

Creation and evaluation of timelines for longitudinal user posts

Anonymous ACL submission

Abstract

001 There is increasing interest to work with user
002 generated content in social media and espe-
003 cially textual posts over time. Currently there
004 is no consistent way of segmenting user posts
005 into timelines in a meaningful way that can im-
006 prove the quality and cost of manual annota-
007 tion. Here we propose a set of methods for
008 segmenting longitudinal user posts into time-
009 lines that are likely to contain interesting mo-
010 ments of change in a user’s behaviour based on
011 the content they have shared online and their
012 online activity. We also propose a framework
013 for evaluating the timelines returned in terms
014 of containing candidate moments of change in
015 close proximity to manually annotated time-
016 lines that are dense in such moments of change.
017 Finally, we present a discussion of the linguis-
018 tic content of highly ranked timelines.

019 1 Introduction

020 An increasing body of work considers time-aware
021 models trained on social media data for a number
022 of different tasks, including personal event identi-
023 fication (Li and Cardie, 2014; Li et al., 2014; Chang
024 et al., 2016a), suicidal ideation and suicide risk de-
025 tection (Coppersmith et al., 2014, 2018; Cao et al.,
026 2019; Matero et al., 2019; Sawhney et al., 2020,
027 2021). For such tasks deriving meaningful *time-*
028 *lines* (i.e. relatively short sequences of posts by in-
029 dividuals, containing examples of the phenomenon
030 under study) from large-scale collections, together
031 with associated annotations, is crucial. This is es-
032 pecially important for computational approaches
033 in mental health given surging numbers of those
034 seeking help online (Neary and Schueller, 2018).

035 Earlier work on personal life event detection
036 had considered selecting salient timelines through
037 topic modelling (Li and Cardie, 2014; Li et al.,
038 2014) or through a non-parametric generative ap-
039 proach (Chang et al., 2016a). However, such ap-
040 proaches are not suitable for identifying changes in

mood or mental health more generally. Specifically,
since timelines are selected based on linguistic con-
tent this introduces a sampling bias for downstream
linguistic analysis and annotation (Olteanu et al.,
2019; Mishra et al., 2019). In recent work on suici-
dal ideation detection, timelines are chosen as the
 N most recent posts (Sawhney et al., 2020), which
are not necessarily the most salient for annotation
purposes.

Present Work: We propose a set of methods and
an associated evaluation framework for identifying
salient timelines from the history of social media
users to be annotated for the presence of *Moments*
of Change (MoC). We define a MoC as a particular
point or set of points in time denoting: (1) a shift
in an individual’s mood from positive-to-negative
or vice versa; (2) gradual mood progression. The
aim is to identify methods which can consistently
select timelines that are rich in MoC for large scale
cost-effective annotation. We follow earlier work
in hypothesising that posting behaviour can be used
as a proxy for changes in mental health (De Choud-
hury et al., 2016). Therefore we present methods
for creating timelines based on time-series of post-
ing frequency, such as change-point and anomaly
detection approaches, and evaluate these against
keyword-based methods and randomly selected
timelines. All candidate timelines are evaluated
against manually annotated MoC. We make the
following contributions:

- We present the first approach to extracting time-
lines from users’ posting history on social me-
dia based on change-point detection methods,
anomaly detection and kernel density estima-
tion (see §3).
- We propose a novel evaluation framework for
assessing the quality of annotated timelines,
and timeline extraction methods, on the basis
of manually annotated MoCs (see §4).
- We provide an insightful linguistic analysis into

081	highly ranked (dense in MoCs) timelines and	interested in viral buzzes of mentions of celebri-	129
082	timelines sparse in MoCs (see §5.2).	ties on social media, and as such aims to identify	130
083	2 Related Work	salient dates by simultaneously modelling linguis-	131
084	2.1 Tracking Changes in Mental Health	tic content and frequency based time-series pat-	132
085	Moments of Change (MoC) are an important	terns. While CPD has been explored in news TLS	133
086	concept in work on mental health tracking. Pruk-	(Hu et al., 2011), it remains under-explored for	134
087	sachatkun et al. (2019) identifies a MoC as a posi-	social media data.	135
088	tive change in sentiment for a user with respect to a	2.3 Change-point Detection	136
089	particular distressing topic mentioned in a conver-	In §3, we explore using automatically detected	137
090	sation thread. De Choudhury et al. (2016) investi-	change-points (candidate MoCs) as the salient	138
091	gated shifts to suicide ideation by building models	dates used to select timelines of users on social	139
092	to predict transition of a user posting on a suicide	media for annotation.	140
093	support forum. We consider a more general defini-	Change-points (CPs) are typically defined as	141
094	tion of MoC (see §1, “Present Work”).	points in time where the underlying generative pa-	142
095	Creation of Mental Health Datasets. A large	rameters of a data sequence are predicted to have	143
096	body of work in creating mental health datasets in-	changed (van den Burg and Williams, 2020). CPD	144
097	volves labelling posts for symptoms (Gkotsis et al.,	approaches, therefore, involve learning a predictive	145
098	2017 ; Loveys et al., 2017 ; Cheng et al., 2017) or	model of a data sequence. While there are several	146
099	levels of suicide ideation (Masuda et al., 2013 ; Cop-	continuous models (e.g. a Gaussian model (Adams	147
100	persmith et al., 2016 ; Shing et al., 2018). While	and MacKay, 2007)), we are particularly inter-	148
101	annotations for some of these datasets are obtained	ested in models suited to discrete event-based time-	149
102	through proxy signals (e.g., self-disclosure of diag-	stamped data (Knoblauch and Damoulas, 2018) -	150
103	gnoses, posts on support networks) a question arises	such as points in time where a post/comment is	151
104	as to how to select appropriate data for annotation.	made on social media. In such scenarios Temporal	152
105	Mishra et al. (2019) use keyword based methods	Point Processes (TPPs) (Daley and Vere-Jones,	153
106	to identify posts exhibiting the phenomenon un-	2003) are particularly well suited.	154
107	der study (e.g. suicidal ideation) but this leads to	Temporal Point Processes (TPPs) TPPs are de-	155
108	sampling biases. An alternative is to consider time-	defined as stochastic processes modelling discrete	156
109	line extraction approaches agnostic to the linguistic	events occurring on a continuous time domain.	157
110	content, inspired by Timeline Summarisation (TLS)	They are typically characterized by an intensity	158
111	and Change-Point Detection (CPD).	function, $\lambda > 0$, which represents the instanta-	159
112	2.2 Timeline Summarization (TLS)	aneous rate of event occurrence. TPPs vary in com-	160
113	TLS aims to provide concise chronologically or-	plexity: from the simple homogeneous Poisson	161
114	dered timelines consisting only of the most relevant	process (a model governed by a constant λ), to the	162
115	information for a given topic or entity, summarizing	more flexible Hawkes process (RizoIU et al., 2017)	163
116	the key points in time. While TLS has been most	(which has a conditional λ : dependent on both time	164
117	commonly applied in news topic summarization	and historical events), to the rapidly developing	165
118	(Swan and Allan, 2000 ; Martschat and Markert,	field of neural temporal point processes (Shchur	166
119	2017, 2018 ; Steen and Markert, 2019), there has	et al., 2021 ; Lin et al., 2021) (where λ is modelled	167
120	been growing interest in applying TLS applied on	with highly flexible neural networks, such as RNNs	168
121	social media data (Li and Cardie, 2014 ; Chen et al.,	(Du et al., 2016) or more recently models based on	169
122	2019 ; Ansah et al., 2019 ; Wang et al., 2021).	self-attention (Zhang et al., 2020 ; Zuo et al., 2020)).	170
123	TLS consists of a two-step pipeline, where (1)	In order to use TPPs to model event sequences, and	171
124	date selection is followed by (2) summarisation.	predict associated changes - certain CPD models,	172
125	Salient dates to summarize as a timeline are typ-	such as Bayesian Online Change-point Detection	173
126	ically identified using textual content, as well as	(Adams and MacKay, 2007) require that the TPP	174
127	time-series frequency information in the history	be part of the exponential family of distributions	175
128	of an individual / topic. Chang et al. (2016b,a) is	(e.g. the Gaussian distribution, or the Poisson dis-	176
		tribution). This is so that the intensity λ can be	177
		further modelled from a prior conjugate distribu-	178

