

# Efficient Learning with Exponentially-Many Conjunctive Precursors to Forecast Spatial Events

**Abstract**—Forecasting spatial societal events in social media is significant and challenging. Most existing methods consider the frequencies of keywords or n-grams to be features, but have not explored the exponentially large space of the conjunctions of those features, such as keyword co-occurrence in messages, which can serve as crucial precursor rules. Due to the inherent exponential complexity of ensemble rule learning, existing work typically adopts greedy/heuristic strategies. This means that they cannot guarantee the solution’s optimality, which would require a considerably more sophisticated model for spatial event forecasting, while still suffering from major challenges: 1) Exponentially-dimensional feature learning with distant supervision, 2) Numerical values of conjunctive features, and 3) Spatially heterogeneous conjunction patterns. To concurrently address all these challenges with a theoretical guarantee, we propose a novel spatial event forecasting model which learns numerical conjunctive features efficiently. Specifically, to consider their magnitude, traditional Boolean rules are innovatively generalized to deal with numerical conjunctive features with amenable computational properties. To handle the geographical similarity and heterogeneity in numerical conjunctive feature learning, we propose a new model that implements through a new bi-space hierarchical sparsity regularization for locations and features. Moreover, we propose a new algorithm to optimize the model parameters and prove that it enjoys theoretical guarantees for both the error bounds and time efficiency. Extensive experiments on multiple datasets demonstrate the effectiveness and efficiency of the proposed method.

## I. INTRODUCTION

Currently, user-generated contents such as microblogs have become ubiquitous, which serve as real-time “sensors” for social trends and incidents [26]. People use social media to plan, advertise, and organize future social events such as the planned protests in the “Arab Spring” and “Occupy Wall Street” [25]. The predictive power of microblogs for social event forecasting has been widely explored by a great deal of recent research on topics such as crimes [14], civil unrest [38], and disease outbreaks [2]. These research works share essentially similar workflows. First, the model features are typically defined as the counts of terms (e.g., keywords and hashtags) under the domain of interest. The feature values in the aggregated collections of massive microblogs are considered to jointly reflect the social tendencies. The predictive model is then trained to map the social indicators to the model response, in this case the occurrence of future events.

However, the count for a single keyword may not be sufficiently informative to serve as a precursor for forecasting social events. For example, Figure 1(a) shows that instead of either “teacher” or “reform”, the count of their conjunction in the same tweets reflects the public concern regarding edu-

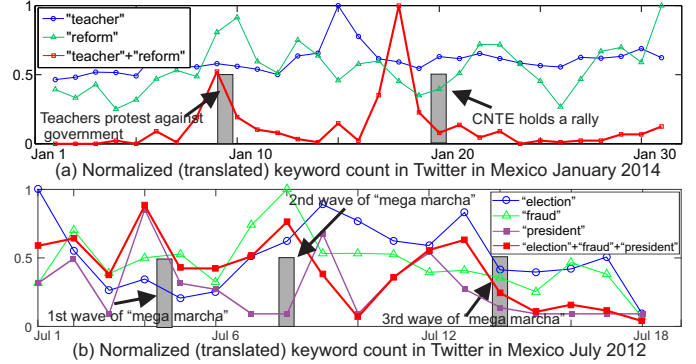


Figure 1: The burstiness of the counts of some keyword co-occurrences preceded the major events. The keyword co-occurrence features could be much more determinative and interpretable than single keywords as precursors for future events.

cational policy. Similarly, as shown in Figure 1(b), the count for anyone of “election”, “president”, and “fraud” individually is a very noisy signal, while the massive co-occurrence of them all in the same tweets is a very informative precursor for the subsequent three waves of protests against the results of the presidential election in Mexico. Therefore, unlike the incidence of single keywords, the co-occurrence of keywords, such as “president+election+fraud” typically conveys much more definite meaning and is thus a significantly more powerful precursor for future protest events. In this paper, we call such new features *conjunctive features*, which in this case refer to two or more co-occurring keywords in Figure 1, and the standalone atomic features as *primitive features*, which here mean single keywords.

*Conjunctive features* are highly informative and thus more interpretable by human observers, which is crucial if decision makers are to understand and utilize the predictive models. However, it is impossible for domain experts to manually provide an extensive set of all the keyword conjunctions that are precursors. Better methods for automatically learning an extensive set of significant keyword conjunctions from the data are clearly required, but due to the exponential complexity in storage and time of computation, this problem is conventionally unfeasible even for a moderate sized of keyword set. For example, among as few as 100 keywords, there are  $2^{100} \approx 10^{30}$  possible combinations of conjunctive features to store and compute, which are far beyond the existing memory capacity and computational power.

Existing methods on supervised rule mining typically utilize greedy or heuristic-based methods [12], where the optimality of the solution cannot be guaranteed. Despite the significance

of this problem, to the best of our knowledge, there has been very little work reported on spatial event forecasting that extensively considers conjunctive features. Even utilizing the simplest settings, several substantial theoretical and practical challenges make this problem unfeasible to solve: **1) Numerical values of conjunctive features.** Existing methods for conjunctive feature learning typically require the conjunctive features to be Boolean-valued if efficient computation is to be feasible. However, in spatial event forecasting, instead of binary values, the numerical frequencies of the keyword conjunctions occurring in the same location at a specific time serves as the indicator, which does not satisfy the Boolean assumption that is universally applied in existing work. **2) Exponentially-dimensional learning with distant supervision.** Due to the immense volume of microblog messages, it is prohibitively labor-intensive to label each individual message. When forecasting spatial events, typically only labels at the aggregate level (e.g., city-level) are available, which can thus only provide distant supervision when learning important conjunctive features. **3) Spatially heterogeneous conjunction patterns.** Different geo-locations may share similar conjunctive features but also have their own exclusive ones within a particular geo-neighborhood. For example, “occupy+street” can be a good indicator for general civil unrest across many different geo-locations, but “occupy+wall+street” typically appears in New York State while “occupy+Texas+state” typically happens in Texas.

In order to simultaneously overcome all the above-mentioned challenges, we propose a novel model named Hierarchical-Task Numerical conjunctive feature Learning (HTNL) for spatial event forecasting. Specifically, a novel kernel is formulated to represent every possible numerical conjunctive feature (NCF) and then this exponentially large set of kernels is correlated via a Directed Acyclic Graph (DAG). The complexity in NCF selection is reduced from exponential to polynomial by utilizing the sparsity structure from the DAG and the favorable computational properties of the proposed kernel. Finally, the similarities between the selected NCFs within geo-hierarchical neighborhoods are enforced to boost model generalizability using the newly proposed bi-space regularization strategy in both feature and location spaces. The major contributions of this paper are as follows:

- **Develop a generic framework for conjunctive precursor learning for spatial event forecasting.** A generic framework is proposed for spatial event forecasting that optimally learns the NCFs, taking into account both geographical similarity and heterogeneity. A number of classic approaches in related research are shown to be special cases of our model.
- **Propose a novel hierarchical multitask model for NCF learning.** First, every possible NCF is formulated as a novel kernel with structured sparsity on a DAG. Then the similarity of sparsity patterns is enforced using a newly proposed bi-space regularization strategy that utilizes geo-hierarchical knowledge to boost up model generalizability.

- **Design an efficient optimization method with a theoretical guarantee of optimality.** The proposed model requires the optimization of an NCF set that is exponentially large and geographically correlated. The new algorithm leverages both the topological sparsity among NCFs and the computational efficiency of the proposed kernel, and provides theoretical guarantees for both error bounds and time complexity.
- **Conduct extensive experiments for performance evaluations.** The proposed method is evaluated on multiple datasets in different domains and found to significantly outperform the existing methods in prediction performance. Moreover, the conjunctive features discovered by the model clearly demonstrate its effectiveness and interpretability.

The rest of this paper is organized as follows. Section II reviews the background and related work and Section III introduces the problem setup. Sections IV and V presents our proposed model and an efficient model parameter optimization algorithm, respectively. The experiments on 8 real-world datasets are presented in Section VI, and the paper concludes with a summary of the research in Section VII.

