (h_k) represents the non-linear dependency between these components and past events, which can be used to model $f(t_k, u_k | \mathcal{H}_{<t_k}^c)$. In the following, we adopt the LSTM instantiation of the RNN framework [7]. The idea of an LSTM unit is to reliably transmitting important information many time steps into the future. At every time step, the unit modifies the internal status by deciding which part to keep or replace with new information coming from the current input. We use the shortcut $\mathbf{h}_k = \text{LSTM}(\mathbf{z}_k, \mathbf{h}_{k-1})$ to denote a functional architecture that elaborates an input \mathbf{z}_k and outputs the updated state.



3 PATH: Predicting User Activation from a Horizon

We now introduce PATH (Predicting User Activation from a Horizon), which focuses on simplifying $f(t, u | \mathcal{H}_{<t}^c)$ as the binary response function $\mathbb{I}(t \leq T_h | u, \mathcal{H}_{<t}^c)$ that denotes whether u becomes active in c within T_h . We focus on events $e_k \in \mathcal{H}^c$ where the features of interest \mathbf{x}_k are limited to the time delay $\delta_k = t_k - t_{k-1}$ relative to the previous activation within the cascade. This allows us to capture the property that cascades may have intrinsically different diffusion speeds causing some of them to concentrate users' activations in a short timeframe while others in a more extended interval. Given a partially observed cascade $\mathcal{H}_{<t_l}^c$ (with $t_l < T_h^c$ representing the timespan of the observation window), our objective is to predict, for a given entity $u \notin \mathcal{U}_{t_l}^c$, whether $u \in \mathcal{U}_{T_h^c}$.

In order to uncover all the characteristics of the activations within cascades, we consider a model built on all possible prefixes of the available cascades. Notice that, in our reconstruction, we do not consider the first element within the cascade, which we assume becomes “spontaneously” active. Finally, for each $u \notin \mathcal{U}^c$, we associate the cascades $\mathcal{H}_{\leq t_j}^c \cup \{(t_j, u, \delta_j)\}$ (with $1 \leq j \leq m_c - 1$) and $\mathcal{H}^c \cup \{(T_h^c, u_i, (T_h^c - t_{m_c}))\}$ with negative labels. Again, the intuition is that, since u is not active no partial cascade provides the sufficient intensity to activate u within the given time horizon. Adding negative examples in the data preparation represent an effective data augmentation process, which enlarges the training data by inferring new inputs in the training set. This is crucial to let the approach better fine tune separation between active and inactive users, as well as better characterize the true activation time of active users. Let $\mathcal{T}_{\mathcal{C}}$

denote the set of all pairs partial sequence/associated label that can be built from the above discussion. Our idea is to exploit the embedding and LSTM tools described in the previous section to solve the supervised problem at hand. Figure 1 illustrates the basic architecture of the model. Given a pair $\langle \mathcal{H}_i, y_i \rangle \in \mathcal{T}_C$ with $|\mathcal{H}|=n$ and by considering $e_k = (t_k, u_k, \delta_k) \in \mathcal{H}$ (with $1 \leq k \leq n$), the architecture of the network can be captured by the following equations:

$$\mathbf{a}_k = \mathbf{W}_e \mathbf{u}_k \quad (3.1)$$

$$\mathbf{h}_k = \text{LSTM}([\mathbf{a}_k, t_k, \delta_k], \mathbf{h}_{k-1}) \quad (3.2)$$

$$\hat{y}_i = \sigma(\mathbf{W}_o \mathbf{h}_n) \quad (3.3)$$

$$\tilde{y}_i = \exp \left\{ - \left\| \mathbf{a}_n - \sum_{k=1}^{n-1} \mathbf{a}_k \right\|^2 \right\} \quad (3.4)$$

Here, \hat{y}_i represents the probability that y_i is positive, as provided by the network: that is, it encodes the probability that u_n becomes active within t_n ; \tilde{y}_i encodes the affinity between u_n and all users preceding it within \mathcal{H}_i . The distance $\|\mathbf{a}_n - \sum_{k=1}^{n-1} \mathbf{a}_k\|$ plays a crucial role here: since the target user is on the tail of the cascade, the embedding should emphasize the similarities with the predecessors that trigger an activation, and by the converse minimize the similarities with those ones which do not trigger it. The loss is a combination of cross-entropy and the embedding loss previously described:

$$\mathcal{L} = \sum_{\substack{\langle \mathcal{H}, y \rangle \in \mathcal{T}_C \\ |\mathcal{H}|=n}} \{y(\gamma \log(\hat{y}) + \beta \log \tilde{y}) + (1-y)(\gamma \log(1-\hat{y}) + \beta \log(1-\tilde{y}))\} \quad (3.5)$$

where γ and β are weights balancing cross-entropy and embedding.

4 Experiments

We validate our approach by analysing the algorithm on real-life datasets. In particular, we analyse the capability of the algorithm at predicting the activation time of users within an information cascade. The implementation we use in the experiments can be found at <https://github.com/gmanco/PATH>.