tion, making it possible to construct the likelihood of the chosen predictive model in a closed form. TPPs part of the exponential family of distributions, specifically the Poisson-Gamma predictive model, therefore form a class of computationally inexpensive models that are scalable to large datasets, making them particularly attractive for our task.

3 Approach

Task. Our principal aim is to select timelines for annotation that are rich in MoC. To achieve this, we test a series of timeline extraction methods presented in this section, which we then evaluate using a novel evaluation framework in §4.

Selecting Candidate Timelines. To select timelines for annotation, we extract candidate timelines as a span of timestamps $S_{i,u}$ from a user’s u history H_u . To do so we first propose identifying *Candidate Moments of Change* (CMoC), which are dates predicted to be surrounded by many MoCs (§3.1). Subsequently, we extract the user’s posts surrounding these CMoC within a fixed time window, as timelines to be returned for annotation (§3.2).

3.1 Identifying Candidate MoCs (CMoC)

We explore different approaches for identifying CMoC, as detailed below:

(1) Change-point Detection (CPD): In a recent evaluation involving experiments with both synthetic and real-world change-points, [van den Burg and Williams \(2020\)](#) showed that Bayesian Online Change-point Detection (BOCPD) was the best performing model for a variety of CPD tasks. BOCPD functions by learning a predictive model on a data sequence, and when changes in the model’s underlying generative parameters are identified, a change-point is declared. The models which BOCPD is typically fit with continuous (e.g. the Gaussian distribution). However, it is also possible to use temporal point processes (§2.3) which are more appropriate for modelling discrete event-based data ([Knoblauch and Damoulas, 2018](#)).

Since we hypothesize that changes in posting behaviour coincide with changes in mood (see “Present Work” in §1), we use BOCPD to identify changes in individuals’ posting frequency. As such we consider the daily frequency of posts made by a user as a Temporal Point Process, and use the homogeneous Poisson-Gamma (PG) point process model with BOCPD ([Knoblauch and Damoulas, 2018](#)) to fit and identify changes in the daily fre-

quency of posts by a user u from their entire associated history H_u . Note that we investigate this hypothesis by evaluating the density of changes in mood from timelines selected this way in our results section (§5.2), and also investigate changes in posting activity coinciding with changes in mood and sentiment in the same section (table 2).

By using a PG model with BOCPD, we assume that each point in a user’s posting frequency is sampled from a Poisson distribution with a discrete intensity λ . Here λ represents the expected number of posts by a user within a given time interval. As we use this conjugate Bayesian model, λ is further assumed to be drawn from a Gamma distribution with a set of priors α_0 and β_0 , that act as initial hyper-parameters in our model, where α_0/β_0 , α_0/β_0^2 denote the prior mean and variance over the intensity of the time-series of the data. BOCPD has an additional hyper-parameter which is the hazard, h_0 where $1/h_0$ expresses a prior belief about the probability of change-points (CPs) occurring at a given time t , provided that a CP has not recently occurred: a low h_0 results in the over-generation of change-points while a large h_0 is more conservative and returns very few change-points (ideal in our scenario, to ensure that we do not waste annotation resources, by avoiding annotating too many timelines generated by noise).

Since BOCPD computes a full probability distribution over the location of the CPs, quantifying probable CPs along with their associated uncertainty, we use the maximum a posteriori (MAP) segmentation of the probability distribution to return exact point estimates for CPs ([Fearnhead and Liu, 2007](#); [van den Burg and Williams, 2020](#)). These predicted points in time can represent CMoCs. An illustration of identifying CMoCs from a given user’s history in our implementation of BOCPD is provided in Fig. 1. Here change-points define CMoCs.

(2) Anomaly Detection (AD): Here we aim at identifying (a) days of abnormally high user activity and (b) abnormally long time periods of no user activity at all. We hypothesize that such points in time can be used to select salient timelines. We experiment using different features to fit our model, including the daily frequency of a user’s posts and the number of comments they receive for those corresponding posts by others. Using either activity type, we scan over the user’s entire history.

For (a) we explore the use of *Kernel Density*

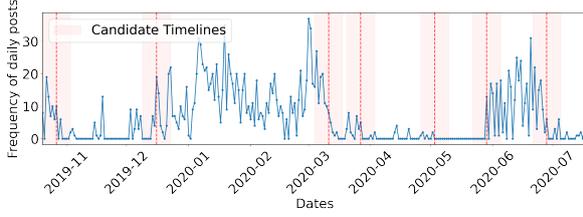


Figure 1: Using change-points in an example user’s posting behaviour to define candidate moments of change $M_u^{(c)}$ (dashed red line). Candidate timelines are then created centred on each $M_u^{(c)}$, with a radius $r=7$.

Estimation (KDE) (Rosenblatt, 1956; Scott, 2015) to estimate the probability density function of the user’s activity. For (b), we focus on time periods in the user’s history lasting at least 14 days during which the user had no activity (posts/comments) at all. If, given the past 90 days of a user’s activity, the probability on a particular day of seeing either (a) such a high volume of activity or (b) a long period of ‘silence’ is lower than .01, then we mark the start of this period as an ‘anomaly’. In both (a) and (b), we treat detected anomalies as CMoCs.

(3) Keywords: We further experiment with keyword-based methods based on the *suicide risk severity lexicon* (Gaur et al., 2019). Each keyword present in the lexicon corresponds to different levels of suicide risk severity such as “I’m tired of this suffering”, and “I’m going to kill myself”. We hypothesize that the presence of such phrases in a user’s post may be indicative of a MoC. The keywords-based methods we evaluate against simply return CMoCs for the timestamps of posts by a given user that contain a keyword from the full lexicon or a sub-lexicon.¹

(4) Random: “Random single day” is a baseline method we evaluate against, which selects a single date from a uniform distribution over all days in a user’s u posting history H_u as the CMoC $M_u^{(c)}$.

“Every day” is another baseline we experiment with, which simply returns every day as a CMoC. We experiment with it to see how well our methods are at avoiding the over-generation of candidate timelines. We seek to avoid over-generating timelines as we want to only return timelines with a high density of MoC, since this aligns with our goal of improving annotation efficiency.

¹Upon inspection of the phrases included in the sub-lexica, we excluded the “suicidal_indicator” sub-lexicon as it produced a lot of false positives.