## II. RELATED WORK

**Event Detection and Forecasting in Social Media.** A considerable amount of work has been done on detecting ongoing events, including disease outbreaks [28], earthquakes [27] and various other types of events [36]. Generally, for event detection, either classification or clustering is utilized to extract tweets of interest and then the spatial [27], temporal [28], or spatiotemporal burstiness [11] of the extracted tweets is examined to identify the potential occurrence of an ongoing event. However, these approaches typically uncover events only after they have commenced. To forecast future events, several event forecasting methods have been proposed, most of which focus on temporal events and ignore the underlying geographical information, such as the forecasting of elections [32], stock market movements [6], disease outbreaks [2], box office ticket sales [4], crimes [34], and others [37], [9]. These works typically utilize linear/nonlinear regression models [4], [6], [16] or time series-based methods [2]. Few existing approaches provide true spatiotemporal resolution for predicted events. In [14], Gerber utilized a logistic regression model for spatiotemporal event forecasting using topic-related tweet volumes as features, while Ramakrishnan et al. [25] built separate LASSO models for different locations to predict the occurrence of civil unrest events. Zhao et al. [37] proposed a multi-task learning framework for event forecasting that jointly learns multiple related spatial locations. But it requires extra knowledge on dynamic features.

**Rule Ensemble Learning (REL).** Given a set of basic propositional features describing the data, the goal of REL is to supervisedly learn a set of feature conjunctions with good predictability. To handle the inherent exponential complexity of this problem, many REL methods have been proposed majorly in three categories: 1) *filter-based methods*, which assume that important conjunctive features must be frequent and thus

only retain frequent instances for classification [8], [10], [29]. However, frequency and the predictability of features are not equivalent because predictability is dependent on the specific prediction task while frequency is not; 2) *heuristic/boosting-based methods*, where researchers address the challenge in Category 1, by learning the feature conjunctions and the predictive model concurrently [26], [23], [12]. To ensure computational efficiency, heuristic strategies based on greedy or boosting methods are generally utilized. And only sub-optimum or local optimum can be found. 3) *optimization-based methods*: To address the problem in Category 2 and ensure efficiency, recently few methods have been proposed for conjunctive feature selection with theoretical guarantee on the error bound to the global optima [5], [18]. This is achieved by utilizing an active set algorithm to scale down the solution space. To check the optimality of current active set efficiently, the “product-of-sum” property [5] of Boolean rules must be exploited. However, existing optimization-based methods do not apply in more general situations where the rules are numerical because they violate the “product-of-sum” property. Classic approaches such as discretizing the numerical values into multiple binary features arbitrarily scale up the number of basic propositional features and thus exponentially enlarge solution space. In order to address this problem, our paper proposes a new method that can directly handle numerical rules efficiently without discretization.

**Multi-task learning:** Multi-task learning (MTL) learns multiple related tasks simultaneously to improve generalization performance [3], [22]. Many MTL approaches have been proposed over the last decade [39]. In [19], Kim et al. proposed a regularized MTL which constrained the models of all tasks to be close to each other. The task relatedness can also be modeled by constraining multiple tasks to share a common underlying structure, e.g., a common set of features [35], or a common subspace [1]. MTL approaches have been applied in many domains, including computer vision and biomedical informatics.

### III. PROBLEM FORMULATION

In this section, the problem in this paper is formulated. Section 3.1 poses the problem of “precursor rule learning for event forecasting”. Denote  $X = \{X_{s,t}\}_{s,t}^{S,T}$  as a collection of microblog data, where  $T$  is the set for time intervals and  $S$  is the set of the spatial locations.  $X_{s,t}$  denotes the data for  $t$ th time interval (e.g.,  $t$ th date) at location  $s$  such that  $X_{s,t} \in \mathbb{Z}^{n_{s,t} \times |V|}$ , where  $n_{s,t}$  denotes the number of microblog messages sent during time interval  $t$  at location  $s$ , and  $|V|$  denotes the size of the vocabulary  $V$ , which is a set of *primitive features* that can include occurrences of specific keywords, hashtags, and hyperlinks.  $X_{s,t}$  is defined as a matrix whose element  $[X_{s,t}]_{i,v} \in \{0, 1\}$  denotes the occurrence (with value 1) or not (with value 0) of the primitive feature  $v$  in the  $i$ th message in location  $s$  during time interval  $t$ . The important notations in this paper are listed in Table I.

As explained earlier, a *conjunctive feature* (or *feature conjunction*) is defined as the conjunction of a set of distinct

Table I: Important Notations

Notations	Explanations
$X_{s,t}$ and $Y_{s,t}$	Input data and event occurrence in location $s$ at time $t$
$\mathcal{V}$	The set of all the conjunctive features
$\phi_v(X_{s,t})$	The frequency of $v \in \mathcal{V}$ in location $s$ at time $t$
$D(v)$ and $A(v)$	The sets of descendants and ancestors of $v \in \mathcal{V}$
$W_s$	Weight vector for the conjunctive features in location $s$
$G$	The set of geographical neighborhoods
$\Theta_{d,r}$	Boundary and interior of a generalized $d$ -simplex
$\alpha_{s,t}$	dual variable for location $s$ at time $t$
$p, \hat{p}$ , and $\bar{p}$	$p, \hat{p}$ , and $\bar{p}$ -norms, $\bar{p} = \hat{p}/(\hat{p} - 1)$ , $\hat{p} = p/(2 - p)$

primitive features such as keywords that co-occur in the same message. Hence, the set of all the possible conjunctive features is denoted as  $\mathcal{V} = \{v | v \subseteq V\}$ , whose size is  $|\mathcal{V}| = 2^{|V|}$ .  $\phi_v(X_{s,t}) \in \mathbb{R}^+ \cup \{0\}$  denotes the frequency of the conjunctive feature  $v$  in location  $s \in S$  at time  $t \in T$ . Therefore, instead of assigning this a Boolean value, in our problem the conjunctive feature is generalized to a numerical value, referred to as the **numerical conjunctive feature (NCF)**.

NCFs have topological relationships with each other. We denote these relationships using a *directed acyclic graph (DAG)* known as a *feature conjunction lattice*:  $\mathcal{G}(\mathcal{V}, \mathcal{E})$ , as illustrated in Figure 2(b). In a feature conjunction lattice, the top node (i.e., Level 0) is an empty conjunctive feature while the nodes in Level 1 are the primitive features  $V$ . A node  $v_1 \in \mathcal{V}$  is called the *parent* of another node  $v_2 \in \mathcal{V}$  if  $v_2 \subset v_1$  and  $|v_2| + 1 = |v_1|$ ; hence  $v_2$  is a *child* of  $v_1$ . Let  $D(v)$  and  $A(v)$  denote the set of descendants and ancestors of  $v \in \mathcal{V}$ , respectively. We assume that both  $D(v)$  and  $A(v)$  include the node  $v$ . For a subset of nodes  $\mathcal{U} \subset \mathcal{V}$ , we define the *hull* and *sources* of  $\mathcal{U}$  as  $H(\mathcal{U}) = \bigcup_{v \in \mathcal{U}} A(v)$  and  $S(\mathcal{U}) = \{v | A(v) \cap \mathcal{U} = \{v\}\}$ , respectively.  $|\mathcal{U}|$  denotes the number of NCFs in set  $\mathcal{U}$  while  $\bar{\mathcal{U}}$  denotes the complementary set of  $\mathcal{U}$ , namely all the NCFs that are in  $\mathcal{V}$  but not in  $\mathcal{U}$ .

Define  $Y = \{Y_{s,t}\}_{s,t}^{S,T}$  as the event occurrences, where  $Y_{s,t} \in \{1, -1\}$  such that  $Y_{s,t} = 1$  means there is an event in location  $s$  at time  $t$ , otherwise  $Y_{s,t} = -1$ . The following is our problem definition of the forecasting task for each location  $s \in S$ , given the microblog data  $D_s = \{D_{s,t}\}_t^T$  for location  $s$ , and the primitive feature set  $V$ , our goal is to discover the set of NCFs  $\mathcal{U}_s \subseteq \mathcal{V}$  that are crucial precursors for a future event in each location  $s$ , and thus learn a mapping function for event forecasting:

$$f : \{\phi_v(X_{s,t})\}_{v \in \mathcal{U}_s} \rightarrow Y_{s,\tau} \quad (1)$$

where  $\tau = t + q$ , and  $q$  is the lead time for forecasting. Among all of the  $|\mathcal{V}| = 2^{|V|}$  candidate NCFs, typically only a few are useful precursors for forecasting.