Datasets. We evaluated the prediction capability of PATH by exploiting two real-world datasets containing propagation cascades crawled from the timelines of **Twitter**² and **Flixster**.³ In particular, **Twitter** includes $\sim 32\text{K}$ nodes with $\sim 9\text{K}$ cascades while **Flixster** includes $\sim 2\text{K}$ nodes with $\sim 5\text{K}$ cascades.

² <http://www.twitter.com/>

³ <http://www.flixster.com/>

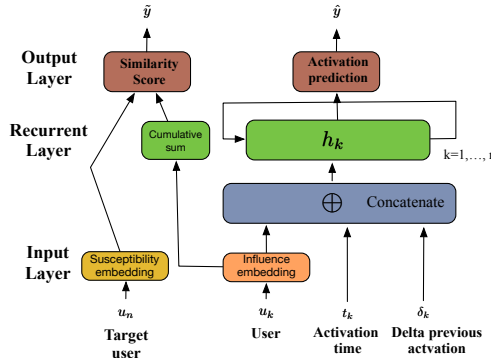


Fig. 1: Overview of PATH.

The information propagation mechanism on **Twitter** is expressed by retweeting, in other words a chain of repetitions and transmissions of a tweet from a set of users to their neighbors in a recursive process. Each activation corresponds to a retweet. An activation in **Flixster** happens when a user rates a movie, while a cascade is composed by all the activations related to the same movie. The two datasets differ essentially for the following characteristics: **Twitter** includes a larger number of users and shorter delays than **Flixster**. In addition, retweets intuitively highlight two relevant aspects, namely the importance of the topic and the single influence of the individual from which the retweet is performed. By contrast, movie ratings are more likely to exhibit a cumulative effect: popular movies are more likely to be considered than unpopular ones.

Evaluation Methodology. We evaluate **PATH** against two baseline models, both relying on Survival Analysis [8]. The first instantiation implements a Cox proportional hazard model (*CoxPh* in the following). We implement the model using the `lifelines`⁴ package and extract, for each event $e_k \in \mathcal{H}^c$, the following features: (1) size of the prefix; (2) last activation time; (3) average delay for each active user so far; (4) number of neighbors in the history, and (5) coverage percentage of them within the history; (6) the activation time of the most recent neighbor, if any; (7) correlation between the activation of the current user and its neighbors within the history, computed in previous cascades. This model represents an intuitive baseline where features are manually engineered and include a mix of external information (coming from the underlying network neighborhood) and information derived from the cascade itself. The second instantiation is given by the *Survival Factorization* (*SF* in the following) framework described in [1]. The comparison is important since *SF* relies on the same guiding ideas of **PATH** (influence/susceptibility) with the difference that there is no cumulative effect of $\mathcal{H}_{<t_k}^c$, but instead an influential user has to be detected for each activation.

To evaluate the approaches we proceed as follows: given training and test sets \mathcal{C}_{train} and \mathcal{C}_{test} , we train the model on \mathcal{C}_{train} and measure the accuracy of the predictions on \mathcal{C}_{test} . The two sets are obtained by randomly splitting the original dataset by ensuring that there is no overlap among the cascades of the two sets, but there is no entity in the test that has not been observed in the training. For the evaluation, we chronologically split each cascade $c \in \mathcal{C}_{test}$ into c_1 and c_2 such that, for each $u \in c_1$ and $v \in c_2$, we have that $u \prec_c v$. Next, we pick a random subsample $c_3 \subseteq \mathcal{V} - \mathcal{U}^c$. Then, given a target horizon T_h^c , we measure *TP*, *FP*, *TN* and *FN* by feeding the models on c_1 and then predicting the activation within T_h^c for each element in $c_2 \cup c_3$. The choice of T_h^c can follow different strategies; Fixed horizon (FH): setting T_h^c as the maximum observed activation time $T_h^{test} = \max\{t | t \in \mathbf{t}^c, c \in \mathcal{C}_{test}\}$; Variable horizon (**VH**): varying T_h^c from the smallest to the largest activation time and computing the activation probabilities associated to each possible value; Actual Time (**AT**): a particular case of the VH strategy, where $T_h^c \triangleq T_h^{u,c}$ is relative to the true activation time in c of each user $u \in c_2 \cup c_3$.

⁴ see <http://lifelines.readthedocs.io> for details.

We plot the ROC and the F-Measure curves relative to the above alternatives and report the AUC and F values. For *PATH*, the encoding of sequences as described in section 3 already presumes that users are evaluated on intermediate timestamps prior to their actual activation. Thus, VH and AT roughly coincide in this case. Since both *CoxPh* and *SF* are capable of inferring, for each (u, \mathcal{H}) pair, the probability $S_u(t|\mathcal{H})$, the comparison with *PATH* is done by computing $1 - S_u(\tilde{t}|\mathcal{H})$ where \tilde{t} is the horizon timestamp. The parameter space for *PATH* was explored by grid-search, measuring the loss on a separate portion of the training set by 5-fold cross-validation. We report 5 different instantiations, which differ from the number of cells in the LSTM (32/64), the dimensionality of the embedding (32/64) and the batch size in the training (128/256/512). Concerning *SF*, the number of factors was set to 16 for both datasets.

