



HAL
open science

Opinion mining from Twitter data using evolutionary multinomial mixture models

Md Abul Hasnat, Julien Velcin, Stephane Bonnevey, Julien Jacques

► **To cite this version:**

Md Abul Hasnat, Julien Velcin, Stephane Bonnevey, Julien Jacques. Opinion mining from Twitter data using evolutionary multinomial mixture models. 2015. hal-01204613v1

HAL Id: hal-01204613

<https://inria.hal.science/hal-01204613v1>

Preprint submitted on 24 Sep 2015 (v1), last revised 27 Feb 2019 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

OPINION MINING FROM TWITTER DATA USING EVOLUTIONARY MULTINOMIAL MIXTURE MODELS

*Laboratoire ERIC, Université de Lyon - Lumière **

BY MD. ABUL HASNAT*, JULIEN VELCIN*, STEPHANE
BONNEVAY* AND JULIEN JACQUES*

Image of an entity can be defined as a structured and dynamic representation which can be extracted from the opinions of a group of users or population. Automatic extraction of such an image has certain importance in political science and sociology related studies, e.g., when an extended inquiry from large-scale data is required. We study the images of two politically significant entities of France. These images are constructed by analyzing the opinions collected from a well known social media called *Twitter*. Our goal is to build a system which can be used to automatically extract the image of entities over time.

In this paper, we propose a novel evolutionary clustering method based on the parametric link among Multinomial mixture models. First we propose the formulation of a generalized model that establishes parametric links among the Multinomial distributions. Afterward, we follow a model-based clustering approach to explore different parametric sub-models and select the best model. For the experiments, first we use synthetic temporal data. Next, we apply the method to analyze the annotated social media data. Results show that the proposed method is better than the state-of-the-art based on the common evaluation metrics. Additionally, our method can provide interpretation about the temporal evolution of the clusters.

1. Introduction. We define an *image* as a multi-faceted representation that aggregates a set of opinions or general impressions regarding an entity. By entity, we mean a politician, a celebrity, a company, a brand, etc. In this research, we are particularly interested to use annotated social media data to extract the image of two French politicians and observe its changes/evolution over time. We consider the annotated data from the *ImagiWeb* project (Velcin et al., 2014) which are extracted before and after the 2012 French presidential election. The annotation provides a compact and meaningful representation for each tweet. Our goal is to develop a temporal/evolutionary clustering technique, which groups the annotated opinions and then extracts the image of an entity over time from the clustering results. Subsequently, we want to explain/interpret the temporal changes of the image created from each group of users.

In the recent years, the social media plays a significant role in many aspects of our daily activity. There exist numerous popular social media such as Twitter or Facebook, where the users (people) often provide their opinions about particular entity, e.g., persons (politician, actor), products consumed in the daily life, etc. A common method to analyze such data is to use a clustering method that naturally groups the users/opinions, and then investigate each group independently. An important property of these data is that they may change *over time* due to changes of the attributes, and appearance/disappearance of users. Moreover, users may change their opinion about the targeted entity.

An ordinary clustering method is unlikely to adapt with such temporal dynamics of the data, as it does not consider any relevant information such as history and temporal effects. The notion of evolutionary clustering (Chakrabarti, Kumar and Tomkins, 2006; Xu, Kliger and Hero Iii, 2014; Chi et al., 2009; Xu et al., 2012) appears in such situations, where the method should be specialized in clustering temporal data by taking care of the historic information and current data altogether. Numerous methods exist, which address these issues appropriately and cluster temporal data. These methods are based on different strategies, such as spectral clustering (Chi et al., 2009; Xu, Kliger and Hero Iii, 2014) and probabilistic generative model (Blei and Lafferty, 2006; Xu et al., 2012; Kim et al., 2015). However, it remains an important issue - how to interpret the evolution of the clusters. In this research, we are motivated by this issue and propose a novel method based on the Multinomial mixture model (Bishop et al., 2006) to cluster the temporal data as well as interpret the evolution of the clusters through some prior belief. Therefore, we propose a novel method which simultaneously performs evolutionary clustering and interpreting the evolution.

Multinomial Mixture (MM) model based clustering strategy is a popular method for clustering discrete data (Meilă and Heckerman, 2001; Silvestre, Cardoso and Figueiredo, 2014; Hasnat, Alata and Trémeau, 2015; Agresti, 2002). Most recently, it has been exploited to perform evolutionary clustering (Kim et al., 2015). In this research, we consider MM as the core model for the data and propose an evolutionary clustering method by deriving appropriate link between the parameters of MM at different time.

Parametric link among probability distributions has been used in the context of transfer learning (Biernacki, Beninel and Bretagnolle, 2002; Jacques and Biernacki, 2010; Beninel et al., 2012), where the goal is to adapt a clustering model from a source population to a target one. In the context of continuous features, Biernacki, Beninel and Bretagnolle (2002) proposed

a parametric link between the Normal distributions. Jacques and Biernacki (2010) extended it for the binary features using Bernoulli distribution. However, no such formulation exists for the Multinomial distribution. Moreover, such parametric link-based methods are never considered in the context of evolutionary clustering. This research addresses both of these issues.

This research proposes a novel evolutionary clustering method for extracting *image* of political entities. The highlights of our contributions include: (a) propose a formulation for a parametric link among Multinomial distributions; (b) develop a novel evolutionary clustering method by exploiting the link parameters and (c) provide interpretation of the link parameters to interpret cluster evolutions. First, we use synthetic data to evaluate and compare the proposed method w.r.t. the state-of-the-art methods. Next, we apply it to analyze the temporal dynamics of social media data obtained from the *ImagiWeb* project (Velcin et al., 2014). Results in Sec. 4 show that the proposed method is better than the state-of-the-art methods.

In the rest of the paper, we present the data in Sec. 2, describe our proposed method in Sec. 3, present the experimental results in Sec. 4, provide analysis of the political data in Sec. 5, and finally draw conclusions in Sec. 6.

2. The Imagiweb project and the political opinion dataset. We collected data from the *political opinion dataset* of the ImagiWeb¹ (IW-POD) project, see Velcin et al. (2014) for further details of data collection, relevant statistics and representation. IW-POD consists of manually annotated tweets, from May 2012 to January 2013, related to two French politicians: Francois Hollande (FH) and Nicolas Sarkozy (NS). First, these tweets are annotated into 11 different aspects, such as Attribute (Att), Person (Per), Entity (Ent), Skills (Skl), Political line (Pol), Balance (Bal), Injunction (Inj), Projet (Pro), Ethic (Eth), Communication (Com) and No aspect detected (N/A). Afterward, each aspect is annotated with 6 opinion polarities, such as very negative (-2), negative (-1), no polarity (0), Null, positive (+1) and very positive (+2). For example, the tweet - *Sarko is more rational (orig: Sarko est plus rationnel)* is annotated with the aspect called *Person* and polarity +1. It is about NS and indicates that the user provides positive opinion with an emphasis on the personal attribute. Another example, the tweet - *Nicolas Sarkozy, the worst president of the Fifth Republic (Orig: Nicolas Sarkozy, le plus mauvais président de la Vème République)* is annotated with the aspect called *Skill* and polarity -1. It is a negative opinion about NS and indicates that the user emphasizes on the skill of NS.

¹<http://mediamining.univ-lyon2.fr/velcin/imagiweb/dataset.html>

In order to use these tweets for clustering, they are regrouped within the specified time epoch. Moreover, similar polarities are merged, e.g., two positives (+1 and +2) are merged into one as only positive (+). Therefore, each aspect consists of four polarities, such as positive (+), negative (-), zero (0) and undefined/null (\emptyset). As a consequence, finally each regrouped tweet represents the opinion of an user about a particular politician which is a 44 (11×4) dimensional vector of discrete data. In our experiment, we group opinions from IW-POD into three time² epochs: t_1 , t_2 and t_3 , see Table 1 for details of the temporal data. Moreover, since the true number of clusters is unknown, we run clustering for different numbers of clusters ranging from 3 to 9.

TABLE 1

Details of the IW-POD dataset which is divided into three time periods. Each observation consists of a 44 dimensional discrete valued vector that encodes information about 11 different aspects each having 4 polarities.

Time stamp	Time period	Significance	Num. opinions N. Sarkozy	Num. opinions F. Hollande
t1	03/12 - 06/12	Before and After Election	1018	1168
t2	07/12 - 10/12	After Election	1067	1079
t3	11/12 - 01/13	After Election	1079	708

3. Parametric Link Based Evolutionary Clustering. We adopt the parametric link approach (Biernacki, Beninel and Bretagnolle, 2002; Jacques and Biernacki, 2010) for evolutionary clustering by assuming that the source samples are equivalent to the samples at time epoch t and target samples represent sample of time $t+1$. With this assumption, we incorporate linear link between Multinomials at different time epoch. The algorithm for the proposed clustering method is presented in Algorithm 1.

3.1. *Related work.* Evolutionary Clustering (ECL), also called *clustering over time*, aims to cluster the data that dynamically evolves over time (Chakrabarti, Kumar and Tomkins, 2006). Ordinary clustering methods are not appropriate as they group/partition the data samples only based on the certain properties of the data. In contrary, ECL methods cluster the data by additionally considering the temporal smoothness to reflect the long-term trends of the data while being robust to the short-term variations (Chakrabarti, Kumar and Tomkins, 2006; Xu, Kliger and Hero Iii,

²The first round of the presidential election was held in 22/04/2012 and the second round run-off was held on 06/05/2012. Therefore, the data collected during this election period belong to time epoch t_1 .

2014; Chi et al., 2009). ECL should maintain four properties (Chakrabarti, Kumar and Tomkins, 2006) such as consistency, noise removal, smoothing and cluster correspondence. The demand and application of such clustering method are increasing rapidly due to the significant growth of the dynamic data in numerous domains. It has been successfully applied to analyze news (Xu et al., 2012), social media (Kim et al., 2015), stock price (Xu, Kliger and Hero Iii, 2014), photo-tag pairs (Chakrabarti, Kumar and Tomkins, 2006), and documents (Blei and Lafferty, 2006).

Temporal/evolutionary data clustering has been addressed from several viewpoints in the literature, which naturally raises several task-specific notions about ECL. A distinction among them can be as follows: (1) clustering (2) monitoring and (3) interpreting. In the following paragraphs, we review relevant literature based on this distinction.

Following the definition of Chakrabarti, Kumar and Tomkins (2006), the ECL method clusters data by considering the historic information and current data. Based on this definition, in this research we do not consider the methods which do not take into account the historic information. Besides, in order to limit our focus on the parametric methods, we do not consider the methods from non-parametric Bayesian based approaches (Xu et al., 2008; Dubey et al., 2013; Kharratzadeh, Renard and Coates, 2015).

Numerous methods based on different techniques have been proposed in the literature (Chakrabarti, Kumar and Tomkins, 2006; Xu, Kliger and Hero Iii, 2014; Chi et al., 2009; Xu et al., 2012; Kim et al., 2015; Blei and Lafferty, 2006). Chakrabarti, Kumar and Tomkins (2006) provided a generic framework for this problem and proposed evolutionary version of k-means and hierarchical agglomerative clustering methods. Their proposed framework is based on optimizing a global cost function that consists of snapshot (static clustering) quality and history cost (temporal smoothness). This is considered as the first work for evolutionary clustering and has been subsequently extended by other researchers. Chi et al. (2009) proposed two evolutionary clustering methods based on spectral clustering strategy. In their approach, they added terms within the clustering cost functions in order to regularize the temporal smoothness. Xu, Kliger and Hero Iii (2014) recently proposed AFFECT, which performs adaptive evolutionary clustering by estimating an optimal smoothing parameter. This approach is extended with several static clustering methods, such as k-means, hierarchical and spectral. A common property of these methods is that they specialized for continuous data and hence may not be an appropriate choice for clustering categorical data that is our concern in this research.