There are three technical challenges involved in solving this problem. **First, exponential solution space.** This problem is extremely difficult to solve even for a modest size of  $V$  because of the exponentially large size of  $|\mathcal{V}|$ , which causes intractability in both memory and computation. **Second, numerical values of conjunctive features.** To ensure an efficient solution, the state-of-the-art methods require a conjunctive feature to have a Boolean value. However, in our problem the conjunctive feature value  $\phi_v(X_{s,t})$  must be numerical and therefore the Boolean assumption is not satisfied, creating a

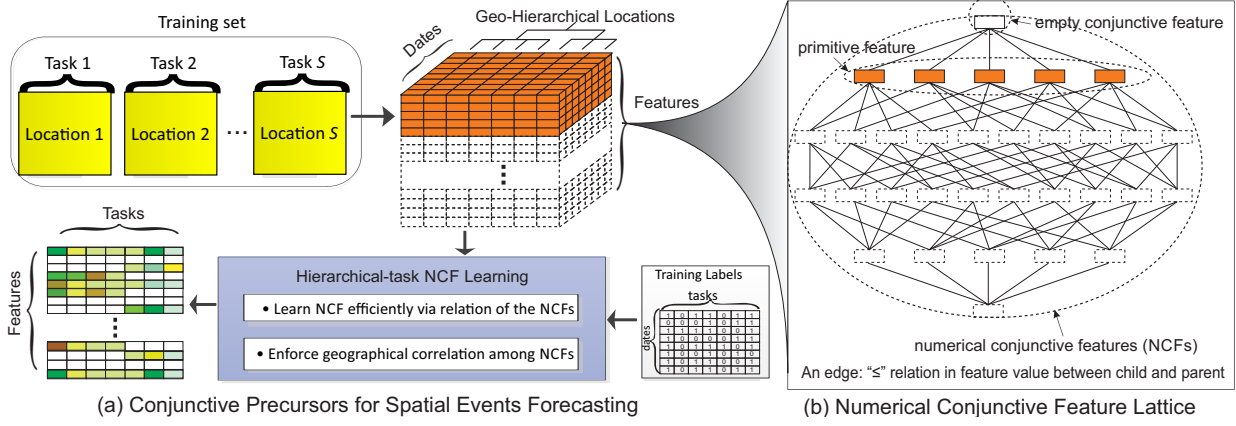


Figure 2: The proposed hierarchical-task numerical conjunctive feature learning (HTNL). (a) The flowchart of proposed HTNL. (b) The NCF lattice where the edge denotes the “ $\leq$ ” relation among NCFs

serious challenge for the model efficiency. **Third, geographical influences in NCF learning.** For spatial event forecasting, both the geographical relationship and heterogeneity in the conjunctive feature learning are crucial and must be considered, a combination that has never been addressed by existing methods. To address all three of these challenges, we have developed the new model presented in the following section.

#### IV. MODEL

In this section, we propose a new model, HTNL, to address the challenges described above. First, NCF is mathematically defined and analyzed. Second, geographical relationships and heterogeneity are considered for NCFs in spatial event forecasting problem. Third, the objective function and its dual form are proposed.

##### A. Computational Properties of NCFs

The state of the art requires a Boolean value for conjunctive features to ensure efficiency, because this is the only way the conjunctive features can be efficiently computed by multiplying the primitive features that they consist of. However, for our problem of spatial event forecasting, the Boolean assumption is not satisfied and a new and generic version, namely NCF, is required. To ensure the computational efficiency is retained in such a generalized setting, the unique formulation and properties of NCF are explored in the following.

###### 1) Calculation of the NCF.

As noted in Section III, the value of NCF  $v \in \mathcal{V}$  in location  $s$  at time  $t$ , namely  $\phi_v(X_{s,t})$ , is defined as the spatiotemporally accumulated occurrence of NCF  $v$ . Given that  $[X_{s,t}]_{v,i} \in \{0, 1\}$ ,  $\phi_v(X_{s,t})$  can then be calculated as:

$$\phi_v(X_{s,t}) = \sum_i \left( \bigwedge_{j \in v} [X_{s,t}]_{i,j} \right) = \sum_i \prod_{j \in v} [X_{s,t}]_{i,j} \quad (2)$$

where  $\bigwedge_{j \in v} [X_{s,t}]_{i,j}$  is the logical “and” among the values of the primitive features.

###### 2) The kernel that induces the feature mapping $\phi_v(X_{s,t})$ .

The computation of NCF  $\phi_v(X_{s,t})$  is a nonlinear mapping from the input. In the following, we prove that  $\phi_v(X_{s,t})$  is induced by a kernel and thus can benefit from efficient computation through kernel methods and kernel hierarchy.

**Lemma 1.**  $k_v(X_{s,t}, X_{s',t'}) = \phi_v(X_{s,t}) \cdot \phi_v(X_{s',t'})$  is a kernel.

*Proof.* We first prove  $k_v([X_{s,t}]_{i,v}, [X_{s',t'}]_{j,v})$  is a kernel by considering kernel’s properties. Then  $k_v(X_{s,t}, X_{s',t'})$  is proved to be a kernel using the theory of convolution kernels [13]. A detailed proof is provided in the supplementary materials [30].  $\square$

The predictive mapping  $f$  in Equation (1) can be instantiated as the linear combination of a subset of NCFs:  $f(W, \{\phi_v(X_{s,t})\}_v^{\mathcal{U}}) = \sum_v^{\mathcal{U}} W_{s,v} \phi_v(X_{s,t}) + b$ , where  $W_{s,v}$  represents the weight of the NCF  $v$  for location  $s$ . Thus, learning such a mapping function is equivalent to optimizing a subset of  $\mathcal{U}$  and their corresponding weights  $W = \{W_{s,v}\}_{s,v}^{\mathcal{U}}$ . Mathematically, this can be achieved by jointly optimizing the empirical risk term and regularization term:

$$\min_{W_{s,v}} C \sum_{s,t} \mathcal{L}(Y_{s,t}, f(W_{s,\cdot}, \{\phi_v(X_{s,t})\}_v^{\mathcal{U}})) + \sum_s \frac{1}{2} \Omega^2(W_s) \quad (3)$$

where  $\mathcal{L}(\cdot)$  is the loss function, which is convex and proper. To address the classification problem, this could be a hinge loss.  $\Omega(\cdot)$  is the regularization term that enforces sparsity so that only a few  $W_{s,v}$  will retain nonzero values to form the subset  $\mathcal{U}$ . Due to the property in Lemma 1, the efficient *representer theorem* [13] can be utilized to formulate the predictive function  $f(W, \{\phi_v(X_{s,t})\}_v^{\mathcal{U}})$  as a linear combination of hierarchical kernels. We denote  $W_{s,\cdot} = \{W_{s,v}\}_v^{\mathcal{V}}$ . The major computational challenge in solving Equation (3) comes from the large size of  $D(v)$ , which is exponential to  $|V| - |v|$ . Thanks to the favorable properties of our proposed kernel in Lemma 1, this computation can be reduced to be linear with  $|V|$ , which will be proved in Theorem 2 in Section V.

##### B. Geographical relationships of NCFs

The geographical relationships of NCFs include both geographical similarity and geographical heterogeneity.

###### 1. Geographical similarity: general conjunctive precursors.

For a domain of interest, the sparsity among NCFs for different locations can be learned jointly because they follow the same DAG relation shown in Figure 2. Specifically, the

majority of the important conjunctive features tend to be the smaller ones (i.e., those consisting of fewer primitive features), while most long conjunctive features can normally be enforced to zeros. To achieve this, we enforce  $\ell_{p,1}$ -norm ( $p \in (1, 2]$ ) on the norms of the descendants of each NCF so that longer conjunctive features will be subject to greater penalties.

$$\Omega(W) = \sum_{v \in \mathcal{V}} d_v \left( \sum_{u \in D(v)} (r_{D(v)}(W_{\cdot, u}))^p \right)^{\frac{1}{p}} \quad (4)$$

where  $d_v = a^{|v|}$ ,  $a > 0$  is the regularization parameter corresponding to NCF with a specific length, and  $p = (1, 2]$  controls the sparsity. An NCF can be selected only when all of its ancestors are selected, in which case  $p = 2$ ; otherwise, an NCF could be selected (i.e., nonzero) even if its ancestors are zeros.