3.2 Extracting Posts

Once a CMoC, $M_u^{(c)}$, is found, a span of timestamps $S_{i,u}$ from the user’s history H_u is then identified within a certain radius r^2 around $M_u^{(c)}$. Subsequently, we return the posts that are posted within the resulting time window as the candidate timeline, $T_{u,i}^{(c)}$. A candidate timeline therefore consists of the associated sequence of posts, corresponding timestamps and comments within $S_{i,u}$.

4 Evaluation of candidate timelines

Objective. We aim to identify the best method for extracting user timelines and also assess how good a given timeline is, while using minimal annotation resources. A good timeline is one that would contain a high proportion of posts that would be annotated as MoC, if manually labelled. As such, we define a good timeline selection method as one that is able to identify CMoC close to *dense regions* of Ground-Truth MoCs (GTMoCs) in an initial trial set of pre-annotated timelines.

4.1 Identifying dense regions in annotated timelines

Medoids. To represent the location of dense regions of GTMoCs, we propose the use of medoids. A medoid is a timestamp of a post, considered to be the centre of a cluster where the distances of all other timestamps of annotated posts in the timeline are minimal relative to it. In our work, medoids are computed for sets of labelled GTMoCs in annotated timelines. We therefore define a medoid $C_{u,i}^{(g)}$ as the timestamp in a timeline which is a GTMoC that has a minimal Euclidean distance $d(\cdot, \cdot)$ in time to all other annotated GTMoCs $M_{u,i}^{(g)}$ within the same annotated timeline $T_{u,i}^{(g)}$. The location of medoid $C_{u,i}^{(g)}$ in an annotated timeline is thus computed as:

$$C_{u,i}^{(g)} = \arg \min_{M_{u,i}^{(g)} \in T_{u,i}^{(g)}} \sum_{M_{u,j}^{(g)} \in T_{u,i}^{(g)}} d(M_{u,i}^{(g)}, M_{u,j}^{(g)}) \quad (1)$$

Density of annotated timelines. We aim to further characterise the locations of dense regions (medoids) by the number of GTMoC they contain. For this purpose we introduce a simple density metric which we assign to medoids. The density $\rho_{u,i}$

²Here we take $r = 7$ which gives a manageable amount of posts while providing context before and after the CMoC.

for a ground truth timeline is defined as:

$$\rho_{u,i} = \frac{|M_{u,i}^{(g)}|}{|p_{u,i}|} \quad (2)$$

, where $|M_{u,i}^{(g)}|$ is the sum total number of GTMoCs within timeline $T_{u,i}^{(g)}$ normalized by $|p_{u,i}|$, the sum total of posts within the same timeline.

In order to weight timelines by how dense they are, a medoid $C_{u,i}^{(g)}$ further inherits the density $\rho_{u,i}$ of the timeline $T_{u,i}^{(g)}$ it represents. We transform the raw density scores to provide a binary distinction between “dense” (+1) and “sparse” (-1) medoids as in equation 3:

$$\rho_{u,i}^{(\text{bin})} \begin{cases} +1 & \text{if } \rho_{u,i} \geq \text{Median}(\rho_{u,i} \forall u, i) \\ -1 & \text{otherwise} \end{cases} \quad (3)$$

A good timeline is therefore one that is identified as “dense” (+1) in equation 3, and the ideal location for a CMoC within it is as close to the medoid timestamp, defined in equation 1.

In an ideal scenario where we have the resources to annotate many timelines, sampled from many candidate methods – then it would be straightforward to compare and rank them based on the number of dense timelines or the average resulting density scores. This would allow us to directly identify the best method to select timelines in the future. However, due to the high cost and time-consuming process of annotation – we instead propose a few additional steps in our evaluation framework that allow us to identify alternative timeline selection methods without the need to annotate those timelines directly. We do this by proposing a scoring system based on distance scores of CMoC relative to dense medoids.

4.2 Scoring timeline selection methods

To assess how good a given method is for selecting desirable timelines, we make use of the evaluation framework in §4.1 to assess the quality of pre-annotated timelines against the CMoC in unannotated candidate timelines.

Assuming an annotated ground-truth timeline, $T_{u,i}^{(g)}$, we aim to assess how close an identified CMoC, $M^{(c)}$, is to a dense region of GTMoCs. Based on how we have defined good timelines, we therefore give preference to methods that identify CMoCs in close proximity to medoids that are

identified to be dense in GTMoC, while penalizing methods that over-generate CMoC. The reason for this is that we want to identify methods that are able to select timelines that will contain a high density of GTMoC when annotated, while avoiding methods that simply annotate the entire history of a user. The latter is infeasible and goes against our original aim of reducing the amount of data needed to be annotated by individuals.

Distance Scores To calculate the proximity of CMoCs to medoids, we compute the minimum absolute distance $d_{i,m}^{\text{min}}$ (in days) between all CMoCs detected by a given model m for a user’s u history H_u . Subsequently, we compute the following distance score metric per annotated timeline:

$$d_{i,m}^{(\text{score})} = (d_{i,m}^{(\text{min})} + \epsilon) * \text{sign}(\rho_{u,i}^{(\text{bin})})$$

where $\epsilon = 0.001$, to preserve the sign of each medoid’s $\rho_{u,i}^{(\text{bin})}$ in the case that $d_{i,m}^{(\text{min})} = 0$. The $d_{i,m}^{(\text{score})}$ is then used to denote the proximity of CMoCs generated by method m (in days) to a ground truth medoid $C_{u,i}^{(g)}$ with density $\rho_{u,i}^{(\text{bin})}$.

Since we want to generate timelines that are close to regions that are dense in terms of GTMoC, we aim for low positive scores of $d_{i,m}^{(\text{score})}$.

Voting procedure. We aim to reward methods that identify CMoC in close proximity to a “dense” ground-truth medoid (low positive $d_{i,m}^{(\text{score})}$), and penalize methods which over-generate CMoC - for example in locations that contain a low density of GTMoC. We thus assign votes to each method, to assess how well we achieve this objective.

Votes are assigned to each method, m for each computed distance score, $d_{i,m}^{(\text{score})}$, as follows:

$$v_{m,i} = \begin{cases} +1 & \text{if } 0 \leq d_{i,m}^{(\text{score})} \leq t_+ \\ 0 & \text{otherwise} \end{cases}$$

where t_+ is a threshold set to 10 days after experimentation. This score gives a positive vote to a method generating a CMoC that falls within a margin of t_+ days to a ground truth timeline. Setting a threshold is common in the field of change point detection (van den Burg and Williams, 2020). Votes, v , are then normalized per timeline and method:

$$v_{m,i}^{(\text{scaled})} = \frac{v_{m,i}}{|M_{u,i}^{(c)}|}$$

where $|M_{u,i}^{(c)}|$ is the total number of CMoCs generated by method m .

Scoring of methods. Timeline selection methods are subsequently scored and ranked by summing the votes $v_{m,i}^{(\text{scaled})}$ for each method m over all ground truth timelines, as shown in the results of table 1. The minimum score a given method can receive is 0, and scores can only be positive - while the maximum score is the total number of "dense" (+1) medoids in the dataset (190 in our case).