Dynamic Topic Model (DTM) is a well-known probabilistic method for

analyzing temporal categorical data (Blei and Lafferty, 2006). It was originally developed to analyze time evolution of topics in large document collections. DTM extends the popular topic modeling method called Latent Dirichlet Allocation (LDA) (Blei, Ng and Jordan, 2003). It uses Dirichlet prior based smoothing, which sometime over-smooth the data. As a consequence, it may cluster the data samples with non co-occurring features in the same group (Kim et al., 2015). This eventually causes DTM to underperform to cluster some classical non-textual temporal categorical data. Recently, Kim et al. (2015) address this issue and proposed a probabilistic generative model based evolutionary clustering method, called Temporal Multinomial Mixture (TMM). TMM extends the classical Multinomial Mixture (MM) model by incorporating temporal dependency into the relation between data components of current time epoch and the clusters of the previous time epoch. MM is a well-known standard probabilistic model, which has been widely used to cluster static discrete/categorical data (Meilă and Heckerman, 2001; Silvestre, Cardoso and Figueiredo, 2014). Similar to MM, TMM estimates model parameters using an Expectation Maximization (EM) algorithm. Although both DTM and TMM provide reasonable results to cluster temporal categorical data, they are unable to detect and provide any interpretation of the cluster evolutions, which is one of the main foci of this research. Indeed, TMM is more related to our proposed approach as we aim to establish parametric link among MMs at different time epochs.

The evolution monitoring task (Spiliopoulou et al., 2006; Oliveira and Gama, 2010; Ferlez et al., 2008; Lamirel, 2012) tracks the evolution of clusters by identifying the birth, death, split, merge and survival of clusters at different time. An external clustering method is first used at each time to cluster the data, e.g., Spiliopoulou et al. (2006) and Oliveira and Gama (2010) used the k-means method, whereas Lamirel (2012) used the neural clustering method. Afterward, the association and mapping among the clusters at different time is examined based on several heuristics. For example, Oliveira and Gama (2010) used cluster centroid related statistics, called comprehensive representation of clusters. This approach is very similar to the notion of detecting recurrent concept drifts in a semi-supervised context, see Li, Wu and Hu (2012) for an example. A different method, called label-based diachronic approach (Lamirel, 2012), exploits the MultiView Data Analysis paradigm among the cluster labels at different time. In this approach, each feature is analyzed individually to compute recall, precision and F-measure. These information are used to construct heuristics for monitoring evolution. Our approach is different than the above methods, because: (a) we do not aim to propose a cluster monitoring method explicitly and (b) we do not use

a static clustering method. Besides the above methods, [Ferlez et al. \(2008\)](#) proposed a joint clustering-monitoring method which uses the cross association algorithm to cluster data and a bipartite graph to monitor evolution. For data clustering, they group the distinct features (word) in each cluster and hence features do not coexist in different clusters. This is different than us as we exploit all the features in order to provide a feature level interpretation for the evolution.

The task of evolution interpretation aims to explain the reason for the evolution of clusters at different time. It can be accomplished by explicitly analyzing the features. To this aim, [Lamirel \(2012\)](#) used the F-measures from individual features of the matched clusters (of different time) and construct a similarity report. In our work, this interpretation can be directly obtained from the link parameters by applying threshold on the link parameters values. Therefore, our method is different from [Lamirel \(2012\)](#) as the link parameters computation is an integral part of the clustering task.

Based on the above distinctions from several viewpoints (clustering, monitoring and interpretation), we find that our method is more similar to the evolutionary clustering methods rather than the evolution monitoring methods. Therefore, we compare our method only with the relevant state-of-the-art evolutionary clustering methods, such as [Xu, Kliger and Hero Iii \(2014\)](#), [Blei and Lafferty \(2006\)](#) and [Kim et al. \(2015\)](#).

Now we focus on the literature related to our proposal. The idea of parametric link in a transfer learning context ([Beninel et al., 2012](#)) is inherited from the concept for Generalized Discriminant Analysis (GDA) ([Biernacki, Beninel and Bretagnolle, 2002](#)). GDA adapts the classification rule from a source population to a target population through a linear link map of their descriptive parameters. This is different than standard discriminant rules which assumes a similarity between the source and target populations. [Biernacki, Beninel and Bretagnolle \(2002\)](#) proposed several models with associated estimated parameters for GDA within the context of multivariate Gaussian distribution. Later, [Jacques and Biernacki \(2010\)](#) extends the work of [Biernacki, Beninel and Bretagnolle \(2002\)](#) for binary data using Bernoulli distribution ([Bishop et al., 2006](#)). We observe that these approaches can be exploited for developing an evolutionary clustering method by replacing the notion of source/target with different time epochs $t - 1/t$. Besides, such development requires the derivation of the linear link for the Multinomial distribution. Afterward, the link parameters naturally allow us to interpret the evolution of the clusters at different time.

Categorical data/observations consists of the responses from a certain number of categories. Different types (nominal and ordinal) of categorical

data are observed in numerous studies (Agresti, 2002), such as social science, biomedical science, genetics, education and marketing. Moreover, data from different tasks, such as text retrieval and visual object classification, are often converted to the categorical form. For example, text data can be converted to this form by considering the unique words of the vocabulary as an independent category/term and then each sentence/paragraph/document is represented as a discrete count vector (Zhong and Ghosh, 2005). The Multinomial distribution is a standard probability distribution for modeling and analyzing the discrete categorical data (Agresti, 2002).

The Multinomial Mixture (MM) is a statistical model based on the Multinomial distribution. It has been used for cluster analysis with discrete data (Meilă and Heckerman, 2001; Agresti, 2002; Zhong and Ghosh, 2005; Silvestre, Cardoso and Figueiredo, 2014; Hasnat et al., 2015). Meilă and Heckerman (2001) studied several Model-Based Clustering (MBC) methods with MM and experimentally compared them using different criteria such as clustering accuracy, computation time and number of selected clusters. Silvestre, Cardoso and Figueiredo (2014) proposed a MBC method for MM which integrates both model estimation and selection task within a single EM algorithm. In their work, they extended the MBC strategy previously proposed by Figueiredo and Jain (2002) and provided a formulation to compute the Minimum Message Length (MML) criterion for model selection. Most recently, Hasnat et al. (2015) proposed a MBC method which performs simultaneous clustering and model selection using the MM. Their strategy performs similar task as Silvestre, Cardoso and Figueiredo (2014) in a computationally efficient manner which has been previously proposed for the Gaussian distribution (Garcia and Nielsen, 2010) and Fisher distribution (Hasnat, Alata and Trémeau, 2015). Moreover, similar to Meilă and Heckerman (2001), they provided a comparison among different model initialization and selection strategies. Following all of the above approaches (Meilă and Heckerman, 2001; Silvestre, Cardoso and Figueiredo, 2014; Hasnat et al., 2015), in this research we exploit the MBC framework to cluster discrete data with MM.

MBC (Fraley and Raftery, 2002; Melnykov and Maitra, 2010) is a well-established method for cluster analysis and unsupervised learning. It assumes a probabilistic model (e.g., mixture model) for the data, estimates the model parameters by optimizing an objective function (e.g., model likelihood) and produces probabilistic clustering. The Expectation Maximization (EM) (McLachlan and Krishnan, 2008) is mostly used in MBC to estimate the model parameters. EM consists of an Expectation step (E-step) and a Maximization step (M-step) which are iteratively employed to maximize the

log likelihood of the data.

Initialization of the EM algorithm has significant impact on clustering results (McLachlan and Krishnan, 2008; Baudry and Celeux, 2015). The EM algorithm is sensitive to its initialization, because with different initializations it may converge to different values of likelihood function, some of which can be local maxima (i.e., sub-optimal results). In order to overcome this, numerous different initialization strategies are proposed and experimented in the relevant literature (Biernacki, Celeux and Govaert, 2003; Meilă and Heckerman, 2001; Baudry and Celeux, 2015; Hasnat et al., 2015). Following recommendations, we use the small-EM (Biernacki, Celeux and Govaert, 2003; Biernacki et al., 2006; Baudry and Celeux, 2015; Hasnat et al., 2015) method to initialize the MM parameters.

MBC has been commonly exploited to identify the best model for the data by fitting a set of models with different parameterizations and/or number of components and then applying a statistical model selection criterion (Fraley and Raftery, 2002; Biernacki, Celeux and Govaert, 2000; Figueiredo and Jain, 2002; Melnykov and Maitra, 2010; Hasnat, Alata and Trémeau, 2015). In this paper, we apply this model fitting and selection strategy for two purposes: (a) to identify the parametric submodels (Section 3.4) and (b) to automatically select the number of components (Section 3.7).

3.2. Statistical model for evolutionary data samples. Let S^t be a set of samples corresponding to time t and S^{t+1} be a set from the next time $t + 1$. We assume that while the cluster labels for S^t are known to us (estimated from $t - 1$), labels of S^{t+1} are unknown.

Let S^t be composed of N^t pairs $(\mathbf{x}_1^t, \mathbf{z}_1^t), \dots, (\mathbf{x}_{N^t}^t, \mathbf{z}_{N^t}^t)$ where $\mathbf{x}_i^t = \{x_{i,1}^t, \dots, x_{i,D}^t\}$ is the D dimensional count vector of order V , i.e., $\sum_{d=1}^D x_{i,d}^t = V$ and \mathbf{z}_i^t is the associated class label such that $\mathbf{z}_{i,k}^t = 1$ if the data belongs to cluster k with $k = 1, \dots, K$ and $\mathbf{z}_{i,k}^t = 0$ otherwise. We assume that any sample \mathbf{x}_i^t of S^t is an independent realization of the random variable \mathbf{X}^t of distribution:

$$\mathbf{X}^t \sim \mathcal{M}(V, \boldsymbol{\mu}_k^t), \quad k = 1, \dots, K$$

with $\mathcal{M}(V, \boldsymbol{\mu}_k^t)$ is the V -order Multinomial distribution with parameter $\boldsymbol{\mu}_k^t = (\mu_{k,1}^t, \dots, \mu_{k,D}^t)$ which is formally defined as (Bishop et al., 2006):

$$(3.1) \quad \mathcal{M}(\mathbf{x}_i|V, \boldsymbol{\mu}_k) = \binom{V}{x_{i,1}, x_{i,2}, \dots, x_{i,D}} \prod_{d=1}^D \mu_{k,d}^{x_{i,d}}$$

here, $\boldsymbol{\mu}_k$ is the parameter of the Multinomial distribution of k^{th} class with $0 \leq \mu_{k,d} \leq 1$ and $\sum_{d=1}^D \mu_{k,d} = 1$. Therefore, samples of the entire set S^t

can be modeled with a mixture of k Multinomials, also called Multinomial Mixture (MM) model, which has the following form:

$$(3.2) \quad f(\mathbf{x}_i | \Theta_K) = \sum_{k=1}^K \pi_k \mathcal{M}(\mathbf{x}_i | V, \boldsymbol{\mu}_k)$$

In Eq. (3.2), $\Theta_K = \{(\pi_1, \boldsymbol{\mu}_1), \dots, (\pi_K, \boldsymbol{\mu}_K)\}$ is the set of model parameters, π_k is the mixing proportion with $\sum_{k=1}^K \pi_k = 1$ and $\mathcal{M}(\mathbf{x}_i | V, \boldsymbol{\mu}_k)$ is the density function (Eq. (3.1)). Besides, we assume that the class label \mathbf{z}_i^t is an independent realization of a random vector \mathbf{Z}^t , distributed according to 1-order Multinomial:

$$\mathbf{Z}^t \sim \mathcal{M}(\mathbf{1}, \boldsymbol{\pi}^t)$$

where $\boldsymbol{\pi}^t = \pi_1^t, \dots, \pi_K^t$ is the mixing proportion of the model in Eq. (3.2).

The assumption of MM is similar for the samples of S^{t+1} with random variable \mathbf{X}^{t+1} and parameter $\boldsymbol{\mu}_k^{t+1}$. However, for S^{t+1} the labels \mathbf{z}_i^{t+1} of N^{t+1} pairs $(\mathbf{x}_1^{t+1}, \mathbf{z}_1^{t+1}), \dots, (\mathbf{x}_{N^{t+1}}^{t+1}, \mathbf{z}_{N^{t+1}}^{t+1})$ are unknown. In the context of evolutionary clustering, our goal is to estimate the unknown labels \mathbf{z}_i^{t+1} for $i = 1, \dots, N^{t+1}$ using the information from S^t and S^{t+1} by establishing a link between $\boldsymbol{\mu}_k^t$ and $\boldsymbol{\mu}_k^{t+1}$.

