# TSExplain: Explaining Aggregated Time Series by Surfacing Evolving Contributors

## Technical Report[*]

Yiru Chen
yiru.chen@columbia.edu
Columbia University
New York, NY, USA

Silu Huang
silu.huang@microsoft.com
Microsoft Research
Redmond, WA, USA

## ABSTRACT

Aggregated time series can be generated effortlessly everywhere, e.g., "total confirmed covid-19 cases since 2019", "S&P500 during the year 2020", and "total liquor sales over time". Understanding "how" and "why" these key performance indicators(KPI) evolve over time is critical to making data-informed decisions. Existing explanation engines focus on explaining the difference between two relations. However, this falls short of explaining KPI's continuous changes over time, as it overlooks explanations in between by only looking at the two endpoints. Motivated by this, we propose TSExplain, a system that explains aggregated time series by surfacing the underlying *evolving* top contributors. Under the hood, we leverage the existing work on *two-relations diff* as a building block and formulate a *K-Segmentation* problem to segment the time series such that each segment after segmentation shares consistent *explanations*, i.e., contributors. To quantify consistency in each segment, we propose a novel within-segment variance design based on top explanations; to derive the optimal *K-Segmentation* scheme, we develop a dynamic programming algorithm. Experiments on synthetic and real-world datasets show that our explanation-aware segmentation can effectively identify evolving explanations for aggregated time series and outperform explanation-agnostic segmentation. Further, we proposed an optimal selection strategy of $K$ and several optimizations to speed up TSExplain for interactive user experience, achieving up to 13× efficiency improvement.

## 1 INTRODUCTION

Time series data is gaining increasing popularity these days across sectors ranging from finance, retail, IoT to DevOps. Time series analysis is crucial for uncovering insights from time series data and helping business users make data-informed decisions. A business analyst typically focuses on three questions: "what happened" to understand the changes in key performance indicator (KPI), "why happened" to reason why KPI changes, and "now what" [3] to guide what actions should be taken. "What" questions have been extensively studied both academia-wise [12] and industry-wise [11, 31, 43]. "Why" questions are starting to attract wide attentions [2, 47, 51]. Existing explanation engines focus on explaining (1) one aggregated value [11, 19] or (2) differences between two given relations: a test relation and a control relation [1, 2, 11, 16, 27, 30, 33, 37–39, 41, 44, 47, 51].

However, KPIs are typically monitored continuously, reporting some aggregated time series. Simply explaining one aggregated value overlooks the trend of time series, e.g. "why up/down"; merely



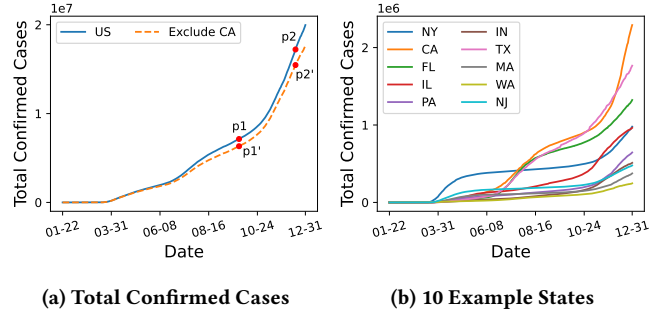**(a) Total Confirmed Cases**     **(b) 10 Example States**

**Figure 1: COVID-19 Total Confirmed Cases [20]**

focusing on its two endpoints and explaining their difference overlooks the evolving explanations in between. As evidence, although `key influencer` feature which explains the difference between two given relations is well received in PowerBI community, a highly voted feature request in PowerBI Idea Forum is called *key influencer*[1] *over time* [34]. Below are some quotes from the user: *"the mix of factors/influencers tends to be more dynamic than static over time...It would be nice to add a time dimension to the Key Influencers analysis to understand how the top key influencers evolve over the evaluated period" [34]*. In time-series scenarios, it is often more desired to explain the evolving dynamics over time than only to consider two end timestamps.

**Disclaimer.** We remark that identifying the root cause of "why" questions in general is only plausible when combining human interpretations with tools. Quoted from Tableau `ExplainData` homepage [44]: "The tool uncovers and describes relationships in your data. It can't tell you what is causing the relationships or how to interpret the data." Following the literature [1, 2, 47, 51], the `explanation` in our work does not equate to "cause", instead it corresponds to the data slice that contributes most to the overall change as we will describe in our Background Section (Definition 3.1).

**Motivating Examples.** We now describe three application scenarios of explaining changes in continuously evolving KPIs.

• `COVID-19`. Figure 1a depicts the total number of covid-19 confirmed cases during year 2020. This time series is obtained by performing a groupby-aggregate query over the original table [20], which consists of attributes like `state`, `total_confirmed_cases`, and `daily_confirmed_cases`. By looking at Figure 1a, users can get an understanding of how overall trend evolves over time. One natural

---

[1]Influencer corresponds to explanation in our notion.

followup question is "what makes the increase" – how different states contribute to the increase as time went along? Manual drill-down and browsing are laborious and overwhelming especially when there is a large number of attributes and each attribute is of high cardinality. Figure 1b illustrates a sample drill-down view along attribute `state`: first, looking at these sampled 10 states is already distracting, let alone full 58 states in the US; more importantly, it is still not clear what the answer to the above question is, though we do observe that each state contributes differently as time moves along. For instance, `state=NY` drives the initial outbreak in the US, while `state=CA` contributes most during the end of 2020.

● `S&P500`. `S&P500` [49] is a stock market index tracking the performance of around 500 companies in the US. In a nutshell, `S&P500` is calculated as the weighted average of these companies' stock prices. After seeing how `S&P500` evolves during 2020, users might be interested in explaining the movement of `S&P500` by stock category. Intuitively, different stock categories drive the drop and rebound of `S&P500`. For instance, `Category=Financial` plays a significant role in the drop during early covid-19 outbreak, but does not contribute to the rebound that much during the second quater of 2020.

● `Liquor-sales`. The `liquor-sales` time series corresponds to query `SELECT date, SUM(Bottles_Sold) FROM Liquor GROUP BY date`, where each row in relation `Liquor` represents a liquor purchase transaction with attributes including `date`, `Bottle_Volume(ml)`, `Pack`, `Category_Name`, and `Bottles_Sold`. Data analysts may wonder why the total bottles sold turns up since mid-January-2020 and how drinking behavior changes during pandemic. As we will show in our experiment, it turns out that people favors large pack liquor during pandemic, leading to sharp sales increase of `Pack=12` and `Pack=24`; and that people increase the purchase of large volume liquor such as `Bottle_Volume(ml)=750` and `Bottle_Volume(ml)=1750`.