2. Geographical heterogeneity: regional conjunctive precursors.

Although different locations may share similar general textual expressions, the strength of this similarity typically varies. The textual expressions within the same spatial neighborhood tends to be more similar than those far away. For example, events that occur at neighboring locations at around the same time could well involve similar topics, so the texts from neighboring locations may share a number of common keyword conjunctions that are related to the events.

To take into account such geographical heterogeneity, we propose a new hierarchical multi-task learning strategy for each NCF's norm  $r_{D(v)}(W_{\cdot, u})$  by enforcing an  $\ell_2$ -norm on each geographical neighborhood  $g$ .

$$r_{D(v)}(W_{\cdot, u}) = \sum_g^G \|\{W_{j,u}\}_j^g\|_2, \quad u \in D(v) \quad (5)$$

where  $G$  is the set of all the geographical neighborhoods. Combining Equations (4) and (5), we propose the following novel bi-space regularization term:

$$\Omega(W) = \sum_{v \in \mathcal{V}} d_v \left( \sum_{u \in D(v)} \left( \sum_g^G \|\{W_{j,u}\}_j^g\|_2 \right)^p \right)^{1/p} \quad (6)$$

where the sparsity of the NCFs are enforced by considering the hierarchical structures in two spaces. Specifically, the conjunction relation in DAG is modeled by the outer  $\ell_{p,1}$ -norm and the geographical hierarchy is modeled by the inner  $\ell_{2,1}$ -norm.

The regularization term in Equation (6) is non-smooth and multilevel. To simplify its form, the following elegant equivalent transformation is proposed and proved in Lemma 2:

**Lemma 2.** Define  $\Theta_{d,r} = \{x \in \mathbb{R}^d | x \geq 0, \sum_i x_i^r \leq 1\}$ . The regularization term  $\Omega(W)$  defined in Equation 6 can be transferred to an equivalent problem as follows:

$$\Omega(W) = \min_{\gamma, \lambda, \mu} \sum_{v \in \mathcal{V}} \frac{d_v^2}{\gamma_v} \sum_{u \in D(v)} \frac{1}{\lambda_{u,v}} \sum_g^G \frac{1}{\mu_{u,v,g}} \sum_j^g \|W_{j,u}\|^2 \quad (7)$$

where  $\gamma \in \Theta_{|\mathcal{V}|,1}$ ,  $\lambda_v = \Theta_{|D(v)|,\hat{p}}$ , and  $\mu_{u,v,g} = \Theta_{|G|,1}$ .  $\hat{p} = p/(2-p)$ .

*Proof.* The proof iteratively utilizes Holder's inequality. Further details are provided in the supplementary materials [30].  $\square$

The equivalent form is an elegant quadratic form of the weights  $W_{j,u}$ , making it possible to utilize the representer theorem to solve the empirical risk minimization problem. By introducing this equivalent form into Equation (3), we obtain:

$$\min_{\gamma, \lambda, \mu, W} C \sum_{s,t}^{S,T} \mathcal{L}(Y_{s,t}, f(W_s, \{\phi_v(X_{s,t})\}_v^{\mathcal{V}})) + \sum_{u \in \mathcal{V}} \sum_g^G \Psi_{u,g}^{-1}(\gamma, \lambda, \mu) \sum_j^g \|W_{j,u}\|^2 \quad (8)$$

where  $\Psi_{u,g}(\gamma, \lambda, \mu) = \left( \sum_{v \in A(u)} \frac{d_v^2}{\gamma_v \lambda_{v,u} \mu_{v,u,g}} \right)^{-1}$ . The above problem is convex in  $\gamma, \lambda, \mu$ , and  $W$ .  $\mathcal{L}(\cdot)$  is the hinge loss.

### C. Relationships to previous models

In this section, we show that our proposed model HTNL is the general form of several state-of-the-art models:

1. Generalization of multiple kernel learning. Let  $p = 2$ ,  $\Omega(\cdot) = \|W\|_2$ ,  $|G| = 1$ , and  $n_{s,t} \equiv 1, \forall s, t \in S, T$ . Our model in Equation (8) is reduced to multiple kernel learning [15]:

$$\min_w C \sum_i^n \mathcal{L}(y_i, f(w, x_i)) + \|w\|_2^2$$

where  $w$  is the set of feature weights,  $n$  is the number of samples,  $x_i$  and  $y_i$  are the  $i$ th input and output of the model.

2. Generalization of hierarchical kernel learning. Let  $p = 2$ ,  $|G| = 1$ , and only allow Boolean conjunctive features, i.e.,  $n_{s,t} \equiv 1, \forall s, t \in S, T$ . Our model in Equation (8) is thus reduced to hierarchical kernel learning [5]:

$$\min_{\gamma, w} C \sum_i^n L(y_i, f(w, \{\Pi_{j \in v} x_{i,j}\}_v^{\mathcal{V}})) + \sum_{u \in \mathcal{V}} \Psi_{u,1}^{-1}(\gamma, 1, 1) \|w_u\|_2^2$$

where  $x_{i,j}$  is the binary value of the  $j$ th primitive feature in the  $i$ th input.

3. Generalization of generalized hierarchical kernel learning. Let  $|G| = 1$  and only allow Boolean conjunctive features, i.e.,  $n_{s,t} \equiv 1, \forall s, t \in S, T$ . Our model in Equation (8) is thus reduced to the generalized hierarchical kernel learning [17]:

$$\min_{\gamma, w} C \sum_k^m \sum_i^{n_k} \mathcal{L}(y_{k,i}, f(w_k, \{\Pi_{j \in v} x_{k,i,j}\}_v^{\mathcal{V}})) + \sum_{u \in \mathcal{V}} \Psi_{u,1}^{-1}(\gamma, \lambda, 1) \sum_k^m \|w_{k,u}\|^2 \quad (9)$$

where  $x_{k,i,j}$  is the binary value of the  $j$ th primitive feature in the  $i$ th input of the  $k$ th task and  $w_{k,u}$  is the weight value of the  $u$ th feature of the  $k$ th task.

## V. OPTIMIZATION ALGORITHM

In this section, we propose a new efficient algorithm to solve the objective function in Equation (8). First, the dual form is proposed, simplified, and then solved by the proposed algorithm, which we have named hierarchical-multitask numerical conjunctive feature learning. Then theoretical analyses of the convergence and time complexity are presented.

---

**Algorithm 1** Hierarchical-multitask NCF Learning

---

**Require:**  $X, Y, C, G$ , and  $\mathcal{V}$ .

**Ensure:** solution  $W$  and  $b$ .

```
1: Initialize  $\mathcal{U} = S(\mathcal{V})$ ,  $\mathcal{W}_m, \Gamma, \Phi = \mathbf{0}$ .
2: repeat
3:   repeat
4:     Normalize  $\eta \leftarrow \eta / \sum_v \eta_v$ 
5:      $\eta_u(\beta) \leftarrow \left( \sum_{v \in A(u)} d_v^p \beta_v^{(1-p)} \right)^{1/(1-p)}$ 
6:     Initialize  $\xi_{g,u} = 1/|G|$ ,  $u \in \mathcal{U}$ 
7:     repeat
8:        $\alpha \leftarrow$  solve Equation (14) given  $\xi$ 
9:        $\xi \leftarrow$  solve Equation (14) given  $\alpha$ 
10:    until convergence
11:    step size  $d \leftarrow \sqrt{\log(\mathcal{U})/k} / \|\nabla H(\eta)\|_\infty$ 
12:     $\eta \leftarrow \exp \mathbf{1} + \log \eta - d \cdot \nabla H(\eta)$ 
13:  until Convergence
14:  if Equation (15) is satisfied then
15:    break
16:  else
17:    Add the nodes violating Equation (15) to  $\mathcal{U}$ 
18:  end if
19: until Forever
```

---

### A. Duality form

The primal form in Equation (8) of the objective function can be reformulated into the following duality form:

$$\begin{aligned} \min_{\gamma, \lambda, \mu} \max_{\alpha} \sum_{t \in T, s \in S} \alpha_{s,t} - \frac{1}{2} \sum_{u \in \mathcal{V}} \sum_g^G \Psi_{v,u,g}(\gamma, \lambda, \mu) h(g, u) \\ \text{s.t.} \quad \sum_{t \in T} \alpha_{s,t} \cdot Y_{s,t} = 0, \forall s \in S \\ 0 \leq \alpha_{s,t} \leq C, \forall s \in S, \forall t \in T \end{aligned} \quad (10)$$

where  $\alpha = \{\alpha_{s,t}\}_{s,t}^{S,T}$  is the dual variable.  $h(g, u) = \sum_j^G \sum_{i,k} \alpha_{j,i} Y_{j,i} \phi_u(X_{j,i}) \phi_u(X_{j,k}) Y_{j,k} \alpha_{j,k}$ . The above function is convex in  $\gamma, \mu$ , and  $\lambda$  and concave in  $\alpha$ . However, the problem as stated involves too many variables and is thus difficult to solve efficiently. To address this problem, Theorem 1 proposes a simplified equivalent formation.