3.3. Parametric link/relationship among temporal data. For random variables Y^t and Y^{t+1} distributed according to the Gaussian distribution, a linear distributional link exists (under weak assumptions) (Biernacki, Beninel and Bretagnolle, 2002), which has the form: $Y^{t+1} \sim DY^t + b$, where D and b are the link parameters among the samples of different time epoch. For binary data the following distributional linear link among Bernoulli parameters (α^{t+1} and α^t with $0 \leq \alpha \leq 1$) is derived by Jacques and Biernacki (2010):

$$(3.3) \quad \alpha^{t+1} = \Phi(\delta \Phi^{-1}(\alpha^t) + \lambda \gamma)$$

where $\delta \in \mathbb{R}^+ \setminus \{0\}$, $\lambda \in \{-1, 1\}$ and $\gamma \in \mathbb{R}$ are the link parameters. Φ is the cumulative Gaussian function of mean 0 and variance 1, see Fig. 3.1. We can use the above formulation for Multinomial parameters by considering two issues: (1) Multinomial parameter $\boldsymbol{\mu}_k$ has equivalent property as $\boldsymbol{\alpha}_k$ except $\sum_{d=1}^D \mu_{k,d} = 1$ and (2) samples from X are not necessary to be binary, which makes λ useless. Considering these issues we can derive parametric link between $\boldsymbol{\mu}_t$ and $\boldsymbol{\mu}_{t+1}$ as:

$$(3.4) \quad \mu_{k,d}^{t+1} = \frac{\Phi(\delta_{k,d} \Phi^{-1}(\mu_{k,d}^t) + \gamma_{k,d})}{\sum_{r=1}^D \Phi(\delta_{k,r} \Phi^{-1}(\mu_{k,r}^t) + \gamma_{k,r})}$$

where $\delta_{k,d} \in \mathbb{R}^+ \setminus \{0\}$ and $\gamma_{k,d} \in \mathbb{R}$ are the link parameters. In Eq. (3.4), the combination of parameters $\delta_{k,d}$ and $\gamma_{k,d}$ for $\forall k, d$ is called a full model which is over-parameterized and may leads to ambiguity. Instead, we consider several sub-models with certain constraints on the parameters, see the following section.

3.4. Parametric sub-models. The idea of defining sub-models is frequent in Model-Based Clustering (MBC) (Fraley and Raftery, 2002). We fit the evolutionary clustering model (Eq. (3.4)) with different sub-models and then select the best model using the Bayesian Information Criteria (Schwarz et al., 1978):

$$(3.5) \quad BIC = -2L(\Theta) + \nu \log(N^{t+1})$$

where $L(\Theta)$ is the log-likelihood (Eq. (3.6)) value associated to the MM parameters of $t + 1$, ν is the number of free parameters of the sub-model. These sub-models provide sufficient interpretation about the change in parameters from time t to $t + 1$. Definition and interpretation of several basic sub-models, defined as pair $(\delta_{k,d}/\gamma_{k,d})$ are given below:

(M1) 1/0: This model is constrained with $\delta_{k,d} = 1$ and $\gamma_{k,d} = 0$ for $\forall k, d$, i.e., $\nu = 0$. It indicates that the observations X^{t+1} can be modeled with $\mu_{k,d}^t$ and hence no evolution occurred.

(M2) 0/ $\gamma_{k,d}$: This model is constrained with $\delta_{k,d} = 0$ for $\forall k, d$, i.e., $\nu = K * D$. It indicates that the observations X^{t+1} should be modeled without considering $\mu_{k,d}^t$. This model should be selected when a new cluster evolved independently and does not consider any historical information. This is the most general model that can certainly fit the observations X^{t+1} to a MM most efficiently subject to a good initialization of the alternative iterative method. Several possible variations³ of this model are: $0/\gamma$, $0/\gamma_k$ and $0/\gamma_d$.

(M3) $\delta_{k,d}/0$: This model is constrained with $\gamma_{k,d} = 0$ for $\forall k, d$, i.e., $\nu = K * D$. It indicates that $\mu_{k,d}^{t+1}$ are evolved through $\mu_{k,d}^t$ in a specific transformation space (inversed cumulative Gaussian). This model should be selected when true evolution occurred which can be explained in detail through certain belief on observed features and obtained clusters. Moreover, such a model can be plugged in with any other method in order to describe the cluster evolution. Several possible variations of this model are: $\delta/0$, $\delta_k/0$ and $\delta_d/0$. This model is equivalent to the fundamental unconstrained model assumed by Biernacki, Beninel and Bretagnolle (2002).

³Subscript k means cluster dependent and d means feature dependent. No subscription means a constant value for all clusters and features.

(M4) $1/\gamma_{k,d}$: In this model, $\delta_{k,d} = 1$ for $\forall k, d$, i.e., $\nu = K * D$. This model does nearly similar task as model M3. It is relatively easier to fit through the additive term in the inverse cumulative Gaussian space. On the other hand, it is less expressive in terms of interpretation. Several possible variations of this model are: $1/\gamma$, $1/\gamma_k$ and $1/\gamma_d$.

3.5. Parameter estimation. In our proposed formulation of evolutionary clustering, we estimate two different types of parameters (see Eq. (3.4)): (1) MM model parameters: μ and π and (2) temporal link parameters: δ and γ . We estimate them in two steps. The first step consists of estimating μ and π (only for $t = 1$) for the observed samples of time t . In the second step, we estimate δ and γ . At any time epoch, we estimate the class labels \mathbf{z}_i by *maximum a posteriori*.

3.5.1. Multinomial Mixture (MM) Parameters. At time $t = 1$, we estimate the MM parameters using an Expectation Maximization (EM) algorithm that maximizes the log-likelihood value which has the following form:

$$(3.6) \quad L(\Theta) = \sum_{i=1}^N \log \sum_{j=1}^K \pi_j \mathcal{M}(\mathbf{x}_i | \boldsymbol{\mu}_j)$$

where $N = N^1$ is the number of samples. In the Expectation step (E-step), we compute posterior probability as:

$$(3.7) \quad \rho_{i,k} = p(z_{i,k} = 1 | \mathbf{x}_i) = \frac{\pi_k \prod_{d=1}^D \mu_{k,d}^{x_{i,d}}}{\sum_{l=1}^K \pi_l \prod_{d=1}^D \mu_{l,d}^{x_{i,d}}}$$

In the Maximization step (M-step), we update π_k and $\mu_{k,d}$ as:

$$(3.8) \quad \pi_k = \frac{1}{N} \sum_{i=1}^N \rho_{i,k} \quad \text{and} \quad \mu_{k,d} = \frac{\sum_{i=1}^N \rho_{i,k} \mathbf{x}_{i,d}}{\sum_{i=1}^N \sum_{r=1}^D \rho_{i,k} \mathbf{x}_{i,r}}$$

The E and M steps are iteratively employed until certain convergence criterion (difference of the log-likelihood values of successive iterations) is satisfied. The estimation of $\mu_{k,d}$ using Eq. (3.8) is only applicable for $t = 1$ due to the unavailability of any temporal information. For any time $t + 1$, when the link parameters are available, $\mu_{k,d}$ is estimated with Eq. (3.4).

3.5.2. Link parameters. Estimation of link parameters $\delta_{k,d}$ and $\gamma_{k,d}$ uses $\mu_{k,d}^t$ and the observed samples at time $t + 1$. Similar to [Jacques and Biernacki](#)

(2010), we use again an EM algorithm, but in which the M step is not explicit. Consequently, we employ an external optimization method such as an alternative iterative algorithm which consists of a succession, componentwise of the simplex method⁴ (Nelder and Mead, 1965). In general, the starting point of the alternative algorithm corresponds to the case when $\mu_{k,d}^{t+1} = \mu_{k,d}^t$, i.e., $\delta_{k,d} = 1$ and $\gamma_{k,d} = 0$. However, in order to obtain a better estimate and save computation time⁵, we apply an efficient approach, see Section 3.6.2.

Algorithm 1: Algorithm for clustering using parametric link among multinomial mixtures (PLMM).

Input: $\chi = \{S^t\}_{t=1,\dots,T}$, $S^t = \{\mathbf{x}_i\}_{i=1,\dots,N^t}$, $\mathbf{x}_i = \{x_{i,d}\}_{d=1,\dots,D}$, $x_{i,d} \in \mathbb{N}$
Output: Evolutionary clustering of χ with K classes and link parameters: $\delta_{k,d}^t$ and $\gamma_{k,d}^t \forall k, d, t$.

```

foreach  $t$  do
  if  $t = 1$  then
    | Initialize  $\pi_{j,k}$  and  $\mu_{j,k}$  for  $1 \leq j \leq k$  using the small-EM procedure, see
    | Section 3.6.1;
  end
  while not converged do
    | {Perform the E-step of EM};
    foreach  $i$  and  $j$  do
    | | Compute  $\rho_{ik} = p(z_{i,k} = 1 | \mathbf{x}_i)$  using Eq. (3.7)
    end
    | {Perform the M-step of EM};
    for  $k = 1$  to  $K$  do
    | if  $t = 1$  then
    | | Update  $\pi_k$  and  $\mu_k$  using Eq. (3.8)
    | else
    | | Update  $\pi_k$  using Eq. (3.8)
    | | Compute  $\delta_{k,d}$  and  $\gamma_{k,d}$ , see Sec. 3.5.2
    | | Update  $\mu_k$  using Eq. (3.4)
    | end
    end
  end
end

```

3.6. *Parameters initialization.* In the proposed clustering method (Algorithm 1), we need to initialize both the MM parameters $\Theta_K^{init} = \{(\pi_1^{init}, \mu_1^{init}), \dots, (\pi_K^{init}, \mu_K^{init})\}$ for time t_1 and the link parameters (δ and γ).

⁴For the implementation, we used *neldermead* function of *nloptr* R package (Ypma, 2014). The lower and upper bounds were set to -2.5 and $+2.5$ respectively only for the $\gamma_{k,d}$ parameters.

⁵The simplex method requires a large number of iterations to converge.

3.6.1. *Multinomial Mixture (MM) Parameters.* Generally, the MM parameters are initialized randomly (Meilă and Heckerman, 2001; Hasnat et al., 2015). However, with both synthetic and real data it has been demonstrated by Hasnat et al. (2015) that, random initialization has its limitation w.r.t. the clustering performance and stability. Therefore, following Hasnat et al. (2015), we initialize the model parameters using the small-EM procedure. This small-EM procedure consists of running multiple short runs of randomly initialized EM and then selecting the one with the maximum likelihood value. Here, short run means that the EM procedure does not need to wait until convergence and it can be stopped when a certain number of iterations is completed.

3.6.2. *Link parameters.* We propose an initialization procedure based on the predictive parameters set for next time epoch $\Theta_K^{pred} = \{(\pi_1^{pred}, \boldsymbol{\mu}_1^{pred}), \dots, (\pi_K^{pred}, \boldsymbol{\mu}_K^{pred})\}$. Let $\Theta_K^t = \{(\pi_1^t, \boldsymbol{\mu}_1^t), \dots, (\pi_K^t, \boldsymbol{\mu}_K^t)\}$ is the set of parameters for the current time (t) epoch. Our initialization procedure consists of the following steps:

- Step 1: estimate Θ_K^{pred} using data samples of next time X^{t+1} and an EM algorithm which is initialized with Θ_K^t .
- Step 2: compute $\delta_{k,d}^{init}$ and $\gamma_{k,d}^{init}$ for each k and d as:

$$(3.9) \quad \gamma_{k,d}^{init} = \Phi^{-1}(\mu_{k,d}^{pred}) \quad \text{for model M2}$$

$$(3.10) \quad \delta_{k,d}^{init} = \frac{\Phi^{-1}(\mu_{k,d}^{pred})}{\Phi^{-1}(\mu_{k,d}^t)} \quad \text{for model M3}$$

$$(3.11) \quad \gamma_{k,d}^{init} = \Phi^{-1}(\mu_{k,d}^{pred}) - \Phi^{-1}(\mu_{k,d}^t) \quad \text{for model M4}$$

The Eq. (3.9), (3.10) and (3.11) are simply derived from Eq. (3.4) with the consideration that denominator is equal to 1, i.e., $\sum_{d=1}^D \mu_{k,d} = 1$ for $k = 1, \dots, K$.