Comparison to Previous Work. Our evaluation of timeline selection methods differs from previous work on evaluating change-point detection methods, as we aim to compare distances to *regions* of changes (represented as medoids), rather than distances to *exact* change points (van den Burg and Williams, 2020). Typical measures for evaluating the identification of change points include: clustering metrics - such as the segmentation covering metric (used traditionally in image segmentation (Everingham et al., 2010; Arbelaez et al., 2010)), and classification metrics such as F1 scores as described in (van den Burg and Williams, 2020). Similar to our proposed approach, these metrics capture whether the distance of a predicted change-point to a ground-truth change-point falls within a certain threshold (van den Burg and Williams, 2020).

Our evaluation framework depends on a set of timelines manually annotated with GTMoC. The manually annotated timelines were selected on the basis of a particular method (here BOCPD). While including the method that selected the timelines for manual annotation in the evaluation of methods for generating CMoC and new timelines may appear biased, note that it is theoretically possible for another method to get a higher score. This is because the criteria for manual annotation of GTMoC are different to the assignment of CMoC by the methods. As a result not all annotated timelines are "dense". If a candidate selection method would only return CMoC close to regions where the manually annotated timelines had a high density of GTMoC, it would receive better distance scores and more votes than the method which originally selected the timelines for annotation. This is because the method which originally generated the timelines for annotation would be penalized for predicting a CMoC close to an sparse timeline, annotated with very few GTMoC.

Another advantage of our evaluation setting is that if an alternative method identifies a CMoC

towards the end, or slightly outside, a manually annotated timeline - there is the potential that the resulting candidate timeline will contain a higher density of GTMoC if annotated. In such a scenario the alternative method has the potential to receive better distance scores as it may select a timeline closer to a dense region of GTMoC, if this exists near the edges of the originally manually annotated timeline. Thus our distance scores can potentially help us identify better methods for timeline extraction than the method originally used to select the timelines for manual annotation.

5 Experiments

We empirically evaluate our proposed timeline selection methods (§3), using our proposed evaluation framework (§4) based on ground-truth human annotated data.

5.1 Dataset

We licensed a de-identified dataset from TalkLife³ consisting of 1.1 million users, resulting in 12.3 million posts between August 2011 to August 2020, of which we sampled from based on methods described in §3 to create timelines.

Due to high variance in posting frequency of users, we chose to annotate only timelines that had between 10 and 150 posts - so that there was sufficient amount of context to annotators to assess whether a change had occurred, and that the timelines were not impractically long. The final annotated dataset includes 500 timelines from 500 separate TalkLife users, consisting of 18,702 posts in total where the mean number of posts per timeline is $\mu = 35 \pm 22$. The 500 timelines were selected using a BOCPD Poisson-Gamma model, where the parameters ($\alpha_0:0.01; \beta_0:10; h_0:10^3$) were fixed on the basis of improved model performance compared to 70 initial manually annotated validation timelines which had been generated using the anomaly detection (high activity: posts) method (§3). All timelines within this dataset were manually inspected and filtered according to the details in appendix A.1.

TalkLife is a free-to-use global peer-support social network platform operating primarily as a mobile app. Users are mainly English speakers, where 70% of them are in the age range of 15 to 24 (Sharma et al., 2020a). We chose data from TalkLife for

³<https://www.talklife.com>

538 this work since the content across the entire plat-
 539 form is focused on conversations around mental-
 540 health and daily-life issues and feelings. It is thus
 541 suited to identifying MoC, and is complementary
 542 to recent work which uses TalkLife data for compu-
 543 tationally analysing mental health (Pruksachatkun
 544 et al., 2019; Sharma et al., 2020b; Saha and Sharma,
 545 2020; Kim et al., 2021).

546 TalkLife users make textual posts and others on
 547 the platform may comment on them. While there
 548 are several features available from TalkLife, we
 549 propose to select timelines on the basis of only the
 550 frequency of posts made by users and frequency
 551 of associated comments received, all of which are
 552 timestamped. The context of the posts is only used
 553 in the manual annotation of selected timelines with
 554 GTMoC, used for evaluation. As such, the methods
 555 proposed in this paper are transferable to other
 556 popular platforms such as Twitter and Reddit for
 557 creating timelines for dataset annotation.

558 5.1.1 Annotation Guidelines for GTMoC

559 After extracting timelines (§3) from TalkLife, these
 560 were annotated by 3 English speaking, university
 561 educated annotators (one of them being a native
 562 speaker). Annotation was performed using guide-
 563 lines and an associated annotation interface pro-
 564 posed by (anonymous). The process is described
 565 briefly in this subsection.

566 Annotators were provided with timelines, con-
 567 taining sequences of time-stamped posts by users
 568 along with comments made on those posts. Annota-
 569 tors were asked to label posts containing a "Switch"
 570 (sudden change in mood) or an "Escalation" (grad-
 571 ual mood progression). A label of "None", the
 572 default, is assigned to posts with no MoC. Specif-
 573 ically, a "Switch" is defined in the guidelines as
 574 "a drastic change in mood, in comparison with the
 575 recent past", and the annotator is tasked to label
 576 the first post which has a clearly different mood
 577 compared to previous posts. They are also asked to
 578 specify the duration of the change in mood in terms
 579 of the associated range of posts. An "Escalation"
 580 on the other-hand is defined as "a gradual change in
 581 mood, which should last for a few posts". Annota-
 582 tions are provided for the peak of the escalation and
 583 the range of associated posts (both before and after
 584 the identified peak in mood change). For this paper
 585 we consider all labels of "Switch", "Escalation",
 586 and their corresponding ranges as GTMoC. For the
 587 annotation of GTMoC, posts within timelines were
 588 displayed on a longitudinal basis, thus providing

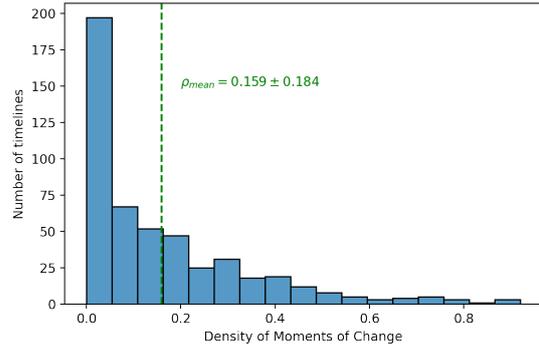


Figure 2: Histogram showing the density of GTMoCs per timeline. All 500 timelines were selected using BOCPD PG ($\alpha_0:0.01$; $\beta_0:10$; $h_0:10^3$).

589 annotators with access to both previous and future
 590 context around each post in the timeline.

591 To obtain GTMoC for our evaluation we aggre-
 592 gate the annotations across all annotators per time-
 593 line as described in (anonymous). The percent of
 594 inter-annotator agreement for the labels "None",
 595 "Switch" and "Escalation" were 0.89, 0.30, and
 596 0.50 respectively based on majority agreement.

597 5.2 Results & Discussion

598 We identify CMoC from the timeline selection
 599 methods in §3.1 on the 500 users for whom we
 600 have GTMoCs, and evaluate these using our ap-
 601 proach in 4. To compare different methods, we
 602 also round all CMoC to the nearest day removing
 603 duplicate predicted dates per method.

604 **Density scores of annotated timelines.** The den-
 605 sity of the final annotated timelines, selected by
 606 our best performing selection method are presented
 607 in Fig. 2. With a mean density of 0.159, this is
 608 comparatively high considering that GTMoCs are
 609 rare events and that many timelines typically do
 610 not contain any GTMoC when annotated.