**Theorem 1.** *The objective function in Equation (10) can be simplified into the following equivalent form.*

$$\min_{\beta} \max_{\alpha} \sum_{t \in T, s \in S} \alpha_{s,t} - \frac{1}{2} \left( \sum_{u \in \mathcal{V}} \eta_u(\beta) \cdot \hat{h}(u)^{\bar{p}} \right)^{1/\bar{p}} \quad (11)$$

where  $\eta_u(\beta) = \left( \sum_{v \in A(u)} d_v^p \beta_v^{(1-p)} \right)^{1/(1-p)}$ ,  $\hat{h}(u) = \max_{g \in G} h(g, u)$ , and  $\bar{p} = \hat{p}/(\hat{p} - 1)$ .

*Proof.* The proof, which is very technical, is shown in the supplementary materials [30].  $\square$

### B. Active Set Algorithm

Because the size of all the possible NCFs  $\mathcal{V}$  is exponential to the size of the primitive features  $V$ , Equation (11) can easily be computationally unfeasible to solve even with a moderate size of  $V$ . To handle this problem, the sparsity of  $\mathcal{V}$  is taken into account. This means that for the optimal solution to Equation

(11), most members of  $v \in \mathcal{V}$  should be 0. Thus, solving the original problem in Equation (11) is equivalent to solving the following subproblem where only the small subset of non-zero variables at the optimal solution of Equation (11) need to be involved. The computational effort required in the latter case will be significantly lower than in the original problem.

$$\min_{\beta} \max_{\alpha} \sum_{t \in T, s \in S} \alpha_{s,t} - \frac{1}{2} \left( \sum_{u \in \mathcal{U}} \eta_u(\beta) \cdot \hat{h}(u)^{\bar{p}} \right)^{1/\bar{p}} \quad (12)$$

where  $\mathcal{U} = \{u | W_u \neq 0, u \in \mathcal{V}\}$  is the set of nonzero-weighted NCF's. The objective function is convex in  $\beta$  but concave in  $\alpha$ .

However, the non-zero variables at the optimum are unknown beforehand. This leads us to leverage the active set algorithm [5], which efficiently updates and optimizes the set  $\mathcal{U}$  until the optimality condition is satisfied. The specific procedures are shown in Algorithm 1. The algorithm initializes  $\mathcal{U}$  as the top node in the DAG; the subproblem in Equation (12) is solved in Lines 3-13, which is elaborated in Section V-C. The solution to the subproblem is then validated against the optimality condition of the original problem via Theorem 2. If the optimality condition is satisfied, then the algorithm is terminated; Otherwise, the NCFs that violate the optimality condition will be added to the current NCF subset  $\mathcal{U}$  for the next iteration.

### C. Solution to the Subproblem

Equation (11) is simplified to a subproblem where only the NCFs with nonzero weights are retained. This permits the active set algorithm to solve the following subproblem by changing the set of nonzero NCFs until the optimality condition is satisfied.

To efficiently solve Equation (12), which is convex and Lipschitz-continuous in  $\beta$ , we employ the mirror descent algorithm [7], which achieves a near-optimal convergence rate when the feasibility set is a simplex such as in our problem.

In general, mirror descent iterations require  $\alpha$  and  $\beta$  to be solved alternately. When fixing  $\alpha$ , the gradient of the objective function in Equation (12) with respect to  $\beta$  is calculated and then  $\beta$  is updated by a descent step  $d$ .  $\beta$  is then fixed, and the updated  $\alpha$  is used to solve the following problem:

$$\max_{\alpha} \sum_{t \in T, s \in S} \alpha_{s,t} - \frac{1}{2} \left( \sum_{u \in \mathcal{U}} \eta_u(\beta) (\max_{g \in G} h(g, u))^{\bar{p}} \right)^{1/\bar{p}} \quad (13)$$

which is difficult to solve due to the “max” function inside the  $\ell_{\bar{p}}$ -norm. To remove the “max” term, an auxiliary variable  $\xi$  is introduced to transform Equation (13) to the following equivalent problem:

$$\max_{\alpha} \min_{\xi} \sum_{t \in T, s \in S} \alpha_{s,t} - \frac{1}{2} \left( \sum_{u \in \mathcal{U}} \eta_u(\beta) \sum_g \xi_{g,u} h(g, u)^{\bar{p}} \right)^{1/\bar{p}} \quad (14)$$

where  $\xi_u \in \Theta_{|G|,1}$ . Thus  $\alpha$  and  $\xi$  can be solved alternately until convergence is achieved. Specifically, when fixing  $\xi$ , solving  $\alpha$  is similar to the  $\ell_{\bar{p}}$ -norm MKL problem [20] with a



different feasibility set for the optimization variables. When  $\alpha$  is updated and fixed,  $\xi$  is easily optimized by straightforward linear programming.

### D. Theoretical Analysis

#### 1. Optimality analysis for convergence criteria.

Algorithm 1 will converge when the current candidate set of NCFs  $\mathcal{U} \subseteq \mathcal{V}$  satisfies the optimality condition. To verify this, in the following, the derivation of the sufficient condition of the optimality is proposed and proved in Theorem 2.

**Theorem 2.** Denote  $(\beta_{\mathcal{U}}, \alpha_{\mathcal{U}})$  as an  $\epsilon_{\mathcal{U}}$ -approximate optimal solution of Equation (12) based on the current active set  $\mathcal{U}$ . It is then an optimal solution for Equation (11) with a duality gap less than  $\epsilon$  if the following condition holds:

$$\begin{aligned} & \max_g \max_{u \in S(\mathcal{U})} \sum_{v \in D(u)} \frac{h(g, v)}{(\sum_{x \in A(v) \cap D(v)} d_u)^2} \\ & \leq \left( \sum_{u \in \mathcal{U}} \eta(\beta_{\mathcal{U}})(\hat{h}(u))^{\frac{1}{\bar{p}}} \right)^{\frac{1}{\bar{p}}} + 2(\epsilon - \epsilon_{\mathcal{U}}) \end{aligned} \quad (15)$$

*Proof.* The proof, which is very technical, is provided in the supplementary materials [30].  $\square$

#### 2. Time complexity analysis.

In the proposed algorithm, the most time-consuming part is the verification of a sufficient condition of convergence because it involves the search of an exponential variable space. Due to the use of the NCF lattice in Figure 2(b) and our proposed kernel, this can be reduced to a polynomial complexity, as proved by Theorem 2:

**Theorem 3.** The sufficient condition can be examined efficiently in polynomial time:  $(\sum_s^S n_s^2) \cdot |\mathcal{U}^*| \cdot e + (\sum_s^S n_s^2) \cdot |\mathcal{U}^*| \cdot e$

*Proof.* To prove this theorem, the properties of our kernel in Lemma 1 are required. The proof process is too technical and provided in the supplementary materials [30].  $\square$

The remaining computation in Algorithm 1 primarily involves the solution of the subproblem in Equation (12). Denote  $|\mathcal{U}^*|$  as the size of the final active set  $\mathcal{U}^*$ . Then Equation (12) is solved  $O(|\mathcal{U}^*|)$  times in the worst case, which requires  $\log(|\mathcal{U}^*|)$  iterations. The dominant computation in each iteration is solving Equation (13), whose conservative complexity estimate is  $O((\sum_s^S n_s^3) \cdot |\mathcal{U}^*|^2)$ , where  $n_s$  denotes the size of the data for location  $s$ . This amounts to  $O(n_{s,t}^3 \cdot S \cdot |\mathcal{U}^*|^3 \log(|\mathcal{U}^*|))$ . After combining this with the time complexity proved by Theorem 2, the overall computational complexity of the proposed algorithm is obtained:  $O(n_{s,t}^3 \cdot S \cdot |\mathcal{U}^*|^3 \log(|\mathcal{U}^*|) + (\sum_s^S n_s^2) \cdot |\mathcal{U}^*| \cdot e + (\sum_s^S n_s^2) \cdot |\mathcal{U}^*| \cdot e)$

## VI. EXPERIMENTS

In this paper, the performance of the proposed model HTNL is evaluated using 9 real datasets from two different domains. First, the datasets and experimental settings are introduced. Then, the effectiveness and efficiency of HTNL are evaluated against several existing methods that are the state-of-the-arts.