3.7. *Varying number of clusters.* The methodology presented in the previous sub-sections assumes the same number of clusters K for each time epoch. In this sub section, we propose an extension of it such that the method can handle varying K at different time, i.e., K_t and K_{t+1} may be different. To this aim, we modify the links initialization strategy (Section

3.6.2) in order to adapt the variability among $\Theta_{K_t}^t$ and $\Theta_{K_{t+1}}^{t+1}$. At time epoch t , this extended method requires additional information, such as: (a) number of clusters K_{t+1} and (b) cluster mapping between $\Theta_{K_t}^t$ and $\Theta_{K_{t+1}}^{t+1}$.

We adopted the method proposed by Hasnat et al. (2015) with L-method (Salvador and Chan, 2004) to select the number of cluster automatically at each time epoch. In order to initialize the link parameters, first we select the number of clusters K_{t+1} and obtain the predictive parameter set $\Theta_{K_{t+1}}^{pred}$. Next, for each cluster k in $\Theta_{K_{t+1}}^{pred}$ we find the corresponding cluster in $\Theta_{K_t}^t$ based on the minimum symmetric kullback leibler divergence (sKLD). sKLD among two clusters a and b is defined as (Hasnat et al., 2015):

$$(3.12) \quad sKLD = \frac{D_{KL}(\boldsymbol{\mu}_a, \boldsymbol{\mu}_b) + D_{KL}(\boldsymbol{\mu}_b, \boldsymbol{\mu}_a)}{2}, \text{ where}$$

$$D_{KL}(\boldsymbol{\mu}_a, \boldsymbol{\mu}_b) = \sum_{d=1}^D \mu_{a,d} \ln \left(\frac{\mu_{a,d}}{\mu_{b,d}} \right)$$

After establishing the correspondences, we use Eq. (3.9), (3.10) and (3.11) to set the initial values of the link parameters. Finally, we estimate the link parameters following Section 3.5.2.

3.8. *Interpretation of cluster evolution.* The link parameters ($\delta_{k,d}$ and $\gamma_{k,d}$) along with the function Φ are the key to interpret the cluster evolution. Let us notice some basic interpretation of the values of these parameters for all feature d and cluster k :

- $\delta_{k,d} = 0$ means that $\mu_{k,d}$ (probability) at $t + 1$ does not depend on t , whereas $\delta_{k,d} = 1$ (with $\gamma_{k,d} = 0$) means identical probability at two different times.
- $\delta_{k,d} \rightarrow 0$ and/or $\gamma_{k,d} \rightarrow \infty$ means that the distribution *tends to uniform* distribution.
- $\delta_{k,d} \rightarrow \infty$ and/or $\gamma_{k,d} \rightarrow -\infty$ means that the distribution tends to be *more concentrated* (Dirac distribution) at time $t + 1$ in the feature which has the highest probability at time t .

In order to get further interpretation, we need to understand the Multinomial parameters $\mu_{k,d}$ and the space spanned by the cumulative Gaussian Φ and its inverse Φ^{-1} . Let us consider an experiment of drawing V balls of $d = 1, \dots, D$ different colors (represent features). After each draw, the color of the ball is recorded in a D dimensional count vector \mathbf{x}_i and the ball is replaced. Therefore, at the end of i^{th} experiment $\mathbf{x}_{i,d}$ reveals the count of drawing the d^{th} colored ball. When a Multinomial distribution is used

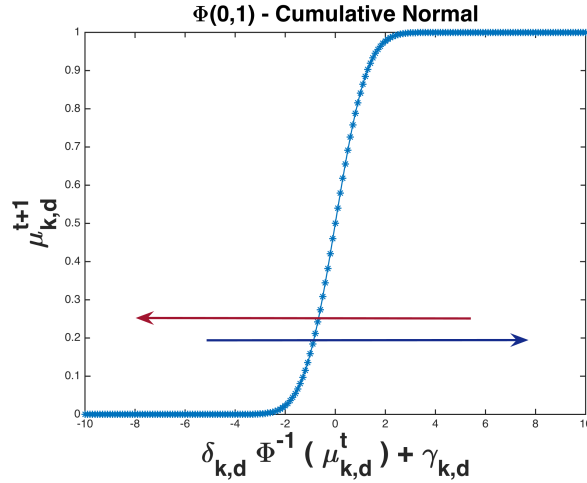


FIG 3.1. Illustrations of Cumulative Gaussian function and its relationship with the parameter change of Multinomial distribution using Eq. (3.4). The arrows indicates the direction of changes in the inverse function space which eventually increase/decrease the probability.

to fit such experimental data, its parameter $\mu_{k,d}$ reveals the probability of drawing the d^{th} colored ball.

Now, let us consider Φ in Fig. 3.1 where the values along the Y-axis represent the possible values of $\mu_{k,d}^{t+1}$ (with $0 \leq \mu_{k,d}^{t+1} \leq 1$) and the X-axis represents the values of $\mu_{k,d}^t$ after transforming through Φ^{-1} function. Now, according to Eq. (3.4), cluster evolutions ($\mu_{k,d}^t \rightarrow \mu_{k,d}^{t+1}$) can be explained through multiplication (using $\delta_{k,d}$) and addition/subtraction (using $\gamma_{k,d}$) operations.

The values of $\gamma_{k,d}$ can certainly indicates the increase/decrease of the probability of certain feature (color) subject to the selection of sub-model **M4**. On the other hand if sub-model **M3** is selected, values of $\delta_{k,d}$ can explain the belief that $\mu_{k,d}^{t+1}$ should decrease if $\mu_{k,d}^t < 0.5$ and increase if $\mu_{k,d}^t > 0.5$. For example, let us consider that in a 2 colors (red and green) ball experiment the probability of the red color ball is changed from 0.8 (at time t1) to 0.7 (at time t2). Such a change can be explained with model **M3** with $\delta_{k,red} = 0.6$, which indicates that the belief is decreased at the next time. From the above discussions it is evident that the proposed method is capable to interpret the cluster evolutions up to the feature level.

4. Numerical experiments. We begin the experiments using simulated evolutionary data samples and evaluate w.r.t. the state-of-the-art methods. A characteristic comparison of different methods is presented in

Table 2. For the simulated samples; we use the Adjusted Rand Index (ARI) (Hubert and Arabie, 1985) as a measure for evaluation. Next, we experiment and compare methods using real data. We use one of the real datasets experimented by Kim et al. (2015). We choose the *political opinion dataset* from the ImagiWeb project (Velcin et al., 2014) as it consists of data from an interesting time period - during and after the election.

TABLE 2

Characteristic comparison of different state-of-the-art evolutionary clustering methods: Parametric Link among Multinomial Mixtures (PLMM, our proposed method), Temporal Multinomial Mixture (TMM) (Kim et al., 2015), Dynamic Topic Model (DTM) (Blei and Lafferty, 2006) and adaptive evolutionary clustering method (AFFECT) (Xu, Kliger and Hero Iii, 2014).

	PLMM	DTM	TMM	AFFECT
Data Type	Discrete	Discrete	Discrete	Continuous
Interpret Evolution	Yes	No	No	No

4.1. *Simulated Data Samples.* Following standard sampling methods we generate different sets $\{S^t\}_{t=1,\dots,T}$ of simulated data for different time epochs. We draw a finite set of categorical samples (discrete count vectors) $S^t = \{\mathbf{x}_i\}_{i,\dots,N^t}$ with different numbers (10, 20 and 40) of features (dimensions) D . These samples are issued from Multinomial Mixture (MM) models of $K = 3$ classes. We consider two different sets of samples:

- Samples with higher order of categorical count (**hos**) with $V \sim 1.5 * D$ with 3 time epochs each having different number of i.i.d. samples: $N^1 = 500$, $N^2 = 100$, and $N^3 = 200$. We also add noisy counts with these samples. These type of samples provides better resemblance with the MM parameters due to sufficient number of count in the observations. Practically, this is similar to the fact when the observations consists of data over longer period of time.
- Samples with lower order of categorical count (**los**) with $V \sim 0.7 * D$ with 5 time epochs each having different number of i.i.d. samples: $N^1 = 50$, $N^2 = 40$, $N^3 = 40$, $N^4 = 30$ and $N^5 = 20$. This type of samples are sparse and often difficult to distinguish among clusters. Practically, this is similar to the fact when the observations consists of data over shorter period of time.

The evolutionary data generation process consists of two steps: (1) determine MM parameters $\mu_{k,d}$ at each time epoch $t = 1, \dots, T$ and (2) sample observations from the specified MM following assumption specified by Blei, Ng and Jordan (2003). For $t = 1$, we sample $\mu_{k,d}$ from a Dirichlet distribution and verify (separation w.r.t. the other clusters parameters (Silvestre,

(Cardoso and Figueiredo, 2014)) it using the symmetric Kullback-Leibler Divergence value. For $t > 1$, we sample $\mu_{k,d}$ from $\mu_{k,d}^{t-1}$ using the MM link relationship defined in Eq. (3.4). This ensures that we maintain the temporal smoothness property (Chakrabarti, Kumar and Tomkins, 2006; Xu, Klinger and Hero Iii, 2014) of the evolutionary data samples. In order to use the link relationship, we use only model M4 for *hos* data samples and randomly select a model among M1, M3 and M4 for *los* data samples. Next, we set the associated link parameters ($\delta_{k,d}$ and $\gamma_{k,d}$) randomly within a pre-specified range of values.

To sample observations, first we choose the order V_k of each cluster. Our sampling procedure for each observation i at each time t follows the steps below:

- Choose a cluster $z_{i,k} = 1$ as: $\mathbf{z}_i \sim \mathcal{M}(1, \pi_1, \dots, \pi_D)$, with, $\pi_d = \frac{1}{k}$.
- Choose the order τ_i of Multinomial for the sample \mathbf{x}_i using Poisson distribution as: $\tau_i \sim \text{Poisson}(V_{z_i})$.
- Draw sample \mathbf{x}_i using Multinomial distribution as: $\mathbf{x}_i \sim \mathcal{M}(\tau_i, \mu_{k,1}, \dots, \mu_{k,D})$.

TABLE 3

Simulated data evaluation and comparison using Adjusted Rand Index (ARI) (Hubert and Arabie, 1985). Methods: PLMM (proposed), Dynamic Topic Model (DTM), Temporal Multinomial Mixture (TMM) and AFFECT with k-means. Datasets consist of different types (hos and los) of samples with different numbers (10, 20 and 40) of features. hos: higher order samples and los: lower order samples. Boldfaced indicate the best result and underlined numbers indicate second best. Values inside the parentheses provide the standard deviation of the ARI values.

	PLMM	TMM	DTM	AFFECT
10, hos	0.91 (0.07)	<u>0.86</u> (0.11)	0.79 (0.14)	0.43 (0.12)
10, los	<u>0.81</u> (0.19)	0.91 (0.1)	<u>0.81</u> (0.1)	0.34 (0.11)
20, hos	0.96 (0.05)	<u>0.91</u> (0.1)	0.81 (0.18)	0.37 (0.11)
20, los	0.90 (0.18)	0.98 (0.04)	<u>0.95</u> (0.11)	0.35 (0.09)
40, hos	0.97 (0.05)	<u>0.92</u> (0.11)	0.48 (0.4)	0.33 (0.11)
40, los	<u>0.93</u> (0.16)	0.97 (0.05)	0.97 (0.1)	0.36 (0.1)

We applied our proposed Parametric Link among Multinomial Mixtures (PLMM, Algorithm 1) clustering method on these simulated data using the basic sub-models defined in Sec. 3.4. Table 3 provides the results using the ARI (Hubert and Arabie, 1985) measure. Moreover, it provides a comparative evaluation w.r.t. other state-of-the-art methods (see comparison in Table 2): (a) Temporal Multinomial Mixture (TMM) (Kim et al., 2015) with smoothness parameter $\alpha = 1$; (b) Dynamic Topic Model (DTM) (Blei and Lafferty, 2006) with hyper-parameter $\alpha = 0.01$ and (c) Adaptive evolu-

tionary clustering method (AFFECT⁶) (Xu, Kliger and Hero Iii, 2014) with k-means and Euclidean distance as a measure of similarity. We compute the average ARI of time $t = 2, \dots, T$ (at $t = 1$ there is no evolution). Results in Table 3 w.r.t. ARI evaluation shows that:

- PLMM (proposed) provides highest ARI for the *hos* samples and TMM (Kim et al., 2015) provides highest ARI for the *los* samples. These results are not surprising as both PLMM and TMM methods are specialized methods to cluster samples which are drawn from Multinomial distributions.
- DTM (Blei and Lafferty, 2006) provides better results for *los* samples and higher dimensional data. This type of data is more likely to extract from text documents for which DTM was originally proposed.
- AFFECT (Xu, Kliger and Hero Iii, 2014) performs poorly compares to others for both types of sample. This is expected because of the similarity measure used in AFFECT is appropriate for continuous data.