**Problem and Challenges.** These motivating examples can all be abstracted as the same problem: given a relation $R$, a set of explain-by attributes from $R$, and a time series aggregated from $R$, identify the top `explanations`, i.e., conjunctions of predicates over explain-by attributes, that contribute to the changes in the given aggregated time series. Let us illustrate the motivating example of `covid-19` using this problem formulation. Given a relation `Covid-19`, where each row records the total number of confirmed cases in a state on a particular date, a set of explain-by attributes, e.g., `[state]`, and a time series aggregated from $R$, e.g., Figure 1a corresponding to query "`SELECT date, SUM(total_confirmed_cases) FROM Covid-19 GROUP BY date`", our goal is to identify `explanations`, e.g., `state=NY`, that explains the surge in Figure 1a. There are two main challenges in solving this problem: *(a)* `explanations` evolve over time; *(b)* interactivity is critical for data exploration and analytics.

Challenge (a): `explanations` tend to change over time. For instance, by looking at Figure 1b, we can observe that the increase in `New York(NY)` is the main reason of the total increase in Figure 1a during `2020-04` and `2020-05`, while `California(CA)` is the driving force during December 2020. Similarly, different stock categories are responsible for the surges and declines of `S&P500` during different periods. Specifically, `technology` and `financial` industries play leading role in the sink of `S&P500` during the initial covid-19 outbreak (around `2020-02-06` to `2020-03-24`); while `technology` industry is the top-contributor for `S&P500`'s bounce back since
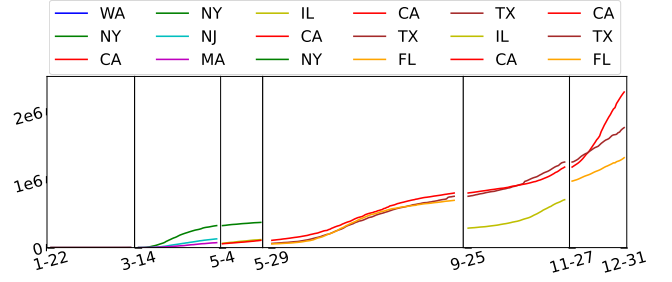


**Figure 2: Evolving Explanations of Figure 1a.**

`2020-03-24`, but `financial` industry is not. Having observed that `explanations` evolve along the time, our first technical challenge lies in how to identify time period with consistent `explanations` and how to derive `explanations` for each consistent time period.

Challenge (b): data analysts typically explore "what" and "why" questions iteratively to understand data and uncover insights. Studies [28] have shown that low latency is critical in fostering user interaction, exploration, and the extraction of insights. It is desired that each query, including both "what" and "why" queries, can get answered in a second to ensure interactivity. Thus, how to reduce the latency for deriving explanations poses another challenge.

**Prior Works.** "Why" questions are gradually gaining attraction both academia-wise and industrial-wise. However, instead of explaining the continuous changes in a time series, existing works focus on explaining either (1) one aggregated value [19], e.g., point $p_1$ in Figure 1a; or (2) the difference between two given relations [2, 27, 33, 39, 44, 51], e.g., a test relation and a control relation corresponding to point $p_1$ and $p_2$ respectively in Figure 1a. Reiterating the "Disclaimer" above, these tools can recommend and expedite answering "why" questions, but the human-in-the-loop interpretation is still required for true root cause analysis. To summarize, no prior works have studied the problem of explaining the evolving changes in time series as depicted in our motivating examples. Please refer to Section 2 for detailed comparison.

**Our Solution.** We propose TSExplain, a system to explain the continuous changes in aggregated time series. In TSExplain, given a relation, users can freely perform OLAP operations, including drill-down, roll-up, slicing, and dicing, and visualize what has happened to some KPI. To explain, users can then specify the time period they are interested as well as a set of explain-by attributes based on their domain knowledge. For the `COVID-19` example, TSExplain returns a trendline visualization in Figure 2, where the whole input time series get partitioned into a few non-overlapping time periods and each time period is associated with the KPI trendlines for top `explanations`. We can see that in early stage, NY and WA are the main contributors to the case increase; while CA, TX, IL,and FL become the main contributors later 2020.

Technically, to tackle challenge (a) of deriving *evolving explanations*, we formulate a *K-Segmentation* problem, aiming to partition the input time period into $K$ smaller time periods such that each time period shares consistent top `explanations`. Furthermore, since $K$ is hard to specify in practice, TSExplain employs "Elbow method" [40] for identifying the optimal $K$. We demonstrate

the effectiveness of our problem formulation with both synthetic and real-world datasets. As for challenge (b) of interactivity, we first analyze the complexity of each step in TSExplain, identifying the bottleneck in the whole pipeline. Next, we propose several optimizations for reducing the bottlenecks in TSExplain. In our experiments, TSExplain has successfully answered all queries within one second.

**Contributions.** The contributions of this paper are as follows:
- We formulate $K$-Segmentation problem for explaining the continuous changes in an aggregated time series. (Section 3)
- We propose a novel within-segment variance design based on top explanations to quantify consistency in each segment and experimentally prove its effectiveness. (Section 4)
- We develop a dynamic program algorithm, perform complexity analysis, propose optimal selection strategy of $K$, and develop several optimizations for improving efficiency. (Section 5)
- We conduct experiments on both synthetic and real-world datasets, demonstrating the effectiveness and efficiency of TSExplain. (Section 4.2 and 7)

## 2 RELATED WORK

**Data Explanation** Existing data explanation engines mainly focuses on explaining (1) one aggregated value [11, 19] or (2) the difference between two relations [1, 2, 11, 16, 27, 30, 33, 37–39, 41, 44, 47, 51]. In academia, SmartDrillDown [19] explains one aggregated value by identifying explanations (called rules in SmartDrillDown) that have high aggregate value. IDIFF [39] identifies the differences between two instances of an OLAP cube. Scorpion [51] and Macrobase [2] aim to find the difference between outlier and inlier data. RSExplain [37] proposes an intervention-based framework to explain why SQL expression's result is high (or low). X-Ray [47] tries to reveal the common properties among all incorrect triples versus correct triples. Abuzaid et al. [1] unified various explanation engines and abstracted out a logical operator called diff. The cascading analysts algorithm [38] provides top *non-overlapping* explanations accounting for the major difference of two specified sets. Li et al. [27] also compares two set differences but with augmented information from other related tables. In industry, `Google Trend` integrates an explanation component for single value and two relation; `Tableau` [43] provides `Explain data` [44] feature; `PowerBI` [31] supports functionality like `Key Influener` [33]; startups like `SisuData` [41] is built to support "why" questions natively and scalably. However, all prior works fall short of explaining aggregated time series because (1) only explaining one aggregated value overlooks the trend of time series -"why up/down"; (2) only focusing on two set differences (i.e., two endpoints in time series) dismisses the explanation in between. Our TSExplain aims to identify the evolving explanations for aggregated time series over time.