611 **Ranking of timeline selection methods.** Table 1
 612 shows that BOCPD with a high h_0 and low α_0/β_0
 613 produces overall timelines closest to the GTMoCs.
 614 Thus this model, which is confident about a low
 615 number of CMoCs, will generate fewer CMoCs and
 616 corresponding timelines. BOCPD is followed by
 617 a standard approach to selecting timelines, which
 618 is to impose a linguistic bias on the user posts and
 619 therefore produce annotated datasets (and hence,
 620 models) based on the presence of certain keywords.
 621 Note that these methods achieve less than half the
 622 top score of BOCPD. Even with a low h_0 and
 623 $\alpha_0/\beta_0 = 1$ (more likely to over-generate CMoCs)

the BOCPD still outperforms most of the anomaly detection methods and the random timeline generation, where a day is chosen at random in a user’s timeline with seven days around it. The anomaly detection method which identifies CMOCs at points of high activity of posts performs similarly to the keyword based methods. All other anomaly detection methods seem to over-generate CMOCs with ones identifying anomalies on low user activity performing worse than the random timeline generation. The floor score for over-generation of CMOCs is provided by considering every day as a CMOC.

Method	Score
BOCPD PG ($\alpha_0: .01; \beta_0: 10; h_0: 10^3$)	27.34
Keywords (three categories)	13.79
Keywords (all)	12.35
AD (high activity: posts)	10.09
BOCPD PG ($\alpha_0: 1; \beta_0: 1; h_0: 10$)	9.83
AD (high & low activity: posts)	8.20
AD (high activity: comments received)	5.21
AD (high & low activity: comments received)	4.92
Random single day	4.00
AD (low activity: comments received)	3.49
AD (low activity: posts)	3.28
Every day	0.25

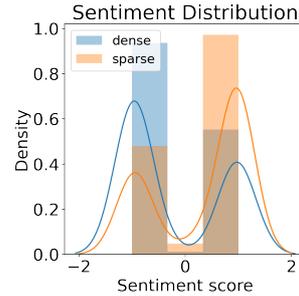
Table 1: Methods (proposed in §3) ranked in descending order by their ability to generate desired timelines, using our evaluation framework in §4.

Linguistic analysis of timelines. To gain some insights into the characteristics of ‘dense’ (high number of GTMoCs) vs ‘sparse’ timelines (low number of GTMoCs), we employ VADER (Hutto and Gilbert, 2014) for sentiment and ‘Twitter-RoBERTa-emotion’ (Barbieri et al., 2020) for emotion recognition⁴ on the post-level of 250 timelines, equally split between ‘dense’ (density score $\rho_{u,i}$ is in upper-quartile of all timelines) and ‘sparse’ (bottom-quartile). The distribution of sentiment scores across these timelines are shown in Fig. 3. For each timeline we extract statistical features (avg, std, min, max) for each emotion/sentiment dimension of the posts within it, and the same statistical features based on their difference across two consecutive posts within the timeline. Using these features, we train a Logistic Regression aiming at predicting ‘dense’ vs ‘sparse’ timelines and extract the coefficients with the highest/lowest values.

Table 2 suggests that sparse timelines frequently

⁴We use the compound sentiment score from VADER, assigning a single sentiment score to each post; Twitter-RoBERTa-emotion assigns one score per emotion: joy, anger, sadness, optimism.

consist of positive posts in sentiment and mood. On the other hand, sadness- and variance-based features have the most positive relationship with predicting a timeline containing many MoCs – a finding that was also empirically confirmed via manual inspection of the most dense timelines. This suggests that future work could also employ methods based on mood or sentiment for extracting user timelines (with the cost of introducing linguistic bias), while highlighting the importance for considering the variation of a user’s mood and sentiment.



Feature	Coef
sadness (avg)	2.29
sadness (std)	1.45
sentiment (std)	1.00
sentiment (avg)	-1.23
optimism (avg)	-1.25
sentiment (min)	-1.31
joy (avg)	-1.58

Figure 3: Distribution of sentiment scores of ‘dense’ vs ‘sparse’ timelines (medians: $-.949$ & $.970$, respectively). Table 2: Coefficients of Logistic Regression trained to classify a timeline as ‘dense’ (1) or ‘sparse’ (-1).

6 Conclusions & Future work

We have introduced methods and an evaluation framework for identifying timelines with many Moments of Change (MoC) in a user’s posting behaviour on social media. Our aim is to use changes in posting behaviour as a proxy for changes in mood, to facilitate the process and maximise the effectiveness of annotation of longitudinal user content. Our methods have been manually evaluated against ground truth MoCs (GTMoCs). Bayesian Online Change Point Detection (BOCPD) with a Poisson-Gamma model shows promise in detecting candidate MoCs close to GTMoCs.

In future work we will explore the incorporation of textual content in the BOCPD Poisson-Gamma model for the distinction between different types of GTMoC. We find that resulting timelines dense in GTMoCs are characterised by a high deviation in sentiment from one post to the next, suggesting that such deviations may be a useful feature for distinguishing between different types of GTMoC.

We expect that the methods proposed in our work will benefit researchers interested in creating longitudinally annotated textual datasets consisting of user posts, particularly when annotating Moments of Change.

7 Ethics Statement

Ethics IRB approval was obtained from the corresponding ethics board of the host University prior to engaging in this research study. Our work involves ethical considerations around the analysis of user generated content shared on a peer support network (TalkLife). A license was obtained to work with the user data from TalkLife and a project proposal was submitted to them in order to embark on the project. The current paper focuses on the identification of periods of interest within the user history, in terms of moments of change. The work on annotation of moments of change (MoC) is separate to this paper but considers sudden shifts in mood (switches or escalations). Annotators were given contracts and paid fairly in line with University pay-scales. They were alerted about potentially encountering disturbing content and advised to take breaks during annotation. The annotations are used to evaluate the work of the current paper, which aims to meaningfully segment timelines in terms of containing likely moments of change. Potential risks from the application of our work in being able to identify moments of change in individuals' timelines are akin to the identification of those in earlier work on personal event identification from social media and the detection of suicidal ideation. Potential mitigation strategies include restricting access to the code base and annotation labels used for evaluation. No data can be shared without permission from the platform or significantly paraphrased. Any examples used from the users' history are anonymised and paraphrased.

References

Ryan Prescott Adams and David J. C. MacKay. 2007. [Bayesian Online Change-point Detection](#). *arXiv:0710.3742 [stat]*. ArXiv: 0710.3742.

anonymous. anonymous.

Jeffery Ansah, Lin Liu, Wei Kang, Selasie Kwashie, Jixue Li, and Jiuyong Li. 2019. A graph is worth a thousand words: Telling event stories using timeline summarization graphs. In *The World Wide Web Conference*, pages 2565–2571.

Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. 2010. Contour detection and hierarchical image segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 33(5):898–916.

Francesco Barbieri, Jose Camacho-Collados, Leonardo Neves, and Luis Espinosa-Anke. 2020. Tweet-eval: Unified benchmark and comparative eval-

uation for tweet classification. *arXiv preprint arXiv:2010.12421*.