Next, we test statistical hypothesis among PLMM, TMM and DTM using *two sample t-test* at the 5% significance level. The null hypothesis is that - the data in two results comes from independent random samples from normal distributions with equal means and equal but unknown variances. Results show that for all *hos* data the hypothesis is rejected with $p\text{-value} < 0.001$. On the other hand, for the *los* data it is rejected only for 10 dimensional samples among the pairs (PLMM, TMM) and (DTM, TMM) with $p\text{-value} < 0.0001$.

Next, we analyze the evolution of the clusters in terms of selected sub-models. Table 4 provides the rate of different selected models. We see that, for the *hos* data samples the model M4 ($1/\gamma_{k,d}$) is mostly selected. On the other hand, for the *los* data samples, different models M1: ($1/0$), M4: ($1/\gamma_{k,d}$) and M3: ($\delta_{k,d}/0$) are selected at certain rate. This observation confirms that PLMM successfully recovers the cluster evolutions with different models which were used to generate the simulated data. Interestingly, we observe that the model M2 ($1/\gamma_{k,d}$) is not selected which reflects the true fact that it was not considered to generate the simulated data samples. Now based on the selected model, we can provide further interpretation using $\delta_{k,d}$ and $\gamma_{k,d}$, see Sec. 3.4.

Finally, we conduct experiments with varying number of clusters K at different time epoch. For this experiment, we use the same MM parameters which were used to generate the *hos* data samples. To ensure different K at different epoch, we randomly select a pair of time epochs and remove a

⁶We experimented AFFECT with hierarchical and spectral clustering also. However, k-means provided the best results.

TABLE 4

Percentage of the selected models for the interpretation of evaluation. **hos**: higher order (categorical count) samples and **los**: lower order samples. **Boldfaced** indicate the highest rate.

	M1: (1/0)	M4: ($1/\gamma_{k,d}$)	M3: ($\delta_{k,d}/0$)	M2: ($0/\gamma_{k,d}$)
10, hos	0 %	94 %	6 %	0 %
10, los	15 %	38 %	47 %	0 %
20, hos	0 %	92 %	8 %	0 %
20, los	14 %	43 %	43 %	0 %
40, hos	0 %	96 %	4 %	0 %
40, los	4 %	37 %	59 %	0 %

cluster from one of them. Then, we generate $N^t = N^{t+1} = 1000$ synthetic data samples from them using the same procedure mentioned before. Applying the extension of PLMM method (Section 3.7) on these data provides the following results (ARI): 0.967 (0.09) for $d = 10$, 0.988 (0.04) for $d = 20$ and 0.986 (0.05) for $d = 40$. These results confirms that our proposed extension can cluster the synthetic data with varying K and provides reasonable accuracy.

4.2. *IW-POD dataset.* We consider three different methods, Dynamic Topic Model (DTM) (Blei and Lafferty, 2006), Temporal Multinomial Mixture (TMM) (Kim et al., 2015) and Parametric Link among Multinomial Mixtures (PLMM), for a comparative evaluation of the performance on IW-POD dataset. These methods are selected based on their specialty to cluster discrete evolutionary/temporal data. We set 100 maximum number of iterations as the convergence criterion for all methods. Besides, we set the threshold log-likelihood difference values as 0.0001 for PLMM and TMM. The smoothness parameter α of TMM was set to 1. The DTM hyper-parameter α was set to 0.01. For the PLMM method, we consider the sub-models mentioned in Sec. 3.4.

IW-POD dataset does not provide ground truth cluster labels, due to which we were unable to evaluate clustering results with the known-labels based metric such as **ARI**. In this context, we evaluate the methods using a well known likelihood related measure called *perplexity* on a held-out test set (Murphy, 2012; Blei, Ng and Jordan, 2003). **Perplexity** is a quantity originally used in the field of language modeling (Murphy, 2012). It measures how well a model has captured the underlying distribution of language. In clustering context, *perplexity* is defined as the reciprocal geometric mean of the per feature (word) log-likelihood of a test set, which is computed using the model parameters learned with a training set. Therefore, the *lower perplexity* value indicates that the estimated (trained) model performs *better*

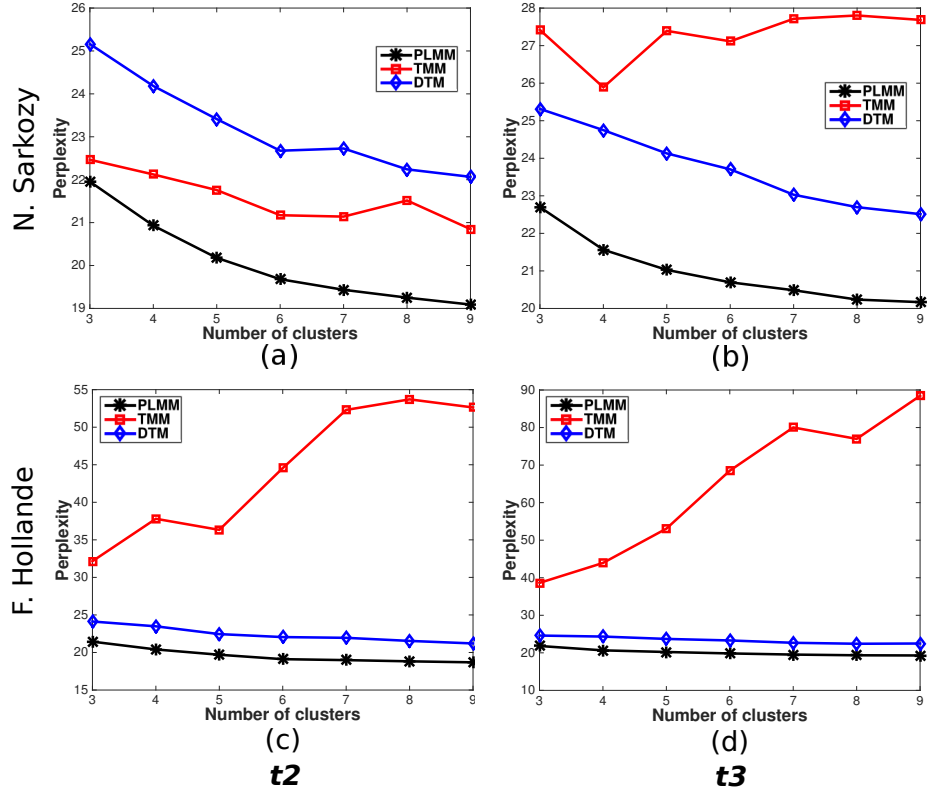


FIG 4.1. Comparison of different methods w.r.t. the perplexity values (*lower is better*) computed from the IW-POD data of two entities (row-1: Sarkozy and row-2: Hollande) and two time epochs (column-1: epoch t2 and column-2: epoch t3). Methods: Dynamic Topic Model (DTM) (Blei and Lafferty, 2006), Temporal Multinomial Mixture (TMM) (Kim et al., 2015) and our proposed Parametric Link among Multinomial Mixtures (PLMM) method.

to fit the test data. *Perplexity* can be formally defined as (Blei, Ng and Jordan, 2003):

$$(4.1) \quad \text{perplexity}(X^{test}) = \exp\left(-\frac{L(\Theta^{train})}{\sum_{i=1}^{N^{test}} V_i}\right)$$

where, V_i is the total number of feature counts (words for document) in observation i , $L(\Theta^{train})$ denotes the log-likelihood of the test data set computed using the trained model parameters Θ^{train} and Eq. (3.6).

In our experiments, for each time epoch t , we compute *perplexity* from 5 folds of training-test data division and then take the average of 5 *perplexity* values as the final measure. For each fold, we used 80% data for training the model and obtain parameters Θ^{train} and the remaining 20% data to compute *perplexity* using Eq. (4.1). Fig. 4.1 illustrates the perplexity values computed from the data of two entities (row-1: Sarkozy and row-2: Hollande) and two time epochs (column-1: epoch $t2$ and column-2: epoch $t3$). Time epoch $t1$ is not considered because it does not reflect the link relationship and temporal aspect of data clustering.

Results in Fig. 4.1 show that, PLMM provides the best *perplexity* compared to DTM and TMM. This means that, compared to other methods, PLMM provides better fitting of the underlying Multinomial distribution to the test data. The next best (3 out of 4) method is the DTM followed by the TMM. Indeed, the results from TMM are intuitive as the fitted models are highly influenced by the other cluster components (Multinomial distributions) from the previous and next time epochs. In contrary, PLMM only consider the link from one cluster in the previous time epoch and fit the data accordingly.

Fig. 4.2 provides a visual illustration of clustering results obtained from the above three methods. This illustration is obtained by using the Multi-dimensional scaling (Kruskal and Wish, 1978) technique where the distance matrix among the observations is computed by first converting the count vectors into probabilities and then using the sKLD (Eq. 3.12) as a measure of distance. The clustering results are obtained with $K = 3$, time epoch $t2$ and the observations associated with the entity NS. From visual comparison among the plots in Fig. 4.2, we can say that PLMM provides better separation than TMM and DTM. Indeed, this observation agrees with the numerical results obtained with the *perplexity* values in Fig. 4.1(a) for $K = 3$.

Next, we apply the extension of PLMM method (Section 3.7) with this dataset and observe the *perplexity* for time epochs $t2$ and $t3$. For the entity NS, we obtain average *perplexity* values as: $t2 : 26.56$ and $t3 : 25.06$ where

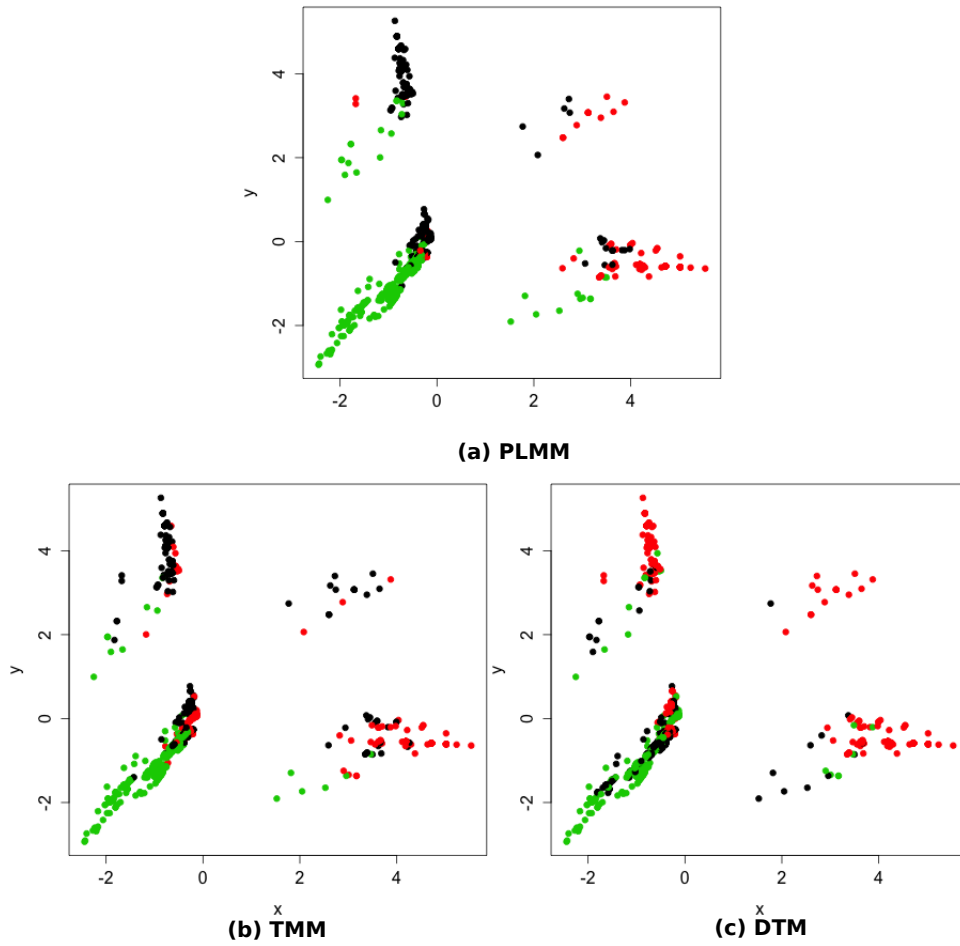


FIG 4.2. Illustration of clustering results visualized with Multidimensional scaling (Kruskal and Wish, 1978). Methods: (a) proposed Parametric Link among Multinomial Mixtures (PLMM); (b) Temporal Multinomial Mixture (TMM) (Kim et al., 2015) and (c) Dynamic Topic Model (DTM) (Blei and Lafferty, 2006).

average K_{t2} is 3 and average K_{t3} is 5. For the entity FH, we obtain average *perplexity* values as: $t2$: 13.08 and $t3$: 5.17 where average K_{t2} is 4 and average K_{t3} is 5. Compared to the results in Fig. 4.1 we see that, *perplexity* values increases (performance decreases) for entity NS and decreases (performance improves) for FH. Based on these observations, we can say that the extension of PLMM provides a good compromise in performance and works well for varying K at different epochs. We do not compare these results with the TMM and DTM methods as they work with fixed K for all time epochs.