**Time Series Segmentation** Time series segmentation has been studied for decades. We note that the word "segmentation" is somewhat overloaded in the literature. One line of segmentation works focuses on visual-based piecewise linear approximation. Specifically, [24, 46] use the sliding windows algorithm, which anchors the left point of a potential segment, then attempts to approximate the data to the right with increasing longer segments. Douglas et al., [8] and Ramer et al., [35] designs top-down algorithms to

partition the visualization. [22] and [14] have used the bottom-up algorithm to merge from the finest segments. Keogh et al. [21] show that the bottom-up algorithm achieves the best results compared with sliding window and top-down, and further introduced a new online algorithm that combines the sliding window and bottom-up to avoid rescanning of the data when streaming.

Another line of work is semantic segmentation including FLUSS [9, 10], AutoPlait [29], NNSegment proposed in LimeSegment [42], which aims to divide a time series into internally consistent subsequences, e.g., segmenting the heartbeat cycles when a person switches from running to walking. Thus, these algorithms require an extra input called subsequence length, e.g., a heartbeat cycle. However, our task is to explain the trends in the aggregated time series (e.g., the trend of total covid cases) instead of explaining or identifying the periodic difference in one time-series instance. Hence, our segmentation does not rely on the period length.

To conclude, unlike all above, TSExplain is the first to segment time series based on segments' explanations.

**Time-series ML Model Explanation** The time-series ML model takes a univariate or multivariate time series as input and outputs a prediction label. Unlike text or image models, there is limited literature on black box Time-series ML Model explainability. FIT [45] is an explainability framework that defines the importance of each observation based on its contribution to the black box model's distributional shift. Similarly, Rooke et al. [36] extend FIT into WinIT, which measures the effect on the distribution shift of groups of observations. Labaien et al. [25] finds the minimum perturbation required to change a black box model output. Recently, LimeSegment [42] selects representative input time series segments as explanations for time-series classifier's output. In these works, the "explain target" is the prediction label and the time series serves as the "explain feature". Contrarily, the "explain target" in TSExplain is the up/down trend in an aggregated time series and the explain-by attributes are our "explain features". Unlike explaining instance-level prediction of the black box ML model, TSExplain is aggregation-level explanation for white-box aggregation.

## 3 PROBLEM OVERVIEW

In this section, we start with existing works on two-relations diff and how they fall short for explaining aggregated time series, followed by formal formulations of our problem: *K-Segmentation*.

### 3.1 Background

*3.1.1 Two-Relations Diff.* Diff operator [1] focuses on identifying the difference between two relations. Given a test relation $R_t$ and a control relation $R_c$, the diff operator returns *explanations* describing how these two relations differ.

DEFINITION 3.1 (EXPLANATION [1]). *Given a set of explain-by attributes $\mathcal{A}$, an explanation $E$ of order $\beta$ is defined as a conjunction of $\beta$ predicates, denoted as $E = (A_1{=}a_1\&...\&A_\beta{=}a_\beta)$ where $A_i \in \mathcal{A}$.*

Explain-by attributes $\mathcal{A}$ can be specified by users based on their domain knowledge; otherwise, dimension attributes from $R_t$ are used. Intuitively, an explanation $E$ corresponds to a data slice satisfying the predicate, and this data slice contributes to the overall difference between $R_t$ and $R_c$. To quantify how well an explanation $E$

| Symb. | Definition | Mathematical Expression |
|---|---|---|
| $\mathcal{A}$ | explain-by attributes | $\mathcal{A}=\{A_1, A_2, ...\}$ |
| $E$ | an explanation | $E=(A_1=a_1..\&A_\beta=a_\beta), A_i \in \mathcal{A}$ |
| $\gamma(E)$ | difference score of $E$ | Definition 3.2 |
| $\tau(E)$ | change effect of $E$ | Definition 3.3 |
| $\mathbb{E}_m$ | m non-overlap $E$ | $\mathbb{E}_m=\{E_1, ...E_m\}$ |
| $\mathbb{E}_m^*$ | top-m non-overlap $E$ | $\mathbb{E}_m^*=\arg\max_{\mathbb{E}_m}[\sum_{E \in \mathbb{E}_m} \gamma(E)]$ |
| $ts$ | time series | $ts=\{p_1, ..p_i, ..p_n\}$ over $[t_1, t_n]$ |
| $c_i$ | $i^{th}$ cutting position | $c_i \in [1, n]$ at time $t_{c_i}$ |
| $P_i$ | segment $i$ | $P_i=[p_{c_i}, p_{c_{i+1}}]$ from $t_{c_i}$ to $t_{c_{i+1}}$ |
| $\mathcal{P}_K$ | K-segment scheme | $\mathcal{P}_K=\{P_1, P_2, .., P_K\}$ |
| $\mathcal{E}$ | evolving explanations | $\mathcal{E}=[\mathbb{E}_m^*(t_{c_1}, t_{c_2}), .., \mathbb{E}_m^*(t_{c_k}, t_{c_{k+1}})]$ |

**Table 1: Notations**

explains the difference between $R_t$ and $R_c$, diff operator [1] provides a *difference metric* abstraction, denoted as $\gamma(E)$. Such abstraction is capable of encapsulating the semantics of many prior explanation engines [2, 47, 51]. Commonly used $\gamma(E)$ include absolute-change, relative-change, risk-ratio. Throughout this paper, we focus on absolute-change. Other metrics can be applied similarly.

DEFINITION 3.2 (Absolute-change). Given a test relation $R_t$, a control relation $R_c$, an aggregate function $f(M, R)$ on some measure attribute $M$ in relation $R$, and an explanation $E$, the absolute change refers to the absolute difference between $[f(M, R_t) - f(M, R_c)]$ before and after removing records that explanation $E$ corresponds to, i.e., $\gamma(E) = |[f(M, R_t) - f(M, R_c)] - [f(M, R_t - \sigma_E R_t)) - f(M, R_c - \sigma_E R_c)]|$, where $\sigma_E R_t$ and $\sigma_E R_c$ denote records satisfying the predicate $E$ in relation $R_t$ and $R_c$ respectively.

As the name indicates, absolute-change only cares about the absolute change of $[f(M, R_t) - f(M, R_c)]$ no matter the change is an increase or decrease. To distinguish an increase from a decrease, we use $\tau(E)$ to denote the change effect caused by including data that $E$ corresponds to: intuitively, if including $E$ leads to an increase in $[f(M, R_t) - f(M, R_c)]$, $\tau(E) = +$; otherwise, $\tau(E) = -$.

DEFINITION 3.3 (CHANGE EFFECT). Following the setting in Definition 3.2, the change effect of an explanation $E$ is defined as $\tau(E) = sign([f(M, R_t) - f(M, R_c)] - [f(M, R_t - \sigma_E R_t)) - f(M, R_c - \sigma_E R_c)])$.