Lei Cao, Huijun Zhang, Ling Feng, Zihan Wei, Xin Wang, Ningyun Li, and Xiaohao He. 2019. Latent suicide risk detection on microblog via suicide-oriented word embeddings and layered attention. *arXiv preprint arXiv:1910.12038*.

Yi Chang, Jiliang Tang, Dawei Yin, Makoto Yamada, and Yan Liu. 2016a. Timeline summarization from social media with life cycle models. In *IJCAI*, pages 3698–3704.

Yi Chang, Makoto Yamada, Antonio Ortega, and Yan Liu. 2016b. [Lifecycle Modeling for Buzz Temporal Pattern Discovery](#). *ACM Transactions on Knowledge Discovery from Data*, 11(2):1–24.

Xiuying Chen, Zhangming Chan, Shen Gao, Meng-Hsuan Yu, Dongyan Zhao, and Rui Yan. 2019. Learning towards abstractive timeline summarization. In *IJCAI*, pages 4939–4945.

Qijin Cheng, Tim MH Li, Chi-Leung Kwok, Tingshao Zhu, and Paul SF Yip. 2017. Assessing suicide risk and emotional distress in chinese social media: a text mining and machine learning study. *Journal of medical internet research*, 19(7):e243.

Glen Coppersmith, Mark Dredze, and Craig Harman. 2014. Quantifying mental health signals in twitter. In *Proceedings of the workshop on computational linguistics and clinical psychology: From linguistic signal to clinical reality*, pages 51–60.

Glen Coppersmith, Ryan Leary, Patrick Crutchley, and Alex Fine. 2018. Natural language processing of social media as screening for suicide risk. *Biomedical informatics insights*, 10:1178222618792860.

Glen Coppersmith, Kim Ngo, Ryan Leary, and Anthony Wood. 2016. Exploratory analysis of social media prior to a suicide attempt. In *Proceedings of the third workshop on computational linguistics and clinical psychology*, pages 106–117.

Daryl J Daley and David Vere-Jones. 2003. *An introduction to the theory of point processes: volume I: elementary theory and methods*. Springer.

Munmun De Choudhury, Emre Kiciman, Mark Dredze, Glen Coppersmith, and Mrinal Kumar. 2016. [Discovering Shifts to Suicidal Ideation from Mental Health Content in Social Media](#). In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, pages 2098–2110, San Jose California USA. ACM.

Nan Du, Hanjun Dai, Rakshit Trivedi, Utkarsh Upadhyay, Manuel Gomez-Rodriguez, and Le Song. 2016. [Recurrent Marked Temporal Point Processes: Embedding Event History to Vector](#). In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1555–1564, San Francisco California USA. ACM.

Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2010. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338.

Paul Fearnhead and Zhen Liu. 2007. On-line inference for multiple changepoint problems. *Journal*

809			
810		<i>of the Royal Statistical Society: Series B (Statistical Methodology)</i> , 69(4):589–605.	
811	Manas Gaur, Amanuel Alambo, Joy Prakash Sain,		
812	Ugur Kursuncu, Krishnaprasad Thirunarayan, Ra-		
813	makanth Kavuluru, Amit Sheth, Randy Welton, and		
814	Jyotishman Pathak. 2019. Knowledge-aware assess-		
815	ment of severity of suicide risk for early interven-		
816	tion. In <i>The World Wide Web Conference</i> , pages 514–525.		
817	George Gkotsis, Anika Oellrich, Sumithra Velupillai,		
818	Maria Liakata, Tim JP Hubbard, Richard JB Dob-		
819	son, and Rina Dutta. 2017. Characterisation of men-		
820	tal health conditions in social media using informed		
821	deep learning. <i>Scientific reports</i> , 7(1):1–11.		
822	Po Hu, Minlie Huang, Peng Xu, Weichang Li,		
823	Adam K Usadi, and Xiaoyan Zhu. 2011. Generating		
824	breakpoint-based timeline overview for news topic		
825	retrospection. In <i>2011 IEEE 11th International Con-</i>		
826	<i>ference on Data Mining</i> , pages 260–269. IEEE.		
827	Clayton Hutto and Eric Gilbert. 2014. Vader: A par-		
828	simonious rule-based model for sentiment analysis		
829	of social media text. In <i>Proceedings of the Interna-</i>		
830	<i>tional AAAI Conference on Web and Social Media</i> ,		
831	volume 8.		
832	Seunghyun Kim, Afsaneh Razi, Gianluca Stringhini,		
833	Pamela Wisniewski, and Munmun De Choudhury.		
834	2021. You don’t know how i feel: Insider-outsider		
835	perspective gaps in cyberbullying risk detection. In		
836	<i>Proceedings of the International AAAI Conference</i>		
837	<i>on Web and Social Media</i> , volume 15, pages 290–		
838	302.		
839	Jeremias Knoblauch and Theodoros Damoulas. 2018.		
840	Spatio-temporal Bayesian on-line changepoint de-		
841	tection with model selection . In <i>Proceedings of the</i>		
842	<i>35th International Conference on Machine Learning</i> ,		
843	volume 80 of <i>Proceedings of Machine Learning Re-</i>		
844	<i>search</i> , pages 2718–2727. PMLR.		
845	Jiwei Li and Claire Cardie. 2014. Timeline genera-		
846	tion: tracking individuals on twitter . In <i>Proceed-</i>		
847	<i>ings of the 23rd international conference on World</i>		
848	<i>wide web - WWW ’14</i> , pages 643–652, Seoul, Korea.		
849	ACM Press.		
850	Jiwei Li, Alan Ritter, Claire Cardie, and Eduard Hovy.		
851	2014. Major life event extraction from twitter based		
852	on congratulations/condolences speech acts. In <i>Pro-</i>		
853	<i>ceedings of the 2014 conference on empirical meth-</i>		
854	<i>ods in natural language processing (EMNLP)</i> , pages		
855	1997–2007.		
856	Haitao Lin, Cheng Tan, Lirong Wu, Zhangyang Gao,		
857	and Stan Z Li. 2021. An empirical study: Extensive		
858	deep temporal point process. <i>arXiv e-prints</i> , pages		
859	arXiv–2110.		
860	Kate Loveys, Patrick Crutchley, Emily Wyatt, and Glen		
861	Coppersmith. 2017. Small but mighty: affective mi-		
862	cropatterns for quantifying mental health from so-		
863	cial media language. In <i>Proceedings of the fourth</i>		
864	<i>workshop on computational linguistics and clinical</i>		
865	<i>Psychology—From linguistic signal to clinical real-</i>		
866	<i>ity</i> , pages 85–95.		
867	Sebastian Martschat and Katja Markert. 2017. Improv-		
868	ing ROUGE for Timeline Summarization . In <i>Pro-</i>		
869	<i>ceedings of the 15th Conference of the European</i>		
		<i>Chapter of the Association for Computational Lin-</i>	870
		<i>guistics: Volume 2, Short Papers</i> , pages 285–290,	871
		Valencia, Spain. Association for Computational Lin-	872
		guistics.	873
	Sebastian Martschat and Katja Markert. 2018. A		874
	Temporally Sensitive Submodularity Framework for		875
	Timeline Summarization . In <i>Proceedings of the</i>		876
	<i>22nd Conference on Computational Natural Lan-</i>		877
	<i>guage Learning</i> , pages 230–240, Brussels, Belgium.		878
	Association for Computational Linguistics.		879
	Naoki Masuda, Issei Kurahashi, and Hiroko Onari.		880
	2013. Suicide ideation of individuals in online so-		881
	cial networks. <i>PLoS one</i> , 8(4):e62262.		882
	Matthew Matero, Akash Idnani, Youngseo Son, Sal-		883
	vatore Giorgi, Huy Vu, Mohammad Zamani, Parth		884
	Limbachiya, Sharath Chandra Guntuku, and H An-		885
	drew Schwartz. 2019. Suicide risk assessment with		886
	multi-level dual-context language and bert. In <i>Pro-</i>		887
	<i>ceedings of the Sixth Workshop on Computational</i>		888
	<i>Linguistics and Clinical Psychology</i> , pages 39–44.		889
	Rohan Mishra, Pradyumn Prakhar Sinha, Ramit Sawh-		890
	ney, Debanjan Mahata, Puneet Mathur, and Rajiv		891
	Ratn Shah. 2019. SNAP-BATNET: Cascading Au-		892
	thor Profiling and Social Network Graphs for Sui-		893
	cide Ideation Detection on Social Media . In <i>Pro-</i>		894
	<i>ceedings of the 2019 Conference of the North Amer-</i>		895
	<i>ican Chapter of the Association for Computational</i>		896
	<i>Linguistics: Student Research Workshop</i> , pages 147–		897
	156, Minneapolis, Minnesota. Association for Com-		898
	putational Linguistics.		899
	Martha Neary and Stephen M Schueller. 2018. State		900
	of the field of mental health apps. <i>Cognitive and</i>		901
	<i>Behavioral Practice</i> , 25(4):531–537.		902
	Alexandra Olteanu, Carlos Castillo, Fernando Diaz,		903
	and Emre Kiciman. 2019. Social data: Bi-		904
	ases, methodological pitfalls, and ethical boundaries.		905
	<i>Frontiers in Big Data</i> , 2:13.		906
	Yada Pruksachatkun, Sachin R. Pendse, and Amit		907
	Sharma. 2019. Moments of change: Analyzing peer-		908
	based cognitive support in online mental health fo-		909
	rums . In <i>Proceedings of the 2019 CHI Conference</i>		910
	<i>on Human Factors in Computing Systems</i> , CHI ’19,		911
	page 1–13, New York, NY, USA. Association for		912
	Computing Machinery.		913
	Marian-Andrei Rizoioiu, Young Lee, Swapnil Mishra,		914
	and Lexing Xie. 2017. A tutorial on hawkes pro-		915
	cesses for events in social media. <i>arXiv preprint</i>		916
	<i>arXiv:1708.06401</i> .		917
	Murray Rosenblatt. 1956. Remarks on Some Nonpara-		918
	metric Estimates of a Density Function . <i>The Annals</i>		919
	<i>of Mathematical Statistics</i> , 27(3):832 – 837.		920
	Koustuv Saha and Amit Sharma. 2020. Causal factors		921
	of effective psychosocial outcomes in online men-		922
	tal health communities. In <i>Proceedings of the Interna-</i>		923
	<i>tional AAAI Conference on Web and Social Media</i> ,		924
	volume 14, pages 590–601.		925
	Ramit Sawhney, Harshit Joshi, Lucie Flek, and Ra-		926
	ajiv Shah. 2021. Phase: Learning emotional phase-		927
	aware representations for suicide ideation detection		928
	on social media. In <i>Proceedings of the 16th Con-</i>		929
	<i>ference of the European Chapter of the Association</i>		930