Finally, let us focus on the interpretations of cluster evolutions in the IW-POD dataset. Table 5 provides the selection rate of different models at different time epochs (see Table 1 for details of time division). Listed rates provide us very interesting observations from which we can say that:

- The opinions about NS were evolving almost similar way during and after the election period. These evolutions can be interpreted through the belief on aspects using models M3:($\delta_{k,d}/0$) (93%) and M4:($1/\gamma_{k,d}$) (7%). This indicates that during $t1$ - $t2$ - $t3$ opinions about NS were changing slowly.
- Model M2:($0/\gamma_{k,d}$) is selected for all clusters of opinions about FH during $t1$ - $t2$. This means that the opinions change significantly between $t1$ and $t2$ period. From $t2$ to $t3$ (both after election period), opinions were evolving, which can be interpreted through the belief on the features with the models M4:($1/\gamma_{k,d}$) (62%) and M3:($\delta_{k,d}/0$) (38%).

TABLE 5

Selection rate of different models (Sec. 3.4) for the IW-POD dataset at different time epochs (see Table 1 for details of time division).

	M1: ($1/0$)	M4: ($1/\gamma_{k,d}$)	M3: ($\delta_{k,d}/0$)	M2: ($0/\gamma_{k,d}$)
NS ($t1$ - $t2$)	0 %	0 %	100 %	0 %
NS ($t2$ - $t3$)	0 %	13 %	87 %	0 %
FH ($t1$ - $t2$)	0 %	0 %	0 %	100 %
FH ($t1$ - $t2$)	0 %	62 %	38 %	0 %

5. Analysis of the political opinion dataset. In this section, we perform analysis on the clustering results only from the PLMM method. In order to visualize the contents, we construct a histogram representation, which helps us to discriminate among different clusters. These histograms are constructed by counting the polarities (in vertical direction) w.r.t. each attribute (in horizontal direction). The color of the bars resembles the color of polarities. Fig. 5.1 illustrates an example of a histogram which is con-

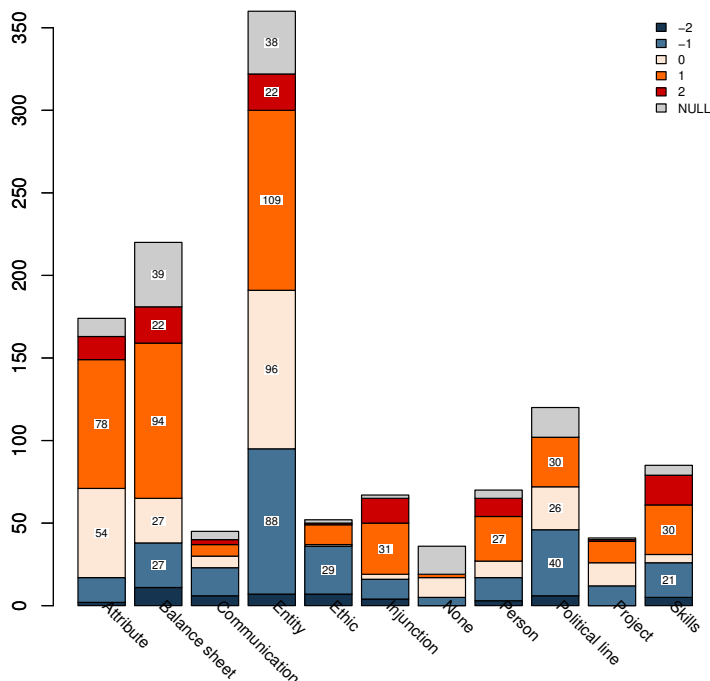


FIG 5.1. Illustration of the clustering results using a histogram constructed from the polarities of different aspects. The aspects are ordered from left to right as: (1) Attribute; (2) Balance sheet; (3) Communication; (4) Entity; (5) Ethic; (6) Injunction; (7) None; (8) Person; (9) Political line; (10) Project and (11) Skills. The polarities are colored and ordered from bottom to top as: -2 (dark blue), -1 (blue), 0 (light orange), 1 (orange), 2 (red) and NULL (grey).

constructed from the tweets of a cluster from time t_2 . Following this illustration, in Fig. 5.2 and 5.3, let us look at the examples of the clusters at different time epochs for the entities NS and FH respectively. These results are obtained by clustering data using PLMM method with $K = 3$. From both figures we observe that, at each time epoch the clusters have different histogram representations. Moreover, during different time epochs each cluster undergoes certain amount of changes in different attributes and associated polarities. This demonstrates that the proposed PLMM method is able to provide sufficient inter-cluster variations (at each time) while respecting the temporal dynamics (for each cluster during different time epochs).

An alternative and compact representation (w.r.t. the MM model parameters) of the clusters for NS is illustrated in Fig. 5.4(a) and 5.4(b). Similar to the examples of Fig. 5.2, this alternative representation demonstrate that, at a certain time epoch different cluster emphasizes on different as-

pects/polarities of an entity. Besides, the temporal changes of the clusters can be identified subsequently during different epochs by observing the increase/decrease of the probabilities. However, from the user’s perspective, this representation may not be convenient to understand. Therefore, we use histograms for further analysis and use this compact representation for a different purpose.

Now, let us explain the semantics obtained from these clustering results. For brevity, here we denote a cluster as *cl.* From Fig. 5.2 (clusters for NS) we see that, while *cl.* 1 and 3 emphasize on the negative (-) and positive (+) polarities respectively, *cl.* 2 emphasizes on a particular attribute. Naively we can say that, there are three groups of peoples: (a) the first group (*cl.* 1) provides negative opinions from various aspects, thus tends to hold a negative image about the entity; (b) the second group (*cl.* 2) particularly emphasizes on *Ethic* of the entity and mostly provide negative opinions and (c) the third group (*cl.* 3) can be seen as a contrary to the first group (*cl.* 1) as it tends to hold a positive image about the entity. Table 6 provides three examples of the tweets for time *t1* and for each cluster about NS. We can realize that these tweets reflect the opinions which truly correspond to the groups obtained by the clustering method.

From temporal viewpoint, we observe several changes w.r.t. different aspects. In order to analyze the changes using histograms, we observe the height of histogram bar for each aspect. This height indicates the number of tweets/opinions corresponding to the related aspect. Let us consider an example of the aspect *Communication* which plays a certain role on clustering. We observe that: (a) for *cl.* 1, the total number of tweets related to the aspect *Communication* remains same during time *t1* and *t2* and reduces during *t2* and *t3*; (b) for *cl.* 2, the total number of tweets related to this aspect reduces continuously and (c) for *cl.* 3, the total number of tweets related to this aspect reduces from *t1* to *t2* and remains same during *t2* to *t3*. Moreover, a closer look on *cl.* 3 from *t2* to *t3* reveals an increase of positive opinions about the *communication* skill of the entity. Another example is the aspect called *Attribute*, whose height reduces continuously with time for both *cl.* 1 and 3. Similarly, from an analysis of the height of histogram bars in Fig. 5.3 (clusters for FH) we see that, the aspects called *Entity*, *Ethic*, *Political line*, *Skills* and *Communication* play certain role to describe the image of FH. For example, the tweet - *Holland would remove the word “race” in the Constitution (orig: Hollande supprimerait le mot “race” dans la Constitution)* from time *t1* and *cl.* 3 is annotated with the aspect called *political line* and polarity *+1*. Another tweet - *Holland and Netanyahu evoke the struggle against anti-Semitism (orig: Hollande et Netanyahu évoquent*

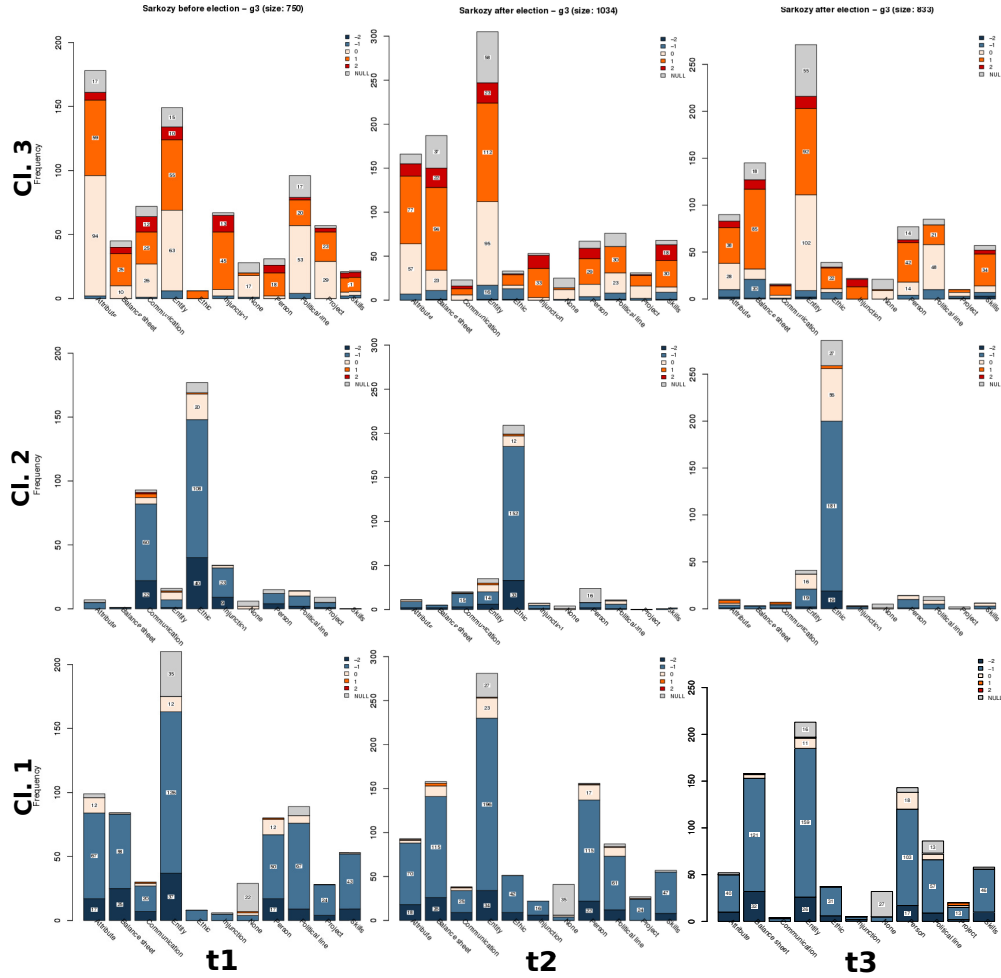


FIG 5.2. Illustration of the clustering results from PLMM methods for NS. Results obtained using $K = 3$ for three time epochs t_1 , t_2 and t_3 . Each cluster is represented as a histogram constructed from the polarities of different aspects. The aspects are ordered from left to right as: (1) Attribute; (2) Balance sheet; (3) Communication; (4) Entity; (5) Ethic; (6) Injunction; (7) None; (8) Person; (9) Political line; (10) Project and (11) Skills. The polarities are colored and ordered from bottom to top as: -2 (dark blue), -1 (blue), 0 (light orange), 1 (orange), 2 (red) and NULL (grey). Each column represents clusters from a particular epoch. Each row represents a particular cluster in different epochs.