Now we have described Absolute-change as an example of $\gamma(E)$. With a difference metric $\gamma(E)$, we can then rank each candidate explanation and return top-m explanations with the highest $\gamma(E)$. However, these top-m explanations may contain overlapping records. Consequently, the effects of these records get duplicated, introducing bias in the top-m explanations. Alternatively, we can constrain these $m$ explanations to be non-overlapping and define *top-m non-overlapping explanations* as in Definition 3.5. Two explanations $E_1$ and $E_2$ are said to be *non-overlapping* if their correspondent records are non-overlapping in any relation $R$, i.e., $\sigma_{E_1} R \cap \sigma_{E_2} R = \emptyset, \forall R$.

DEFINITION 3.4 (m Non-Overlapping Explanations). Given a difference metric $\gamma(E)$ and an explanation order threshold $\bar{\beta}$, let $\mathbb{E}_m$ denote a set of $m$ non-overlapping explanations, i.e., $\mathbb{E}_m =$

$\{E_1, ...E_m\}$, where each $E_i$ has its order $\leq \bar{\beta}$ and is non-overlapping with $E_j$, $\forall E_i, E_j \in \mathbb{E}_m$.

DEFINITION 3.5 (Top-m Non-Overlapping Explanations [2]). Top-m non-overlapping explanations are defined as $\mathbb{E}_m$ with the highest accumulative diff score, i.e., $\mathbb{E}_m^*=\arg\max_{\mathbb{E}_m}[\sum_{E \in \mathbb{E}_m} \gamma(E)]$.

Cascading analyst algorithm [38] is designed for returning top-m non-overlapping explanations. We will use term *top-explanation* for simplicity whenever there is no ambiguity.

*Example 3.1 (Two-Relations Diff).* Consider the two points $p_1$ and $p_2$ in Figure 1a — the underlying data corresponds to $p_1$ and $p_2$ constitute a control relation $R_c$ and a test relation $R_t$ respectively. Two-relations diff aims to explain the difference between $R_c$ and $R_t$. First, users can specify a set of explain-by attributes, e.g., [state, County]. Take explanation $E=(\text{state=CA})$ as an example. It is of order one, i.e., $\beta = 1$ and its difference score $\gamma(E)$ can be calculated as $|(p_2.v - p_1.v) - (p_2'.v - p_1'.v)|$ using absolute-change as shown in Figure 1a. We can then obtain Top-3 non-overlapping explanations for differing $R_t$ and $R_c$ using cascading analyst algorithm [38] — $\mathbb{E}_3^*=\{E_1=(\text{state=CA}), E_2=(\text{state=TX}), E_3=(\text{state=FL})\}$.

*3.1.2 Aggregated Time Series.* Time series is a series of data points indexed in time order. An *aggregated time series* refers to a special type of time series, where each point $p$ is associated with a timestamp $p.t$ and an aggregated value $p.v$, derived by aggregating all records at timestamp $p.t$. Essentially, an aggregated time series corresponds to the result of some group-by query. Consider a relation $R$ with $\{D_i\}$ dimension attributes and $\{M_j\}$ measure attributes, and a group-by query in the form of {SELECT $T$, f(M) FROM R GROUP BY $T$}, where $T$ denotes some time-related ordinal dimension ($T \in \{D_i\}$), and $f(M)$ is some aggregate function on measure $M$ ($M \in \{M_j\}$). The query result can be denoted by an aggregated time series $ts$ with value $f(M)$ over time dimension $T$.

DEFINITION 3.6 (AGGREGATED TIME SERIES). An aggregated time series $ts$ over time $[t_1, t_n]$ is a series of points $\{p_1, ..., p_i, ..., p_n\}$ ordered by time dimension $T$ and each point's value $p_i.v$ is an aggregated number from a list of records with the same time $p_i.t$.

Data analysts often visualize aggregated time series to help understand data's overall trend as time goes along as shown in Figure 1a. A natural follow-up question is "what makes ups and downs". Different from two-relations diff described in Section 3.1.1, this "explain" question focuses on the whole time horizon and the underlying top-explanation tends to evolve dynamically along the time even when the overall trend looks the same visually.

DEFINITION 3.7 (EVOLVING EXPLANATIONS). Given $m$ and an aggregated time series $ts$ over time $[t_1, t_n]$, *evolving explanations* is a sequence of top-explanation at different periods, denoted as $\mathcal{E} = [\mathbb{E}_m^*(t_{c_1}, t_{c_2}), \mathbb{E}_m^*(t_{c_2}, t_{c_3}), ..., \mathbb{E}_m^*(t_{c_k}, t_{c_{k+1}})]$, where $c_1 = 1, c_{k+1} = n$, $\{c_2, c_3..c_k\}$ denote the $(k-1)$ cutting positions in between, and each $\mathbb{E}_m^*(t_{c_i}, t_{c_{i+1}})$ denotes top-explanation from $t_{c_i}$ to $t_{c_{i+1}}$.

*Example 3.2 (Evolving Explanations).* Figure 1a depicts an aggregated time series that corresponds to a groupby-aggregate query

---

[2]Definition 3.5 is defined over $\mathbb{E}_m$. Alternatively, we can define top-m as *at most* $m$ explanations, i.e., $\mathbb{E}_m^*=\arg\max_{\{\mathbb{E}_x | x \leq m\}}[\sum_{E \in \mathbb{E}_x} \gamma(E)]$. Our proposed solution in Section 4 and 5 can work with it in a similar way.

with $f(M)$=SUM(total_confirmed_cases) on table Covid-19. Figure 2 illustrates the underlying evolving explanations for the increase in Figure 1a. We have six different time periods, where each period shares the same intrinsic explanations while neighboring periods have different ones. For instance, the top-3 contributors are $\mathbb{E}_3^*(t_{c_2}, t_{c_3})$={$E_1$=(state=NY), $E_2$=(state=NJ), $E_3$=(state=MA)} during $t_{c_2}$=2020-3-14 and $t_{c_3}$=2020-5-4; $\mathbb{E}_3^*(t_{c_6}, t_{c_7})$={$E_1$=(state=CA), $E_2$=(state=TX), $E_3$=(state=FL)} from 2020-11-27 to 2020-12-31.

## 3.2 Problem Definition

Motivated by the observation that top contributor (i.e., explanation) evolves over time, we study the problem of identifying *evolving explanations* for the continuous changes happened in an aggregated time series. The overall problem of identifying evolving explanations can be decomposed into two sub-problems: *(a).* segmentation; *(b).* find the explanations that contributes most in each segment. These two sub-problems are intertwined with each other: the goodness of a segmentation scheme depends on how cohesive top-explanations are within each segment; meanwhile, top-explanation are derived for each segment after obtaining the optimal segmentation scheme. We remark that performing segmentation without considering explanation information is insufficient, as we will demonstrate experimentally in Section 7.