931 *for Computational Linguistics: Main Volume*, pages
932 2415–2428.

933 Ramit Sawhney, Harshit Joshi, Saumya Gandhi, and
934 Rajiv Shah. 2020. A time-aware transformer based
935 model for suicide ideation detection on social media.
936 In *Proceedings of the 2020 Conference on Empirical
937 Methods in Natural Language Processing (EMNLP)*,
938 pages 7685–7697.

939 David W Scott. 2015. *Multivariate density estimation:
940 theory, practice, and visualization*. John Wiley &
941 Sons.

942 Ashish Sharma, Monojit Choudhury, Tim Althoff, and
943 Amit Sharma. 2020a. [Engagement Patterns of Peer-
944 to-Peer Interactions on Mental Health Platforms](#).
945 *Proceedings of the International AAAI Conference
946 on Web and Social Media*, 14:614–625.

947 Ashish Sharma, Adam Miner, David Atkins, and Tim
948 Althoff. 2020b. [A Computational Approach to Un-
949 derstanding Empathy Expressed in Text-Based Men-
950 tal Health Support](#). In *Proceedings of the 2020 Con-
951 ference on Empirical Methods in Natural Language
952 Processing (EMNLP)*, pages 5263–5276, Online. As-
953 sociation for Computational Linguistics.

954 Oleksandra Shchur, Ali Caner Türkmen, Tim
955 Januschowski, and Stephan Günnemann. 2021.
956 Neural temporal point processes: A review. *ArXiv*,
957 abs/2104.03528.

958 Han-Chin Shing, Suraj Nair, Ayah Zirikly, Meir
959 Friedenber, Hal Daumé III, and Philip Resnik.
960 2018. Expert, crowdsourced, and machine assess-
961 ment of suicide risk via online postings. In *Proce-
962 edings of the Fifth Workshop on Computational Lin-
963 guistics and Clinical Psychology: From Keyboard
964 to Clinic*, pages 25–36.

965 Julius Steen and Katja Markert. 2019. [Abstractive
966 Timeline Summarization](#). In *Proceedings of the
967 2nd Workshop on New Frontiers in Summarization*,
968 pages 21–31, Hong Kong, China. Association for
969 Computational Linguistics.

970 Russell Swan and James Allan. 2000. Automatic gen-
971 eration of overview timelines. In *Proceedings of
972 the 23rd annual international ACM SIGIR confer-
973 ence on Research and development in information
974 retrieval*, pages 49–56.

975 Gerrit JJ van den Burg and Christopher KI Williams.
976 2020. An evaluation of change point detection algo-
977 rithms. *arXiv preprint arXiv:2003.06222*.

978 Shang Wang, Zhiwei Yang, and Yi Chang. 2021. Bring-
979 ing order to episodes: Mining timeline in social me-
980 dia. *Neurocomputing*, 450:80–90.

981 Qiang Zhang, Aldo Lipani, Omer Kirnap, and Em-
982 ine Yilmaz. 2020. [Self-Attentive Hawkes Processes](#).
983 *arXiv:1907.07561 [cs, stat]*. ArXiv: 1907.07561.

984 Simiao Zuo, Haoming Jiang, Zichong Li, Tuo Zhao,
985 and Hongyuan Zha. 2020. Transformer hawkes
986 process. In *International Conference on Machine
987 Learning*, pages 11692–11702. PMLR.