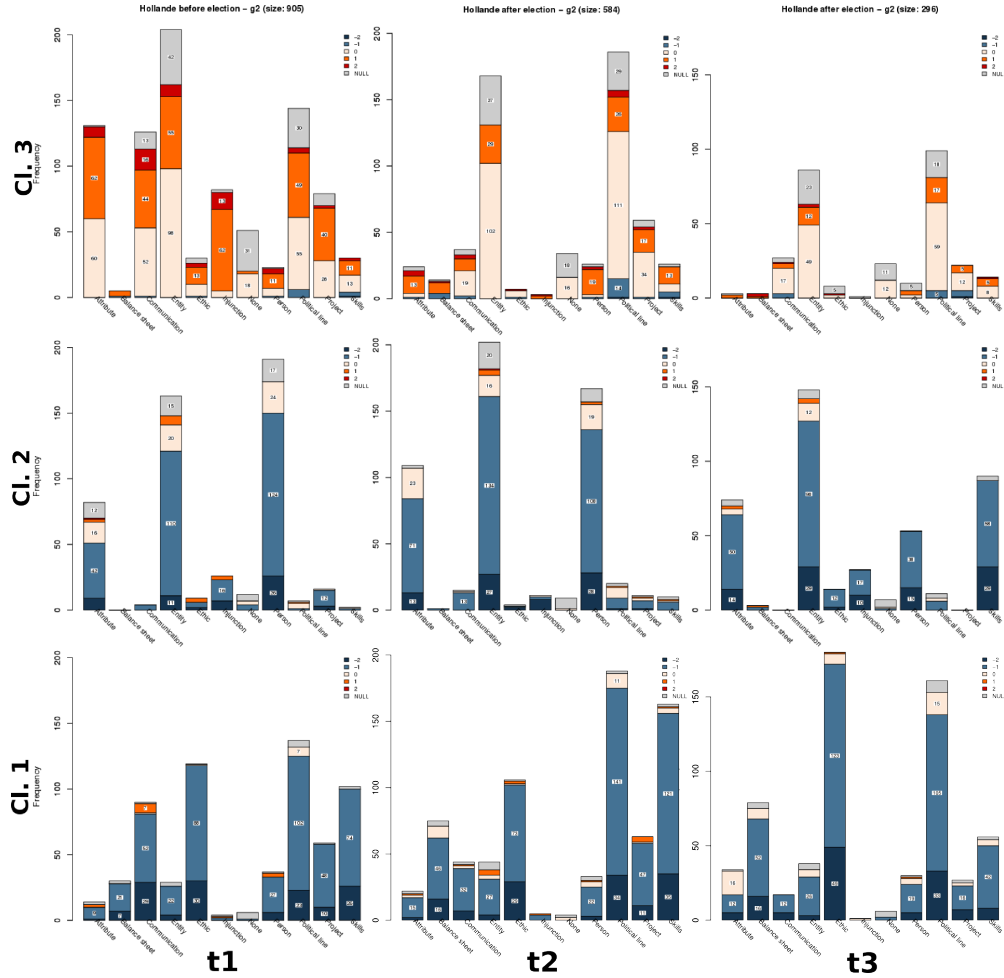


FIG 5.3. Illustration of the clustering results from PLMM methods for FH. Results obtained using $K = 3$ for three time epochs t_1 , t_2 and t_3 . Each cluster is represented as a histogram constructed from the polarities of different aspects. The aspects are ordered from left to right as: (1) Attribute; (2) Balance sheet; (3) Communication; (4) Entity; (5) Ethic; (6) Injunction; (7) None; (8) Person; (9) Political line; (10) Project and (11) Skills. The polarities are colored and ordered from bottom to top as: -2 (dark blue), -1 (blue), 0 (light orange), 1 (orange), 2 (red) and NULL (grey). Each column represents clusters from a particular epoch. Each row represents a particular cluster in different epochs.

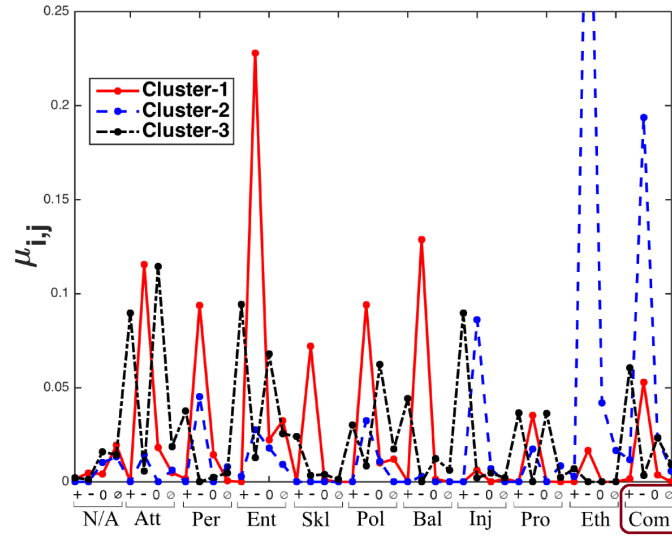
la lutte contre l'antisémitisme) has the same annotation which is from the same cluster but from time $t3$. These two examples reveal the importance of the aspect *political line* for keeping the similar opinions into the same group at different time. The above observations clearly indicate that, for different groups of people different aspects has certain importance at different time. Therefore, an analyst can retrieve the most prominent aspects from people's opinion about an entity at a particular time or within a certain range of time periods.

Besides the above interpretation of the clustering results, an analyst can obtain more information from the PLMM clustering results via the link parameters ($\delta_{k,d}$ or $\gamma_{k,d}$). After analyzing the links among MM parameters we notice that they are able to provide a compact explanation about the temporal changes during two time epochs. Fig. 5.4 illustrates an example for entity *NS* from time $t1$ to $t2$ with 3 clusters, see column 1 and 2 of Fig. 5.2 for corresponding histograms. Fig. 5.4(a) and Fig. 5.4(b) illustrates the MM parameters (probability of aspect-polarity features) and Fig. 5.4(c) provides a compact representation about the cluster evolutions using the values of $\delta_{k,d}$. To better understand this representation in Fig. 5.4(c), we transform the link values as 0 (no change), -1 ($\delta_{k,d} < 0.9$, belief increases) and +1 ($\delta_{k,d} > 1.1$, belief decreases). In the context of the examples from the IW-POD, we can explain belief as: probability of a feature at time $t + 1$ is increased from its probability at time t . Therefore, the belief indicates the relative significance of a particular feature w.r.t. time. An increase in the belief means that users tend to be more attracted by it. Following this, if a feature probability is nearly same at two different times then belief remains unchanged. In Fig. 5.4, we highlight the effect of a particular aspect, called *Communication* (*Com*), and observe its contribution for cluster evolution. From Fig. 5.4 (a) and (b) we see that, from time $t1$ to $t2$ the probabilities are decreased mostly for *cl. 2* and *cl. 3*. This means that, either the users from these clusters loose interest to discuss about *Com* and focus on other aspects, or those users disappeared at time $t2$. Similar to *Com*, we can observe other aspects such as *Eth* (*cl. 1* and *cl. 3*) and *Ent* (*cl. 2* and *cl. 3*) which causes cluster evolution in this example of Fig. 5.4.

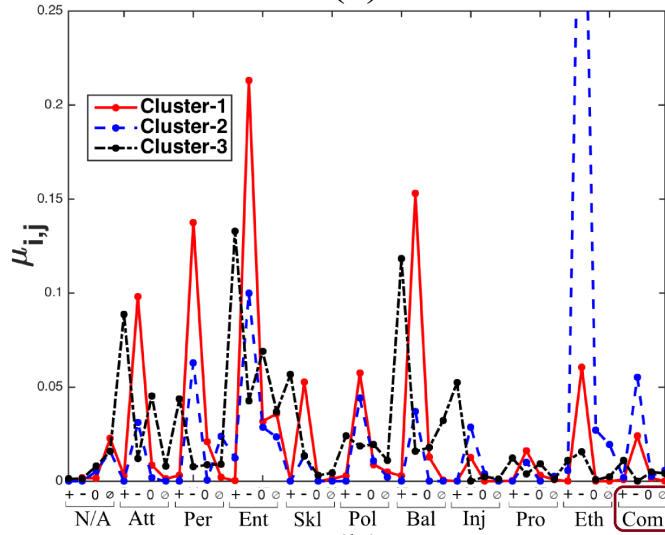
Let us analyze examples from real twitter data and observe them w.r.t. the Fig. 5.4. If we look at *cl. 3* at time $t1$ (before election), the most likely features are often positive and it is clear that it gathers people in favor of NS. The prominent aspects are *Att* (positive and neutral), *Ent* (positive) and *Inj* (positive), such as in the tweet - *40 people @youngpop44 will be present at the great gathering in Place #Concorde for supporting @NicolasSarkozy ! #StrongFrance #NS2012*". This cluster slightly changes later at time $t2$

TABLE 6
 Real twitter data examples of the 3 clusters at time t1 for entity NS. See Fig. 5.2 column 1 for the associated histograms.

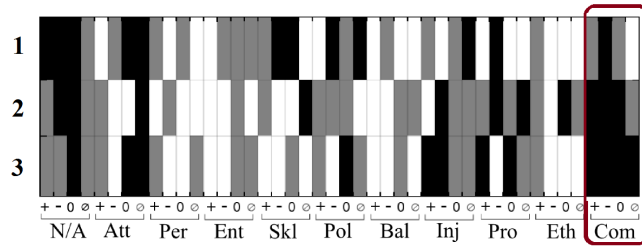
<i>Cluster 1 (Generally Negative)</i>	
<i>Ex. 1</i>	<p>Orig: Il veut desréférendums car... y a pas de pilote dans l'avion, dit-il: quel aveu! #Sarkozy#projet</p> <p>Trans: He wants referendumbecause... there is no pilot in the plane he says: what a confession! #Sarkozy#project</p>
<i>Ex. 2</i>	<p>Orig: Je ne voterais pas #Sarkozy ! ” ” Je ne voterais pas #Sarkozy !</p> <p>Trans: I won't vote for #Sarkozy !” ” I won't vote for #Sarkozy</p>
<i>Ex. 3</i>	<p>Orig:Nicolas Sarkozy, le plus mauvais président de la Vème République</p> <p>Trans: Nicolas Sarkozy, the worst president of the Fifth Republic</p>
<i>Cluster 2 (Negative, specially "Ethic")</i>	
<i>Ex. 1</i>	<p>Orig: Jamais un président n'a été cerné par tant d'affaires! demain ds @lematinch #Bettencourt #Sarkozy</p> <p>Trans: Never before a president was surrounded by so many cases!" tomorrow in @lematinch #Bettencourt #Sarkozy</p>
<i>Ex. 2</i>	<p>Orig: Une liste de condamnés de l'#UMP qui pourrait être bientôt complétée par les noms de #Sarkozy, #Copé, #Woerth</p> <p>Trans: A list of convicted people of #UMP soon completed by names such as #Sarkozy, #Copé, #Woerth (the "Bettencourt case" is a famous case in which Sarkozy was involved)</p>
<i>Ex. 3</i>	<p>Orig: Sarkozy-Kadhafi: la preuve du financement. Et l'urgence d'une enquête officielle #affaireetat</p> <p>Trans: Sarkozy-Kadhafi: the proof of funding. And the urge of an official enquiry #stateaffair (Kadhafi is another case in which Sarkozy was involved in some way)</p>
<i>Cluster 3 (Generally Positive)</i>	
<i>Ex. 1</i>	<p>Orig: N Sarkosy mots clé..challenge, défi, action, travail, réussite, formation, effort, individualisation ..France Forte. Europe Forte #NS2012</p> <p>Trans: N Sarkozy keywords..challenge, défi, action, work, success, training, effort, individualization ..Strong France. Strong Europe #NS2012</p>
<i>Ex. 2</i>	<p>Orig: merci N.Sarkozy pour tout tu restera pour toujours mon Hero merci. merci</p> <p>Trans: Thank you N.Sarkozy for all you will stay my hero forever thanks. thanks</p>
<i>Ex. 3</i>	<p>Orig: Sarko est plus rationnel..</p> <p>Trans: Sarko is more rational..</p>



(a)



(b)



(c)

FIG 5.4. Example of evolution interpretation using link parameter $\delta_{k,d}$ for NS during $t1$ to $t2$ with 3 clusters. (a) MM parameters $\mu_{k,j}^{t1}$ at time $t1$ (b) MM parameters $\mu_{k,j}^{t2}$ at time $t2$ (c) Link parameters $\delta_{k,j}$ between time $t1$ and $t2$. In (c), for each cluster (row-wise), brighter/white color indicates the prior belief about features (aspect-polarity) increases, darker/black color indicates the prior belief about features decreases and grey color indicates the prior belief about features remains same.