**Segmentation.** To explain the continuous change in an aggregated time series $ts$ over time $[t_1, t_n]$, we need to partition the whole time domain $[t_1, t_n]$ into non-overlapping segments, such that each segment $P_i$=$[p_{c_i}, p_{c_{i+1}}]$ during time $t_{c_i}$ and $t_{c_{i+1}}$ shares the same intrinsic explanations while neighboring segments have different ones. This resembles the well-studied clustering problem, whose goal is to minimize within-cluster variance and maximize inter-cluster variance. In particular, given a segment number $K$, we abstract our problem as a *K-Segmentation* problem, adapting the optimization formula from K-Means [13]. Let $\mathcal{P}_K$ denote a $K$-segmentation scheme: $\mathcal{P}_K = \{P_1=[p_{c_1}, p_{c_2}], P_2=[p_{c_2}, p_{c_3}]...P_K=[p_{c_K}, p_{c_{K+1}}]\}$, where $c_1$=1, $c_{K+1}$=$n$, and $\{c_2, c_3..., c_K\}$ are $(K-1)$ cutting positions. In Example 3.2 (Figure 2), $K$=6 and the cutting positions $\{c_2, ..., c_5\}$ correspond to time $\{3\text{-}14, 5\text{-}4, 5\text{-}29, 9\text{-}25, 11\text{-}27\}$. Next, we formulate *K-Segmentation* problem as below:

PROBLEM 1 (*K-SEGMENTATION*). *Given an aggregated time series $ts$ and a segmentation number $K$, identify the optimal segmentation scheme $\mathcal{P}_K^* = \arg\min_{\mathcal{P}_K=\{P_1,...,P_K\}} \sum_{i=1}^K |P_i|\text{var}(P_i)$ where $\text{var}(P_i)$ denotes the variance in segment $P_i$.*

We remark that the design of within-segment variance $\text{var}(P_i)$ is critical to the effectiveness of $K$-Segmentation. No prior works have studied $\text{var}(P_i)$ with the goal of quantifying explanation consistency. This is challenging as we will dive into in Section 4.

**Explain trend in each segment.** Given a fixed segment $P_i = [p_{c_i}, p_{c_{i+1}}]$ from time $t_{c_i}$ to $t_{c_{i+1}}$, we will now describe how to explain the trend in $P_i$ with only one static top-explanation $\mathbb{E}_m^*(t_{c_i}, t_{c_{i+1}})$ — static top-explanation is a special case of evolving explanations $\mathcal{E}$ with segment number $K = 1$. If $P_i$ is cohesive, meaning that $P_i$ has consistent top-explanation during $t_{c_i}$ and $t_{c_{i+1}}$, we can simply focus on its two endpoints and then employ prior works on two-relations diff (Section 3.1). The derived top-explanation $\mathbb{E}_m^*$ explains

the changes from time $t_{c_i}$ to $t_{c_{i+1}}$. However, when $P_i$ is not cohesive, there exists no single static top-explanation $\mathbb{E}_m^*$ that is capable of explaining the whole trend evolvement in $P_i$. Nevertheless, we can still derive some static top-explanation by looking at its two endpoints and using two-relations diff [38], though the explanation quality might be poor. As we will elaborate in Section 4.1, each segment in the optimal K-segmentation $\mathcal{P}_K^*$ is deemed to be cohesive, and the case of incohesive segment would only occur during the exploration phase over candidate segmentation schemes. In all, when given a segment $P_i = [p_{c_i}, p_{c_{i+1}}]$, we first obtain a control relation $R_c$={SELECT * FROM R WHERE $T = t_{c_i}$} and a test relation $R_t$={SELECT * FROM R WHERE $T = t_{c_{i+1}}$} at two endpoints, and then derive top-explanation with cascading analyst algorithm [38].

Now that we can exploit existing works for deriving static top-explanation within each segment, our problem of identifying evolving explanations boils down to a K-segmentation problem. As $K$ is not easy to specify in practice, TSEXPLAIN identifies the optimal $K$ for users by default, as we will discuss in Section 6.

## 4 K-SEGMENTATION

Designing a good within-segment variance to quantify explanation consistency is the key to our work's success and it is definitely non-trivial. In this section, we will start with our design of $\text{var}(P_i)$ in Problem 1, followed by an experiment illustrating its effectiveness.

## 4.1 Design of Within-Segment Variance

Per our definition in Problem 1, K-Segmentation is similar to K-Means clustering. We carefully design the within-segment variance $\text{var}(P_i)$ through making analogy to K-Means algorithm. However, different from K-Means or any existing variance design, $\text{var}(P_i)$ in K-Segmentation should be regarding the variance of top-explanations within each segment $P_i$.

First, let us review the problem formulation of K-Means. Given a set of *objects* $(o_1, o_2, ..., o_n)$ as inputs, where each object is a $d$-dimensional vector, K-Means aims to partition these $n$ objects into $K$ partitions $\mathcal{P}_K$={$P_1, P_2..P_K$} minimizing the within-cluster variance:

$$\arg\min_{\mathcal{P}_K} \sum_{i=1}^K |P_i|\text{var}(P_i) \tag{1}$$

$$\text{var}(P_i) = \frac{1}{|P_i|} \sum_{o \in P_i} dist(o, \mu_i) \tag{2}$$

where $\mu_i$ is the centroid of partition $P_i$ and $dist(o, \mu_i)$ denotes the distance between an object $o$ and the centroid $\mu_i$ in $P_i$, e.g., $L2$ distance. Comparing Problem 1 with Eq. 1, we can see K-Segmentation employs the same optimization objective as K-Means but with a customized $\text{var}(P_i)$. To develop a good $\text{var}(P_i)$ in K-Segmentation, careful thoughts are required around: *(1)* what is an "object"; *(2)* what is the centroid of a segment; and *(3)* how to measure the distance between object and centroid based on explanations.