In addition, qualitative evaluations on the selection of NCFs demonstrates the interpretability of HTNL. All the experiments were conducted on a 64-bit machine with Intel(R) core(TM) quad-core processor (i7CPU@ 3.10GHz) and 16.0GB memory.

### A. Experiment Setup

Table II: Datasets and Labels

Dataset	#Tweets	Label sources <sup>1</sup>	#Events
Argentina	160,564,890	Clarín; La Nación; Infobae	1427
Chile	97,781,414	La Tercera; Las Últimas Noticias; El Mercurio	776
Colombia	158,332,002	El Espectador; El Tiempo; El Colombiano	1287
El Salvador	21,992,962	El Diario de Hoy; La Prensa Gráfica; El Mundo	730
Mexico	197,550,208	La Jornada; Reforma; Milenio	5907
Paraguay	30,891,602	ABC Color; Última Hora; La Nación	2114
Uruguay	10,310,514	El País; El Observador	664
Venezuela	167,411,358	El Universal; El Nacional; Últimas Noticias	3320
U.S.	6,487,623,208	CDC Flu Activity Map	533

1) *Datasets and Labels:* For the datasets on Latin America, the raw data was obtained by randomly sampling 10% (by volume) of the Twitter data from Jan 2013 to Dec 2014 in 8 countries as shown in Table II. The Twitter data for the period from Jan 1, 2013 to Dec 31, 2013 was used for training, while the data for the second half of the period, from Jan 1, 2014 to Dec 31, 2014, was used for the performance evaluation. The event forecasting results were validated against a labeled events set, known as the gold standard report (GSR), exclusively provided by MITRE [24], as shown in Table II. An example of a labeled GSR event was given by the tuple: (CITY="Hermosillo", STATE = "Sonora", COUNTRY = "Mexico", DATE = "2013-01-20"). For the datasets in the United States, the raw data was crawled from Jan 2013 to Dec 2014, as shown in Table II. As in the first dataset, the Twitter data for the period from Jan 1, 2013 to Dec 31, 2013 was used for training while the second half of the period, from Jan 1, 2014 to Dec 31, 2014, was used for the performance evaluation. The forecasting results for the flu outbreaks were validated against the corresponding influenza statistics reported by the Centers for Disease Control and Prevention (CDC) [33]. CDC publishes the weekly influenza-like illness (ILI) activity level within each state in the United States based on the proportion of outpatient visits to healthcare providers for ILI. There are 4 ILI activity levels: minimal, low, moderate, and high, where the level "high" corresponds to a salient flu outbreak and was considered for forecasting. An example of a CDC flu outbreak event is: (STATE = "Virginia", COUNTRY = "United States", WEEK = "01-06-2013 to 01-12-2013").

2) *Parameter Settings and Metrics:* In this experiment, different sets of primitive features were defined for domains of civil unrest and influenza outbreaks, respectively. For the civil

<sup>1</sup>In addition to the top 3 domestic news outlets, the following news outlets are included: The New York Times; The Guardian; The Wall Street Journal; The Washington Post; The International Herald Tribune; The Times of London; Infolatam.

Table III: Evaluation results of all methods in effectiveness and efficiency on 9 datasets

Method	Prediction Performance Area Under the Curve (AUC) of ROC									Runtime (second)
	Argentina	Chile	Colombia	El Salvador	Mexico	Paraguay	Uruguay	Venezuela	Influenza	
LASSO	0.5738	0.5202	0.6441	0.6201	0.6295	0.6013	0.6526	0.5722	0.6738	40
LogReg	0.7268	0.7323	<b>0.7384</b>	0.7315	<b>0.7310</b>	0.7044	0.7274	0.6792	0.4851	18
KDE-LDA	0.7665	<b>0.7647</b>	0.6919	<b>0.7723</b>	0.6454	0.6654	0.7279	0.7214	0.2827	656
MREF	0.7264	0.7625	0.5296	0.5714	0.5613	0.6171	0.6812	0.5887	0.4969	444
TMTL	0.7069	0.7140	0.5633	0.5157	0.5720	0.6129	0.6931	0.6586	0.4989	203
RuleFit	0.7246	0.5300	0.5101	0.7221	0.5648	0.5008	0.5698	0.7080	0.6100	<b>3</b>
gHKL	0.6850	0.6597	0.5198	0.6189	0.5368	0.6067	0.6878	0.6970	0.5000	76
HTNL	<b>0.8264</b>	0.7311	<b>0.7384</b>	0.7386	0.6659	<b>0.7374</b>	<b>0.7538</b>	<b>0.7508</b>	<b>0.6951</b>	132

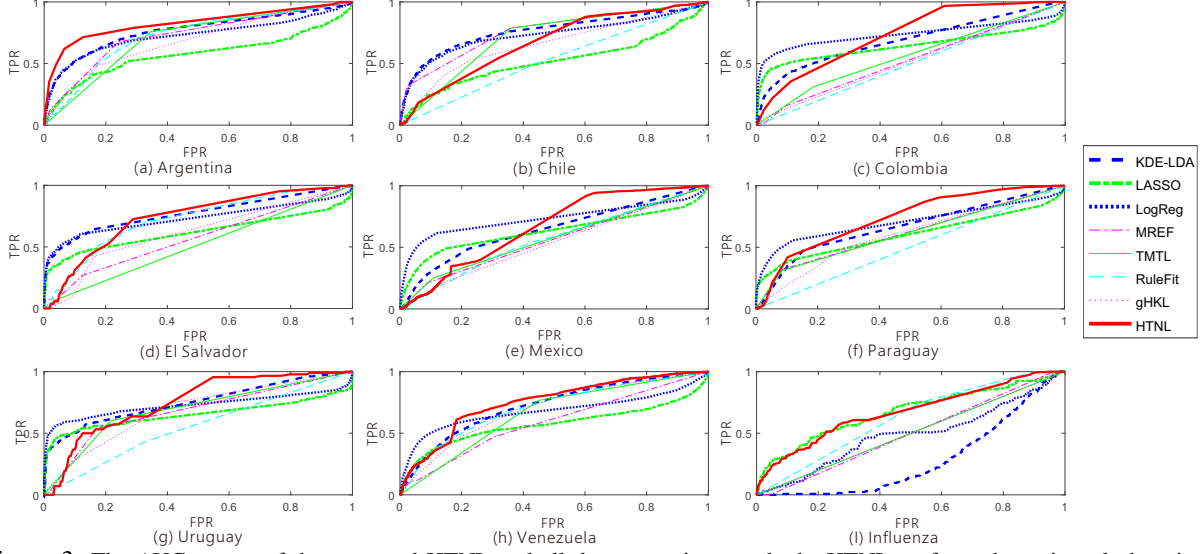


Figure 3: The AUC curves of the proposed HTNL and all the comparison methods. HTNL performed consistently best in general.