(just after election) towards *Att* (positive), *Ent* (positive) and *Bal* (positive). The shift from *Inj* to *Bal* is clearly visible on Fig. 5.4(c), third row: black color for *Inj* means a decrease of attention whereas white color for *Bal* means there are relatively more comments on the balance sheet of NS. Hence, the following message shows some nostalgia felt by many militants: *Whatever the opinion of FH, NS has been a great president. FH can deconstruct all the reforms, we will never forget!*. To sum up, the δ parameter helps us to focus on what are the main changes, even though the observation could have been drawn among the other aspects. Following the same reasoning, all polarities targeting the aspect *Com* are black, which proves that the performances of the politician in the media (e.g., TV, newspapers) are less important once the election is over.

Observations from numerous experiments reveal that, besides performing evolutionary clustering on the temporal data, PLMM also provide reasonable interpretation for the evolutions, thanks to the link parameters. Indeed, this clearly distinguishes PLMM from the rest of the state-of-the-art methods. Moreover, we notice that the interpretability of PLMM (using Eq. 3.9, 3.10 and 3.11) can be separated out and externally plugged in with the results from any other discrete data clustering methods.

6. Conclusion and Future Perspectives. Over the years, a large number of temporal data analysis methods have been proposed in several domains. In this paper, we only focused on the particular clustering methods which have been used for discrete data clustering and which are based on the assumption of the Multinomial distribution.

We proposed an unsupervised method (i.e., no training from labeled data) for analyzing the temporal data. The core element of our proposal is the formulation of parametric links among the multinomial distributions. Computations of these links naturally cluster the evolutionary/temporal data. Furthermore, these links can provide interpretation for cluster evolution and also detect clusters evolution in certain cases. For experimental validation, we extensively used synthetic dataset and evaluated using the *Adjusted Rand Index*. As a practical application, we applied it on a dataset of political opinions and evaluated using *Perplexity* measure. Results show that the proposed method, called PLMM, is better than the state-of-the-art. Moreover, it provides an additional advantage through the link parameters in order to interpret the changes in clusters at different time. We also provide an extension of the proposed method for dealing with varying number of clusters which is not addressed by most of the recent methods.

Monitoring/tracking cluster evolution is an interesting issue which we do

not explicitly and extensively manage in our proposed method, because it is not a primary objective in this paper. Yet, we can partially achieve this task by using certain information (parametric sub-models, see 3.4) which are naturally integrated with our proposed method. That means, our proposed method can be used only as a detector of cluster evolution. At present, we consider the complete monitoring task as a future work. We believe that, an extension of several existing work can be added with our method to completely deal with this issue. For example, we can exploit⁷ MEC (Oliveira and Gama, 2010) which is a cluster evolution monitoring method for continuous data. Besides, we can use *label-based diachronic approach* (Lamirel, 2012) by externally providing our clustering results as an input to it.

Computational complexity is a concern for the proposed method and can be considered as a limitation. From a decomposition of the computational time, we observe that most of the time is consumed by the optimization procedure (*neldermead* simplex method). In future, a better optimization method can be incorporated to address this issue. Moreover, the time can be further reduced by eliminating the parametric sub-models which are experimentally found as redundant.

Although we demonstrated the effectiveness of the proposed method only for political opinion dataset, we believe that it will be equally effective for different datasets that consist of the form of categorical data.

References.

- AGRESTI, A. (2002). *Categorical data analysis*, 2nd ed. John Wiley & Sons.
- BAUDRY, J.-P. and CELEUX, G. (2015). EM for mixtures-Initialization requires special care. *Statistics and Computing* **25** 713-726.
- BENINEL, F., BIERNACKI, C., BOUVEYRON, C., JACQUES, J. and LOURME, A. (2012). *Parametric link models for knowledge transfer in statistical learning. Knowledge Transfer: Practices, Types and Challenges*. Nova Science Publishers.
- BIERNACKI, C., BENINEL, F. and BRETAGNOLLE, V. (2002). A generalized discriminant rule when training population and test population differ on their descriptive parameters. *Biometrics* **58** 387-397.
- BIERNACKI, C., CELEUX, G. and GOVAERT, G. (2000). Assessing a mixture model for clustering with the integrated completed likelihood. *IEEE TPAMI* **22** 719-725.
- BIERNACKI, C., CELEUX, G. and GOVAERT, G. (2003). Choosing starting values for the EM algorithm for getting the highest likelihood in multivariate Gaussian mixture models. *Computational Statistics & Data Analysis* **41** 561-575.
- BIERNACKI, C., CELEUX, G., GOVAERT, G. and LANGROGNET, F. (2006). Model-based cluster and discriminant analysis with the MIXMOD software. *Computational Statistics & Data Analysis* **51** 587-600.

⁷We conducted some initial experiments and found that this approach is applicable up to certain extent and should be further improved to use in our case, e.g., extend it with appropriate distance computation (e.g., using sKLD).

- BISHOP, C. M. et al. (2006). *Pattern recognition and machine learning* **4**. springer New York.
- BLEI, D. M. and LAFFERTY, J. D. (2006). Dynamic topic models. In *Proc. of the Int Conf on Machine Learning* 113–120. ACM.
- BLEI, D. M., NG, A. Y. and JORDAN, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research* **3** 993–1022.
- CHAKRABARTI, D., KUMAR, R. and TOMKINS, A. (2006). Evolutionary clustering. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining* 554–560. ACM.
- CHI, Y., SONG, X., ZHOU, D., HINO, K. and TSENG, B. L. (2009). On evolutionary spectral clustering. *ACM Trans. on Knowledge Discovery from Data* **3** 17.
- DUBEY, A., HEFNY, A., WILLIAMSON, S. and XING, E. P. (2013). A Nonparametric Mixture Model for Topic Modeling over Time. In *SDM* 530–538. SIAM.
- FERLEZ, J., FALOUTSOS, C., LESKOVEC, J., MLADENIC, D. and GROBELNIK, M. (2008). Monitoring network evolution using MDL. In *IEEE Int. Conf. on Data Engineering* 1328–1330. IEEE.
- FIGUEIREDO, M. A. T. and JAIN, A. K. (2002). Unsupervised learning of finite mixture models. *IEEE TPAMI* **24** 381–396.
- FRALEY, C. and RAFTERY, A. E. (2002). Model-based clustering, discriminant analysis, and density estimation. *Journal of the American Statistical Association* **97** 611–631.
- GARCIA, V. and NIELSEN, F. (2010). Simplification and hierarchical representations of mixtures of exponential families. *Signal Processing* **90** 3197–3212.
- HASNAT, M. A., ALATA, O. and TRÉMEAU, A. (2015). Model-based hierarchical clustering with Bregman divergences and Fishers mixture model: application to depth image analysis. *Statistics and Computing* 1-20.
- HASNAT, M. A., VELCIN, J., BONNEVAY, S. and JACQUES, J. (2015). Simultaneous Clustering and Model Selection for Multinomial Distribution: A Comparative Study. In *Advances in Intelligent Data Analysis XIV* Springer.
- HUBERT, L. and ARABIE, P. (1985). Comparing partitions. *Journal of classification* **2** 193–218.
- JACQUES, J. and BIERNACKI, C. (2010). Extension of model-based classification for binary data when training and test populations differ. *Journal of Applied Statistics* **37** 749–766.
- KHARRATZADEH, M., RENARD, B. and COATES, M. (2015). Bayesian topic model approaches to online and time-dependent clustering. *Digital Signal Processing*.
- KIM, Y.-M., VELCIN, J., BONNEVAY, S. and RIZOIU, M.-A. (2015). Temporal Multinomial Mixture for Instance-Oriented Evolutionary Clustering. In *Proc. of the European Conference on Information Retrieval* 593–604.
- KRUSKAL, J. B. and WISH, M. (1978). *Multidimensional scaling* **11**. Sage.
- LAMIREL, J.-C. (2012). A new approach for automatizing the analysis of research topics dynamics: application to optoelectronics research. *Scientometrics* **93** 151–166.
- LI, P., WU, X. and HU, X. (2012). Mining recurring concept drifts with limited labeled streaming data. *ACM Transactions on Intelligent Systems and Technology (TIST)* **3** 29.
- MCLACHLAN, G. J. and KRISHNAN, T. (2008). *The EM algorithm and extensions*, 2. ed ed. *Wiley series in probability and statistics*. Wiley.
- MEILĀ, M. and HECKERMAN, D. (2001). An experimental comparison of model-based clustering methods. *Machine Learning* **42** 9–29.
- MELNYKOV, V. and MAITRA, R. (2010). Finite mixture models and model-based clustering. *Statistics Surveys* **4** 80–116.

- MURPHY, K. P. (2012). *Machine learning: a probabilistic perspective*. The MIT Press.
- NELDER, J. A. and MEAD, R. (1965). A simplex method for function minimization. *The computer journal* **7** 308–313.
- OLIVEIRA, M. D. and GAMA, J. (2010). MEC-Monitoring Clusters' Transitions. In *STAIRS* 212–224.
- SALVADOR, S. and CHAN, P. (2004). Determining the number of clusters/segments in hierarchical clustering/segmentation algorithms. In *IEEE Conf. on Tools with Artificial Intelligence* 576–584.
- SCHWARZ, G. et al. (1978). Estimating the dimension of a model. *The Annals of Statistics* **6** 461–464.
- SILVESTRE, C., CARDOSO, M. G. and FIGUEIREDO, M. A. (2014). Identifying the number of clusters in discrete mixture models. *arXiv preprint arXiv:1409.7419*.
- SPILIOPOULOU, M., NTOUTSI, I., THEODORIDIS, Y. and SCHULT, R. (2006). MONIC: modeling and monitoring cluster transitions. In *Proc. of the ACM SIGKDD Int conf. on Knowledge discovery and data mining* 706–711. ACM.
- VELCIN, J., KIM, Y., BRUN, C., DORMAGEN, J., SANJUAN, E., KHOUAS, L., PERADOTTO, A., BONNEVAY, S., ROUX, C., BOYADJIAN, J. et al. (2014). Investigating the Image of Entities in Social Media: Dataset Design and First Results. In *Proc. of Language Resources and Evaluation Conference (LREC)*.
- XU, K. S., KLIGER, M. and HERO III, A. O. (2014). Adaptive evolutionary clustering. *Data Mining and Knowledge Discovery* **28** 304–336.
- XU, T., ZHANG, Z., YU, P. S. and LONG, B. (2008). Dirichlet process based evolutionary clustering. In *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on* 648–657. IEEE.
- XU, T., ZHANG, Z., YU, P. S. and LONG, B. (2012). Generative models for evolutionary clustering. *ACM Trans. on Knowledge Discovery from Data* **6** 7.
- YPMA, J. (2014). Introduction to nloptr: an R interface to NLOpt.
- ZHONG, S. and GHOSH, J. (2005). Generative model-based document clustering: a comparative study. *Knowledge and Information Systems* **8** 374–384.