*4.1.1 Object in K-Segmentation.* Given an aggregated time series $ts = \{p_1, p_2, ..., p_n\}$, K-Segmentation aims to segment $ts$ into $K$ partitions such that explanations are shared within each partition. Since each single point itself cannot reveal any time series trend, the atomic unit for partitioning is a segment of size two, i.e., $[p_i, p_{i+1}]$. That is, an object in K-Segmentation refers to a segment of size
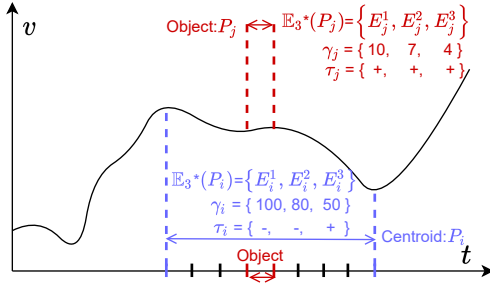
**Figure 3: Object, Centroid, and Top-Explanations**

| Rank | Expl | How well $\mathbb{E}_m^*(P_j)$ explains $P_i$ | | | |
|---|---|---|---|---|---|
| | | Effect (+/-) | | Relevance | Rectified |
| | | on $P_j$ | on $P_i$ | on $P_i$ | Relevance $\bar{\gamma}$ |
| r | $E_j^r$ | $\tau(E_j^r, P_j)$ | $\tau(E_j^r, P_i)$ | $\gamma(E_j^r, P_i)$ | $\gamma(E_j^r, P_i) \times$ $\mathbb{1}_{\tau(E_j^r, P_j)=\tau(E_j^r, P_i)}$ |
| 1 | $E_j^1$ | + | + | $\gamma(E_j^1, P_i)$ | $\gamma(E_j^1, P_i)$ |
| 2 | $E_j^2$ | + | + | $\gamma(E_j^2, P_i)$ | $\gamma(E_j^2, P_i)$ |
| 3 | $E_j^3$ | + | - | $\gamma(E_j^3, P_i)$ | 0 |
| Discounted cumulative gain (DCG): | | $\frac{\gamma(E_j^1, P_i)}{\log_2(1+1)} + \frac{\gamma(E_j^2, P_i)}{\log_2(1+2)} + \frac{0}{\log_2(1+3)}$ | | | |

**Table 2: Example of DCG Between $\mathbb{E}_m^*(P_j)$ and $P_i$**

two as shown in Figure 3 and there is in total $(n-1)$ objects, i.e., $\{o_1 = [p_1, p_2], o_2 = [p_2, p_3], ..., o_{n-1} = [p_{n-1}, p_n]\}$.

*4.1.2 Centroid of a partition.* Different from the general K-Means, objects in K-Segmentation follows a time ordering where $o_1 < o_2 < ... < o_{n-1}$, and a partitioning scheme in K-Segmentation is only valid when objects in each partition form a segment. Hence, a partition in K-Segmentation can be denoted as $P_i = [p_{c_i}, p_{c_{i+1}}]$ with objects $\{o_{c_i}=[p_{c_i}, p_{c_i+1}], o_{c_i+1} .., o_{c_{i+1}-1}=[p_{c_{i+1}-1}, p_{c_{i+1}}]\}$. Naturally, we can use segment $[p_{c_i}, p_{c_{i+1}}]$ as the centroid of partition $P_i$.

*4.1.3 Distance between object and centroid.* In essence, both object and centroid are segments of the input time series *ts*. Next, we focus on the design of distance between segments. Distance between two objects or between an object and a centroid follow naturally.

**High-Level Idea.** As our goal is to group objects with the same top-explanations into one partition, the distance between two segments shall be based on their top-explanations. Given two segments $P_i$ and $P_j$, their top-explanations $\mathbb{E}_m^*(P_i)$ and $\mathbb{E}_m^*(P_j)$ can be derived based on Section 3.1. Each is a ranked list of explanations, i.e., $\mathbb{E}_m^*(P_i) = [E_i^1, E_i^2, ..., E_i^m]$ and $\mathbb{E}_m^*(P_j) = [E_j^1, E_j^2, ..., E_j^m]$ as shown in Figure 3. Strawman approaches like measuring the Jaccard distance between $\mathbb{E}_m^*(P_i)$ and $\mathbb{E}_m^*(P_j)$ fall short when there are multiple explain-by attributes. For instance, it is unclear how to quantify the partial overlap between $E_i^1$={state=WA} and $E_j^2$={state=WA and age>50}. Alternatively, we can measure the distance between $P_i$ and $P_j$ by how well $\mathbb{E}_m^*(P_i)$ explains $P_j$ and how well $\mathbb{E}_m^*(P_j)$ explains $P_i$.

**How well $\mathbb{E}_m^*(P_j)$ explains $P_i$.** We draw inspirations from information retrieval community. Normalized discounted cumulative gain (NDCG) is a commonly used measure for ranking quality in information retrieval and we adapt NDCG to quantify how well $\mathbb{E}_m^*(P_j)$ explains $P_i$. To model our scenario after the web search setting, we can treat each segment $P_i$ as a query and each explanation $E$ as a document. $\mathbb{E}_m^*(P_j)$ corresponds a ranked list of retrieved documents returned by the search engine, while $\mathbb{E}_m^*(P_i)$ is the ideal retrieved document list. The relevance between a segment $P_i$ and an explanation $E$ is quantified by the difference metric $\gamma(E, P_i)$. As illustrated in Table 2 (row in blue), given an explanation $E_j^r \in \mathbb{E}_m^*(P_j)$ with rank $r$, the relevance of $E_j^r$ towards $P_i$ can be calculated as $\gamma(E_j^r, P_i)$. However, explanation $E_j^r$ might make KPI increase in segment $P_i$, but lead to a decrease in segment $P_j$ (see Definition 3.3). Thus, when $E_j^r$ has opposite effects on $P_i$ and $P_j$, we shall rectify

the relevance to zero as our ultimate goal is to measure the distance between $P_i$ with $P_j$. Formally, we denote the *rectified relevance* as $\bar{\gamma}(E_j^r, P_i)$ and we have $\bar{\gamma}(E_j^r, P_i) = \gamma(E_j^r, P_i) \times \mathbb{1}_{\tau(E_j^r, P_j)=\tau(E_j^r, P_i)}$. Take $E_j^3$ in Table 2 as an example, $E_j^3$ contributes the increase in segment $P_j$ but the decrease in $P_i$, thus the relevance is rectified to 0.

Now we have mapped our scenario to query-document retrieval setting, i.e., query $P_i$, retrieved document list $\mathbb{E}_m^*(P_j)$, and the rectified relevance formula $\bar{\gamma}$, $NDCG(P_i, \mathbb{E}_m^*(P_j))$ can be calculated using Eq. 3, 4, and 5. It quantifies how well $\mathbb{E}_m^*(P_j)$ explains $P_i$, with ranges from 0 to 1. At an extreme, when $\mathbb{E}_m^*(P_j)$ is exactly the same as $\mathbb{E}_m^*(P_i)$ and with the same effect on $P_i$ and $P_j$, $NDCG(P_i, \mathbb{E}_m^*(P_j))=1$, meaning $\mathbb{E}_m^*(P_j)$ explains $P_i$ perfectly.

$$DCG(P_i, \mathbb{E}_m^*(P_j)) = \sum_{r=1}^{m} \frac{\bar{\gamma}(E_j^r, P_i)}{\log_2(r+1)} \tag{3}$$

$$DCG(P_i, \mathbb{E}_m^*(P_i)) = \sum_{r=1}^{m} \frac{\bar{\gamma}(E_i^r, P_i)}{\log_2(r+1)} = \sum_{r=1}^{m} \frac{\gamma(E_i^r, P_i)}{\log_2(r+1)} \tag{4}$$

$$NDCG(P_i, \mathbb{E}_m^*(P_j)) = \frac{DCG(P_i, \mathbb{E}_m^*(P_j))}{DCG(P_i, \mathbb{E}_m^*(P_i))} \tag{5}$$