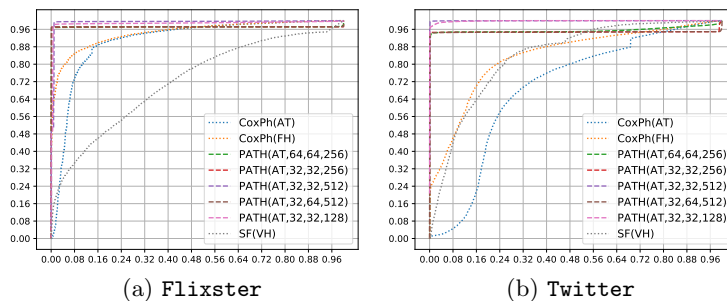


Fig. 2: ROC Curves for *PATH*, *CoxPh* and *SF* on both datasets.

Evaluation Results. Figure 2 reports the ROC curves where we observe that *PATH* outperforms the baselines and in particular exhibits a very good accuracy on all configurations. This is especially true on *Flixster*, where by the converse *SF* does not seem capable of correctly correlating previous activations times. The cumulative influence effect is evident here, as a natural consequence of the underlying domain where cascading effects are more likely as a consequence of a “word of mouth” process. On *Twitter*, where the activation is more likely due to the influence of a single user (as testified by the good performance of *SF*), *PATH* still achieves the best scores, thus proving the capability of the recurrent layer to adapt the influence to a single user. By analyzing Fig. 3, which displays the F-measure curve for varying values of the threshold on the probabilities, we can observe that, contrary to the baselines, higher thresholds do not cause a significant drop of the recall. The only exception is *CoxPh* (FH), which seems more stable on *Flixster*. This is a clear sign that the probabilities associated with active and inactive users in *PATH* differ substantially, and in particular active events are associated with significantly higher probabilities than inactive events.

5 Concluding Remarks and Future Work

We focused on the problem of predicting user activations in a given time horizon and show that the embedding of the user activation history, where users that

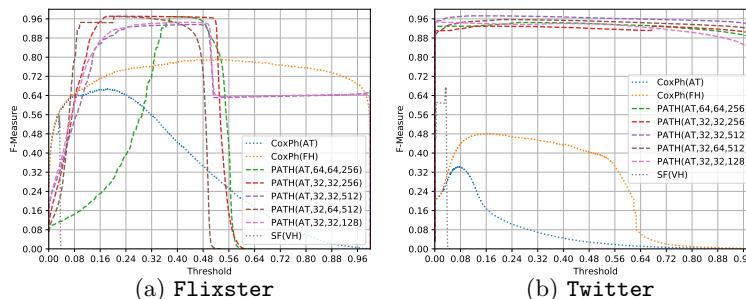


Fig. 3: F-Measure curves for PATH, *CoxPh* and *SF* on both datasets.

become active on the same cascades are placed close, can be effectively learned via recurrent neural networks. Experiments performed on real datasets show the effectiveness of the approach in accurately predicting next activations. It is natural to wonder whether it is possible to cast the intuitions behind our approach in a generative setting, to predict both which user is likely to become active, and the time segment upon which s/he will become active.

References

1. N. Barbieri, G. Manco, and E. Ritacco. Survival factorization on diffusion networks. In *PKDD*, pages 684–700, 2017.
2. Q. Cao et al. Deephawkes: Bridging the gap between prediction and understanding of information cascades. In *CIKM*, pages 1149–1158, 2017.
3. Peng Cui, Shifei Jin, Linyun Yu, Fei Wang, Wenwu Zhu, and Shiqiang Yang. Cascading outbreak prediction in networks: a data-driven approach. In *KDD*, pages 901–909. ACM, 2013.
4. N. Du, H. Dai, R. Trivedi, U. Upadhyay, M. Gomez-Rodriguez, and L. Song. Recurrent marked temporal point processes: Embedding event history to vector. In *KDD*, pages 1555–1564, 2016.
5. A. Guille and H. Hacid. A predictive model for the temporal dynamics of information diffusion in online social networks. In *WWW*, pages 1145–1152, 2012.
6. R. Hadsell, S. Chopra., and Y. LeCun. Dimensionality reduction by learning an invariant mapping. In *CVPR*, pages 1735–1742, 2006.
7. S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
8. J.D. Kalbfleisch and L.P. Ross. *The Statistical Analysis of Failure Time Data*. Wiley Series in Probability and Statistics, 2002.
9. C. Li, J. Ma, X. Guo, and Q. Mei. Deepcas: An end-to-end predictor of information cascades. In *WWW*, pages 577–586, 2017.
10. H Mei and J.M. Eisner. The neural hawkes process: A neurally self-modulating multivariate point process. In *NIPS*, pages 6757–6767, 2017.
11. T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
12. L. Yu, P. Cui, F. Wang, C. Song, and S. Yang. From micro to macro: Uncovering and predicting information cascading process with behavioral dynamics. In *ICDM*, pages 559–568, 2015.