A Appendix 988

A.1 Creating Ground-truth Timelines, by Retaining a Subset of Representative Candidate Timelines 989

990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038

In addition to the details provided in section 3, for selecting candidate timelines, we provide some additional details inline below. As multiple timelines will typically be returned for each user using methods in 3 and annotating all of these can be time-consuming, in order to keep the 500 annotated ground-truth timelines relatively diverse in terms of the types of users - only a single timeline was returned per user to be annotated. Therefore, for each user only a single timeline was randomly sampled per user and these were presented visually in turn to the first author of this paper, with multiple time-scales limiting the x-axis of the visualization returned: (1) the time-scale of the whole user’s history, (2) a radius of 200 days surrounding the CMoC and (3) a radius of 31 days around the CMoC. This was to ensure that the candidate timelines could be inspected in close detail (3), and also observing the timeline in context of the full time-series (1) for that user. These three multiple time-scales for a single user are presented visually in figure 4. A manual binary decision was then made on whether to discard this timeline or retain it to be annotated and thereby create a ground-truth timeline using it. This decision was based on a time-series visualization of the frequency of daily posts for that user and highlighting the location of the timeline to be either retained or discarded. The decision to discard a timeline was based on two criteria: whether the timeline (1) was primarily sparse over the full 15 days of the timelines, or to a lesser degree (2) whether it appeared that the CMoC was generated by noise. It was chosen to discard timelines that were (1) primarily sparse, to ensure that we allow sufficient amount of time to pass between posts such that moments of change can occur. Timelines that appeared to be (2) generated by noise, were discarded such that the ground-truth timelines were representative of timelines that would be generated by a change-point detection algorithm with well chosen hyper-parameters - as the retained timelines were thus timelines that appeared to be generated by realistic change-points. Figure 5 presents a visualisation of a timeline that was discarded as described above, and figure 4 describes a timeline that was included to be annotated as a ground-truth timeline.

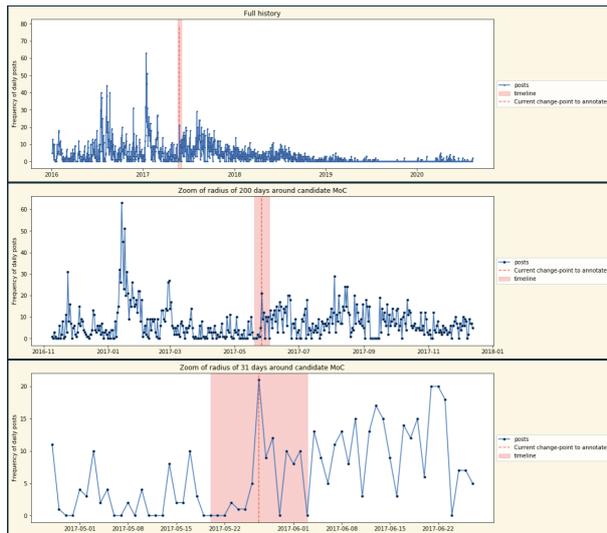


Figure 4: A timeline that was retained, out of the 1,220 timelines manually observed. It was retained as it (1) was not primarily sparse as it contains posts distributed well over the timeline, and (2) appeared to be generated by a plausible change-point rather than noise. Timelines were visualized on 3 time-scales, as shown in this figure, to allow for closer inspection and to compare in context of the full time-series.

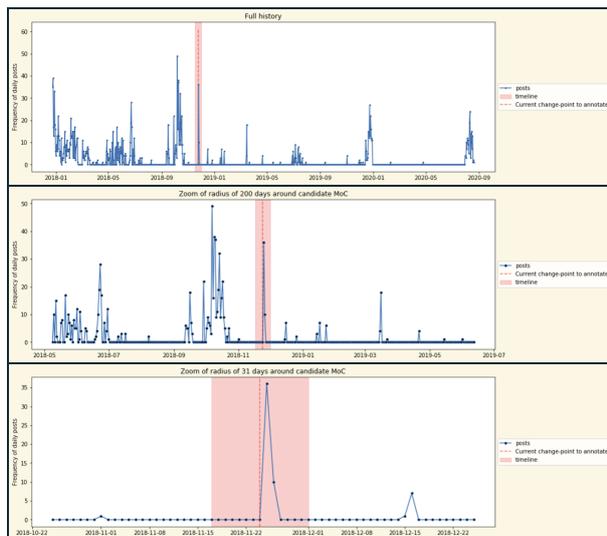


Figure 5: A timeline that was discarded, out of the 1,220 timelines manually observed. It was discarded as it (1) was primarily sparse containing only posts on a few days in the timeline, and (2) appeared to be generated by noise rather than by a realistic change-point.

From the annotated timelines, medoids are returned as the medoid timestamp of the annotated GTMoC after annotations were union aggregated across all annotators as described in (anonymous).

1046
1047
1048
1049

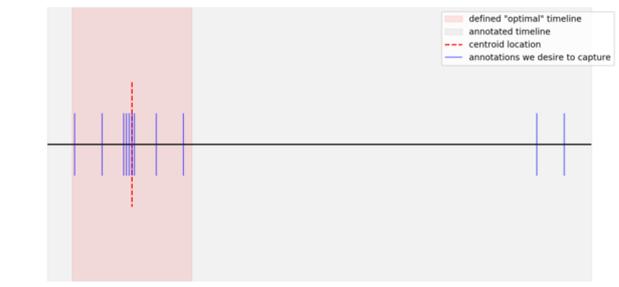


Figure 6: Identifying the position of the medoid, from the timestamps of posts annotated as GTMoCs.

A.2 Annotation Guidelines

1050

The annotation task proposed by (anonymous) was to assign annotators to identify changes in mood, by reading through the posts in chronological order included within the generated timeline of an individual - and annotating the posts which contain a change in the user's mood compared to the recent past.

1051
1052
1053
1054
1055
1056
1057

An example illustrating both a switch, and an escalation are displayed in figure 7. Note, that the example shown in this figure will be paraphrased before the work is published - to further preserve anonymity of this user.

1058
1059
1060
1061
1062

This process of visually deciding whether a randomly sampled candidate timeline should be retained to be converted into a ground-truth timeline was repeated until 500 candidate timelines were retained. This process thus lasted until 1,220 randomly sampled timelines were observed and thus 720 timelines were discarded.

1039
1040
1041
1042
1043
1044
1045

3.3. I feel good today | stopped procrastinating and did smth productive and now i just wanna sleep

SHOW/Hide CONVERSATIONS

Friday, 14 Feb 2020

4.1. I can't sleep.

SHOW/Hide CONVERSATIONS

4.2. I hate myself so much for not having the will to even get up of my bed this day cause i feel like fucking burden

SHOW/Hide CONVERSATIONS

4.3. For the people who don't have a valentine date and are sad just buy chocolate and flowers for yourself u fucking deserve it

Goodnight

SHOW/Hide CONVERSATIONS

4.4. I'm skipping school tomorrow I'm paying money but i think i will feel worse if i go i haven't even done my home works so I'm leaning towards skipping

SHOW/Hide CONVERSATIONS

Saturday, 15 Feb 2020

5.1. I'm useless and a disappointment

SHOW/Hide CONVERSATIONS

5.2. I'm feeling pretty shitty these days...

SHOW/Hide CONVERSATIONS

Figure 7: An example of the annotation interface, displaying a sequence of posts in a timeline shown to an annotator. For these sequence of posts, the annotator annotated a single post as a "switch" and another post as an "escalation". The user has a "switch" at 4.1, drastically changing from a positive mood to a negative mood - where this changed mood persists until 4.4. The "escalation" begins and is at its peak (in this case becoming increasingly negative) at 5.1, and de-escalates up to the post at 5.2."