**Calculating Distance.** We can then define the distance between $P_i$ and $P_j$ as in Eq. 6, where $NDCG(P_i, \mathbb{E}_m^*(P_j))$ quantifies how well $\mathbb{E}_m^*(P_j)$ explains $P_i$ and $NDCG(P_j, \mathbb{E}_m^*(P_i))$ quantifies how well $\mathbb{E}_m^*(P_i)$ explains $P_j$.

$$dist(P_i, P_j) = 1 - \frac{NDCG(P_i, \mathbb{E}_m^*(P_j)) + NDCG(P_j, \mathbb{E}_m^*(P_i))}{2} \tag{6}$$

Eq. 6 averages $NDCG(P_i, \mathbb{E}_m^*(P_j))$ and $NDCG(P_j, \mathbb{E}_m^*(P_i))$ to obtain the similarity between $P_i$ and $P_j$, followed by a complement to get the distance. $dist(P_i, P_j)$ is symmetric with ranges $[0, 1]$.

*4.1.4 Putting all together.* Given a partition $P_i = [p_{c_i}, p_{c_{i+1}}]$, it contains a continuous list of objects $\{P_x = [p_x, p_{x+1}]\}$ where $c_i \le x < c_{i+1}$ and its centroid is $P_i = [p_{c_i}, p_{c_{i+1}}]$. Using Eq. 6 and 2, we can derive our variance of $P_i$ as in Eq. 7.

$$var(P_i) = \frac{1}{c_{i+1} - c_i} \sum_{x=c_i}^{c_{i+1}-1} dist(P_x, P_i), \text{ where } P_x = [p_x, p_{x+1}] \tag{7}$$

Our problem formulation of K-Segmentation is now complete, by instantiating $var(P_i)$ in Problem 1 with Eq. 7.

## 4.2 Effectiveness of Variance

In this subsection, we evaluate the effectiveness of our variance design. We term our variance metric in Equation (7) `tse`. Since real-world datasets lack its ground truth of evolving explanations, we synthesize datasets with ground truths and evaluate how `tse` performs compared with other alternatives. We will evaluate the end-to-end effectiveness in Section 7.

*4.2.1 Synthetic datasets.* We synthesize datasets and generate their ground truth K-Segmentation $\mathcal{P}_K^*$. Each dataset is one relation $R$ with schema: $T$, `sales`, `category`. The aggregated time series represents how the total sales changes along the time $T$ — SELECT $T$, count(sales) FROM R GROUP BY $T$. We set the explain-by attributes $\mathcal{A}$={category} and there are three categories: $a_1$, $a_2$, $a_3$. Each predicate e.g., {category=$a_1$} denotes an explanation $E$.

**Synthesize Procedure** The aggregated `sales` time series can be viewed as a summation of each `category`'s time series. We start by synthesizing each category's time series. In detail, we first randomly pick cutting points $\{c_1^i, ..., c_{j_i}^i\}$ for each category $E_i$'s time series — {SELECT $T$, count(sales) FROM R WHERE category=$a_i$ GROUP BY $T$}. For each segment $[p_{c_k^i}, p_{c_{k+1}^i}], 1 \le k < j_i$ defined by these cutting points, we synthesize either an upward or downward trend in linear shape. We restrict the adjacent segments to have different up or down trends. We can then derive the ground truth segmentation's cutting points of the aggregated time series as the union of each category's cutting points, i.e., $\bigcup_{i=1}^{3}\{c_1^i, ..., c_{j_i}^i\}$. We treat $\bigcup_{i=1}^{3}\{c_1^i, ..., c_{j_i}^i\}$ as our ground truth because (1) each predicate has a consistent up or down trend in each segment; (2) our restriction that adjacent segments have different trend direction guarantees that every cutting point is necessary and $\bigcup_{i=1}^{3}\{c_1^i, ..., c_{j_i}^i\}$ is the minimal coherent segmentation method. We set the time series' length at 100 and synthesize 20 datasets with seven different levels of $SNR_{dB}$. The lower the $SNR_{dB}$ is, the noisier the time series is. We remark that the number K and the length of each segment are diverse in our synthetic datasets, with segment number K varying from 2 to 10 and segment length varying from 6 to 84 as shown in Figure 4.
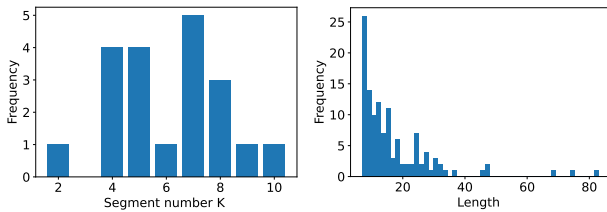


**Figure 4: Distribution of segment number K and length of each segment.**

**Signal-to-Noise Ratio** Real-world data is quite noisy. For our synthetic dataset, we add Gaussian Noise $N(0, \sigma^2)$ to each predicate's time series to simulate noisy time series. We quantify the noise

level using signal-to-noise level, namely SNR [48]. We add different noise $SNR_{dB} = 20, 25.., 50$ to each dataset. The lower the $SNR$ is, the noisier the time series is.

*Example 4.1 (Synthetic Dataset and Ground Truth Segmentation).* In Figure 5, the dash lines illustrate each predicate's time series, and the noise level is SNR = 35. The predicate category = a1 has its cutting points at 52, 76, the predicate category = a2 has its cutting point at 70, 90, and the predicatecategory = a3 has its cutting point at 31. Thus, in this example, we can derive the aggregated time series cutting point as the union of three predicates' – { 31, 52, 70, 76, 90 }.
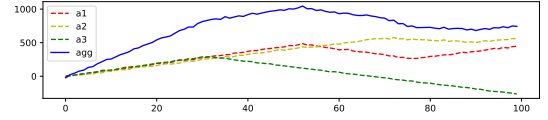


**Figure 5: A synthetic time series with SNR = 35.**

*4.2.2 Effectiveness.* We compare `tse` with other alternatives in terms of effectiveness.