unrest domain, the feature set included over  $10^{30}$  conjunctive features which were all the possible conjunctions of 100 civil unrest related words (such as “protest” and “riot”) and hashtags (such as “#yosoy132”) based on the keyword list in [25]. For the influenza outbreaks, the feature set consisted of over  $10^{54}$  features generated from the combinations of 181 influenza-related words extracted based on the keywords list used in [21]. In the experiment, Twitter data collection was partitioned into a sequence of date-interval subcollections. The event forecasting task was to utilize one day tweet data to predict whether or not there would be an event in the next day for a specific city (for the civil unrest domain), or a specific state (for the influenza outbreaks domain), which means the lead time  $q = 1$ . To perform this task, we created a training set and a test set for each city (or state), where each data sample was the daily tweet observation with the above-mentioned features. The predicted events were structured as tuples of (date, city/state). A predicted event was matched to a real event if both the date and location attributes were matched. To validate the prediction performance, the Area Under the Curve (AUC) of Receiver operating characteristic (ROC) curve were adopted. ROC curve illustrates the performance of a binary classifier as its discrimination threshold is varied. The curve is created by plotting the true positive rate (TPR) against the false positive rate (FPR) at various threshold settings. The AUC measures the area below this curve, which is a well-recognized

metric to reflect the comprehensive performance of a classifier. There are several parameters for our proposed model HTNL. First,  $p$  (with three optional values:  $\{1.1, 1.5, 1.9\}$ ) and  $C$  (with optional values:  $\{0.01, 0.1, 1, 10, 100\}$ ) were determined with a 3-fold cross validation. The parameters were set as  $d_v = 2^{|v|}$  suggested by Jawanpuria et al. [18]. The geographical hierarchy was “state-city” administrative relation for civil unrest datasets while “HHSregion<sup>1</sup>-state” relation for influenza dataset.

### B. Performance

In this section, the proposed HTNL is evaluated quantitatively and qualitatively. First, the effectiveness and efficiency of our HTNL are compared with 7 state-of-the-art methods on spatial event forecasting and predictive rule learning, including *Logistic regression (LogReg)* [9], *LASSO* [25], *Kernel density estimation-based logistic regression (KDE-LDA)* [34], *Tree-guided Group Lasso for Multi-task Learning (TMTL)* [19], *Multi-resolution Spatial Event Forecasting (MREF)* [38], *RuleFit* [12], and generalized hierarchical kernel learning (gHKL) [18]. Their parameter settings are described in our supplementary materials [30]. In addition, the illustration of the selected NCFs by the proposed HTNL is also presented.

1) *AUC on different datasets*: Table III summarizes the effectiveness and efficiency of the proposed HTNL on different

<sup>1</sup>HHSregions: <http://www.hhs.gov/about/agencies/regional-offices/>



Table IV: The NCF precursors (translated in English) discovered for different datasets. (# NCF: The number of selected NCFs; avg. len.: The average length of all the NCFs; The symbol “+” denotes the logical “and” within an NCF)

	Argentina	Colombia	El Salvador	Mexico	Uruguay	Venezuela	Influenza
Top 10 NCFs (length $\geq 2$ )	government+congress government+deputies to+end+hate let's+fight followers+free to+know+results protest+against national+triumph hate+hunger hate+class	government+water water+problem violence+protest national+mayor national+control mayor+control national+government national+water national+freedom power+death	national+prices work+sponsor followers+people security+violence university+freedom government+money security+justice justice+violence government+project government+means	mayor+protests power+to+perform village+help justice+march resources+opportunities national+notice security+protest reform+day reform+protest rights+report	know+hope security+rights project+president power+death power+matches president+hope death+matches project+hope national+university national+fight	government+high violence+order national+support candidate+march national+students patria+control national+government fight+policy students+protest government+violence	bed+flu+home bed+home cold+sick you+flu is+epidemic flu+bed sick+stomach yousick have+today not+well
# NCFs	131	141	108	517	71	325	2017
avg. len.	1.2290	1.2837	1.0648	1.8046	1.0000	1.6892	1.9499

datasets. The AUC measure was adopted to quantify the performance. First, the results shown in Table III demonstrates that the methods that take into account the spatial information, especially the geographical hierarchy, performed better. Specifically, KDE-LDA, MREF,TMTL, and the proposed HTNL typically performed the best in most situations. KDE-LDA performed much better on civil unrest datasets than the influenza dataset. This might be because this method was specially designed to forecast crimes, which are small-scale social events unlike influenza epidemics in “state” level. LogReg also achieved a very competitive performance with AUC larger than 0.73 on three datasets. Second, HTNL outperformed all the other methods in 6 out of the 9 datasets, and achieved the second best in most of the remaining datasets. This is because HTNL not only considers the geo-hierarchy, but more importantly, is to consider the NCFs like the frequencies of keyword co-occurrences as new features that capture crucial precursors for future events. In contrast, the Boolean rule learning methods including RuleFit and gHKL only achieved the AUCs around 0.6 on the datasets, generally worse than the other methods. This is because they can only consider the binary occurrence of keywords on each date instead of the frequencies of keywords. Thus they lose much information of the magnitude of the social indicators. Among all the datasets, the overall performance for Argentina and Chile was among the best while the Influenza outbreaks forecasting was a relatively difficult prediction task with lower AUCs for most of the methods.

2) *Efficiency on running time:* The rightmost column of Table III shows the training time efficiency comparison among HTNL and the competing methods for forecasting influenza outbreaks. The efficiency evaluation results on civil unrest datasets followed a similar pattern and are not provided due to space limitations. The running times on the test set for all the comparison methods were instant (i.e., less than 0.01 second for one prediction) so that are not provided here, either. Table III shows that RuleFit required smallest amount of time of only 3 seconds, because of two reasons 1) it binarizes the numerical frequencies into Boolean values as inputs; and 2) it utilizes an efficient heuristic procedure to obtain a suboptimal solution. Simpler methods like LASSO and LogReg also achieved high efficiency with less than 50 seconds. In addition, even though HTNL need optimize a problem with exponentially large

size of numerical feature set, it still achieved highly efficient computation. This is because of the utilization of the good property of the proposed kernel in Lemma 1, which is proved to reduce the exponential time complexity down to polynomial as proved in Theorem 3.

3) *Event forecasting performance on ROC curves:* In Figure 3, the event forecasting performance in ROC curves for 9 datasets is illustrated. For all these datasets, the proposed HTNL performed consistently among the best whose curves were farthest away from the point (1,0). Specifically, For the the civil unrest datasets like “Argentina”, “Paraguay”, “Uruguay”, and “Venezuela”, HTNL generally performed the best, with ROC curves covering the largest areas above the x-axis. For the datasets of Colombia, El Salvador, and Mexico, HTNL, KDE-LDA, and LogReg were among the best, where HTNL typically performed the best when FPR was larger than 0.5. This merit is important because it indicates that HTNL tends to provide the most extensive true positive alarms among all the methods. Comprehensive detections of sensitive social events are important to many applications such as social emergency management. For the influenza dataset, according to Figure 3(i), HTNL consistently outperformed the other methods with different FPR and TPR values. LogReg and LASSO also achieved quite competitive performance.

4) *Qualitative Evaluation:* Another advantage of our HTNL is its strong interpretability compared to most of the spatial event forecasting models that merely use primitive features. Table IV shows the results on the selection of NCFs by the proposed HTNL for 7 datasets. The results for “Colombia” and “Paraguay” datasets are moved to the supplementary material [30] due to space limitations. Specifically, the top 10 NCFs with length larger than 1, namely precursor rules, are listed. The amount and average length of the selected NCFs are also presented. The original Spanish words were translated into English by Google Translator [31]. The symbol “+” denotes the logical “and” within an NCF. An NCF will be triggered only if all its words connected by “+” co-occur in a tweet.

According to Table IV, the proposed HTNL effectively selected high-quality NCFs robustly for all the datasets in two different domains, namely civil unrest and influenza outbreaks. For civil unrest datasets, the NCFs in Table IV typically represent the motivations or propaganda of the protest events. For example, the high frequency of the NCF “water”+“problem” could probably be one important reason that causes social

unrest in Colombia, while the NCFs like “justice”+“march” and “reform”+“protest” could be the triggers for those future events in Mexico. The NCFs can also be propaganda-related, such as “to”+“end”+“hate” and “let’s”+“fight”. In contrast to civil unrest datasets, those top NCFs for influenza dataset were typically not about the motivation or advertisement of organized social events, but the symptoms or discussions about flu. For example, NCF like “bed”+“flu”+“home” appearing together in a tweet was likely to be a strong indicator for a person’s disease status. “sick”+“stomach” could also be a symptom of stomach flu. Additionally, the numbers of total NCFs optimally selected for different datasets show that for larger countries, the numbers and average lengths of NCFs tend to be larger. This is because larger size of population typically leads to more various social issues and thus the social events could be indicated by more diverse precursors.

## VII. CONCLUSIONS

Forecasting spatial societal events in social media is significant. It is also very challenging because the precursors of the future events are not straightforward and can be sophisticated ensemble of underlying rules. Most existing methods simplified this problem by considering frequencies of keywords or predefined phrases as features due to the challenges such as the inherent exponential complexity of ensemble rule learning, distant supervision, numerical values, geographical relations. To jointly handle all the challenges with theoretical guarantee, we propose a novel spatial event forecasting model named HTNL which learns the NCFs efficiently. an efficient algorithm is proposed to optimize the model parameters and prove its theoretical guarantees for error bound and time efficiency. Extensive experiments on multiple datasets demonstrate the effectiveness and efficiency of the proposed method. Moreover, qualitative analysis on the extracted NCFs explicitly shows the strong interpretability of HTNL.

## REFERENCES

- [1] A. Acharya, R. J. Mooney, and J. Ghosh. Active multitask learning using supervised and shared latent topics. *Pattern Recognition and Big Data*, page 75, 2016.
- [2] H. Achrekar, A. Gandhe, R. Lazarus, S.-H. Yu, and B. Liu. Predicting flu trends using Twitter data. In *IEEE Conference on Computer Communications Workshops*, pages 702–707, 2011.
- [3] B. Ahmed, T. Thesen, K. Blackmon, R. Kuzniecky, O. Devinsky, J. Dy, and C. Brodley. Multi-task learning with weak class labels: Leveraging ieeg to detect cortical lesions in cryptogenic epilepsy. In *Machine Learning for Healthcare Conference*, pages 115–133, 2016.
- [4] M. Arias, A. Arratia, and R. Xuriguera. Forecasting with Twitter data. *ACM Trans. on Intelligent Systems and Technology (TIST)*, 5(1):8, 2013.
- [5] F. Bach. High-dimensional non-linear variable selection through hierarchical kernel learning. *arXiv preprint arXiv:0909.0844*, 2009.
- [6] J. Bollen, H. Mao, and X. Zeng. Twitter mood predicts the stock market. *Journal of Computational Science*, 2(1):1–8, 2011.
- [7] S. Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [8] H. Cheng, X. Yan, J. Han, and C.-W. Hsu. Discriminative frequent pattern analysis for effective classification. In *ICDE 2007*, pages 716–725. IEEE, 2007.
- [9] R. Compton, C. Lee, J. Xu, L. Artieda-Moncada, T.-C. Lu, L. De Silva, and M. Macy. Using publicly visible social media to build detailed forecasts of civil unrest. *Security informatics*, 3(1):4, 2014.

- [10] G. Cong, K.-L. Tan, A. K. Tung, and X. Xu. Mining top-k covering rule groups for gene expression data. In *SIGMOD 2005*, pages 670–681. ACM, 2005.
- [11] X. Dong, D. Mavroudis, F. Calabrese, and P. Frossard. Multiscale event detection in social media. *Data Mining and Knowledge Discovery*, 29(5):1374–1405, 2015.
- [12] J. H. Friedman and B. E. Popescu. Predictive learning via rule ensembles. *The Annals of Applied Statistics*, pages 916–954, 2008.
- [13] T. Gärtner, P. A. Flach, A. Kowalczyk, and A. J. Smola. Multi-instance kernels. In *ICML 2002*, volume 2, pages 179–186, 2002.
- [14] M. S. Gerber. Predicting crime using Twitter and kernel density estimation. *Decision Support Systems*, 61:115–125, 2014.
- [15] M. Gönen and E. Alpaydm. Multiple kernel learning algorithms. *Journal of Machine Learning Research*, 12(Jul):2211–2268, 2011.
- [16] J. He, W. Shen, P. Divakaruni, L. Wynter, and R. Lawrence. Improving traffic prediction with tweet semantics. In *IJCAI 2013*, pages 1387–1393, 2013.
- [17] P. Jawanpuria, S. N. Jagarlapudi, and G. Ramakrishnan. Efficient rule ensemble learning using hierarchical kernels. In *ICML 2011*, pages 161–168, 2011.
- [18] P. Jawanpuria, J. S. Nath, and G. Ramakrishnan. Generalized hierarchical kernel learning. *The Journal of Machine Learning Research*, 16(1):617–652, 2015.
- [19] S. Kim and E. P. Xing. Tree-guided group lasso for multi-task regression with structured sparsity. In *ICML 2010*, pages 543–550, 2010.
- [20] M. Kloft, U. Brefeld, S. Sonnenburg, and A. Zien. Lp-norm multiple kernel learning. *Journal of Machine Learning Research*, 12(Mar):953–997, 2011.
- [21] A. Lamb, M. J. Paul, and M. Dredze. Separating fact from fear: Tracking flu infections on Twitter. In *HLT-NAACL*, pages 789–795, 2013.
- [22] K. Lin and J. Zhou. Interactive multi-task relationship learning. In *ICDM 2016*, pages 241–250, 2016.
- [23] D. M. Malioutov and K. R. Varshney. Exact rule learning via boolean compressed sensing. In *ICML 2013*, pages 765–773, 2013.
- [24] MITRE. <http://www.mitre.org/>.
- [25] N. Ramakrishnan, P. Butler, et al. Beating the news with EMBERS: forecasting civil unrest using open source indicators. In *KDD 2014*, pages 1799–1808. ACM, 2014.
- [26] U. Rückert and S. Kramer. A statistical approach to rule learning. In *ICML 2006*, pages 785–792. ACM, 2006.
- [27] T. Sakaki, M. Okazaki, and Y. Matsuo. Earthquake shakes Twitter users: real-time event detection by social sensors. In *WWW 2010*, pages 851–860, 2010.
- [28] E. Schubert, M. Weiler, and H.-P. Kriegel. Signitrend: scalable detection of emerging topics in textual streams by hashed significance thresholds. In *KDD 2014*, pages 871–880. ACM, 2014.
- [29] G. J. Simon, V. Kumar, and P. W. Li. A simple statistical model and association rule filtering for classification. In *KDD 2011*, pages 823–831. ACM, 2011.
- [30] Supplementary Materials. [https://www.dropbox.com/s/y79z5shz9tjods8/bare\\_conf.pdf?dl=0](https://www.dropbox.com/s/y79z5shz9tjods8/bare_conf.pdf?dl=0).
- [31] G. Translate. <https://translate.google.com/>.
- [32] A. Tumasjan, T. O. Sprenger, P. G. Sandner, and I. M. Welp. Predicting elections with Twitter: What 140 characters reveal about political sentiment. *ICWSM 2010*, 10:178–185, 2010.
- [33] C. F. View. <https://gis.cdc.gov/grasp/fluview/main.html>.
- [34] X. Wang, M. S. Gerber, and D. E. Brown. Automatic crime prediction using events extracted from Twitter posts. In *Social Computing, Behavioral-Cultural Modeling and Prediction*, pages 231–238. Springer, 2012.
- [35] N. J. Yadwadkar, B. Hariharan, J. E. Gonzalez, and R. Katz. Multi-task learning for straggler avoiding predictive job scheduling. *Journal of Machine Learning Research*, 17(106):1–37, 2016.
- [36] L. Zhao, F. Chen, J. Dai, T. Hua, C.-T. Lu, and N. Ramakrishnan. Un-supervised spatial event detection in targeted domains with applications to civil unrest modeling. *PLoS one*, 9(10):e110206, 2014.
- [37] L. Zhao, F. Chen, C.-T. Lu, and N. Ramakrishnan. Spatiotemporal event forecasting in social media. In *SDM 2015*, pages 963–971. SIAM, 2015.
- [38] L. Zhao, F. Chen, C.-T. Lu, and N. Ramakrishnan. Multi-resolution spatial event forecasting in social media. In *ICDM 2016*, pages 689–698, 2016.
- [39] J. Zhou, J. Chen, and J. Ye. Malsar: Multi-task learning via structural regularization. Arizona State University, 2011.