**Alternatives** We design our alternatives by enumerating other reasonable distance metrics and variance structures.

*Change $dist(P_i, P_j)$ definition and keep the variance structure.* Different from our metric in eq. (6), if we only consider how well each object's explanation $\mathbb{E}_m^*(P_j)$ explains the centroid $P_i$, we can derive `dist1` in eq. (8); or if we only consider how well the centroid's $\mathbb{E}_m^*(P_i)$ explains each object $P_j$, we can derive `dist2` in eq. (9).

$$\text{dist1}: dist(P_i, P_j) = 1 - NDCG(P_i, \mathbb{E}_m^*(P_j)) \qquad (8)$$

$$\text{dist2}: dist(P_i, P_j) = 1 - NDCG(P_j, \mathbb{E}_m^*(P_i)) \qquad (9)$$

*Change the variance structure and keep the $dist(P_i, P_j)$ definition.* Instead of comparing the centroid with each object, we compare all possible object pairs. We name this metric `allpair`.

$$\text{allpair}: var(P_j) = AVG\{\sum_{P_x}\sum_{P_y} dist(P_x, P_y)\}, \ P_x, P_y \subseteq P_j$$

$$(10)$$

Based on `tse`, `dist1`, `dist2`, `allpair`, if we further change the second term in the distance metric to its l2 norm, we can derive `Stse/Sallpair`, `Sdist1` and `Sdist2` correspondingly. Till now, we have eight different forms of metrics.

**Evaluation** As shown in eq. (1), our problem is to find the segmentation that minimizes the objective score. Thus, the effectiveness criteria of variance design is whether the ground truth can score the lowest or close to the lowest in noisy datasets.

We compare `tse` with all other alternatives. Given one metric and one dataset $D$ with its ground truth, the K-segmentation search space $\mathcal{P}_K$ is huge and each segmentation scheme has its metric score. We study how the ground truth segmentation's score ranks among all possible segmentation schemes: the higher the rank is, the better the metric is. Because $\mathcal{P}_K$ space is huge, we sample 10000 random segmentation schemes, and rank the ground truth among

all the samples, denoted as *ground truth rank*. On this dataset $D$, we calculate all eight metric's *ground truth rank* following the same method. To compare all the metrics on $D$, we rank across all the eight metrics from *rank* 1 (highest rank) to *rank* 8 (lowest rank) ascendingly based on their own *ground truth rank*. The higher a metric's *rank* is, the fewer segmentation schemes have lower variance than the ground truth, indicating that this metric is more effective on $D$. Figure 6 averages each metric's *rank* over all datasets at the same SNR level and shows how different metrics' *ranks* change along with the noisy level. When SNR = 50dB, we can see that all metrics rank 1st because the ground truth all achieves the lowest score. What's more, tse metric always has the highest *rank* compared with other metrics no matter what level the SNR is.
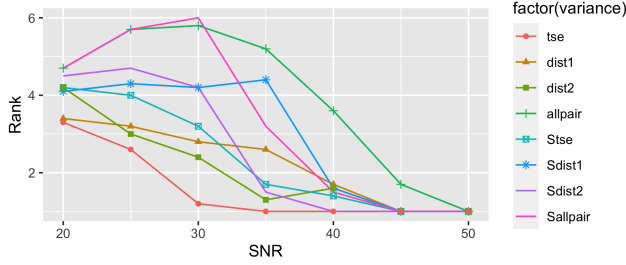


**Figure 6: The** *average rank* **of all metrics at different SNRs.**

**Takeaway.** Compared with other alternative metrics, tse is the most effective one.

## 5 OUR SOLUTION: TSEXPLAIN

Now we have formulated our K-Segmentation problem (Problem 1), together with the designed within-segment variance in Eq. 7. In this section, we will present our solution TSExplain: we will start with a dynamic programming (DP) algorithm for solving K-Segmentation, assuming $\text{var}(P_k)$ is available for each partition $P_k$; next, we will describe our solution pipeline (Figure 7) – steps for preparing $\text{var}(P_k)$ and DP – together with complexity analysis; last, we propose several optimizations for speeding up TSExplain.

### 5.1 DP for K-Segmentation

Different from K-Means, which is computationally intractable (NP-Hard), K-Segmentation is polynomial solvable. At a high-level, the search space in K-Means is $K^n$, while it is $\binom{n-2}{K-1}$ in K-Segmentation as the task is essentially to identify $(K-1)$ cutting positions among the $(n-2)$ non-endpoints of a given time series $ts$.

Intuitively, K-segmentation exhibits optimal substructure — the optimal solution of K-segmentation can be constructed from the optimal solutions of its subproblems. Thus, we can use dynamic programming for solving K-Segmentation. Let $D(j, k)$ denote the minimal total within-segment variance of $k$-segments over time range $[t_1, t_j]$, i.e., $D(j, k) = \min_{\mathcal{P}_k=[P_1,...,P_k]} \sum_{i=1}^{k} |P_i| \text{var}(P_i)$. To derive $D(j, k)$, we can enumerate different positions of the last cut $j'$ and calculate the smallest total variance among all possible $j'$. When the last cut is fixed at position $j'$, the minimal total variance of $k$-segments can be decomposed into two parts: the minimal total variance of $(k-1)$-segments over time range $[t_1, t_{j'}]$ and
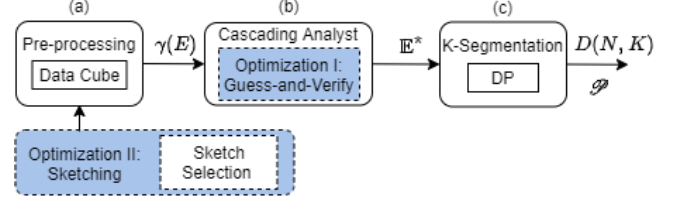


**Figure 7: Solution Pipeline in** TSExplain

the variance of the $k^{th}$ segment during $[t_{j'}, t_j]$, i.e., $D(j', k-1) + |P_k| \text{var}(P_k)$. The DP recursion function is expressed in Eq. 11.

$$D(j, k) = \underset{1 \leq j' \leq j}{\arg\min} [D(j', k-1) + |P_k| \text{var}(P_k)], \quad P_k = [p_{j'}, p_j] \quad (11)$$

$D(j, k)$ is recursively computed with Eq. 11, which involves calculating variance $\text{var}(P_k)$ for all segment $P_k = [p_{j'}, p_j]$ where $1 \leq j' < j \leq n$. Given a segment $P_k = [p_{j'}, p_j]$, its $\text{var}(P_k)$ is computed via the distance $dist(P_k, P_x)$ between centroid segment $P_k$ and each object segment $P_x = [p_x, p_{x+1}]$ where $j' \leq x < j$ as shown in Eq. 7. Now to calculate $dist(P_k, P_x)$, we need to identify the top-explanations $\mathbb{E}_m^*$ in segment $P_j$ and $P_x$. Top-explanations $\mathbb{E}_m^*$ is derived with Cascading Analyst algorithm [38], which requires computing the diff score $\gamma(E)$ for each explanation $E$. In Section 5.2, we will present our solution pipeline consisting of three modules: *(a.) Preprocessing module* for calculating $\gamma(E)$; *(b.) Cascading Analyst module* for deriving $\mathbb{E}_m^*$; and *(c.) K-Segmentation* module for computing $dist(P_k, P_x)$, $\text{var}(P_k)$, and finally $D(j, k)$. We will also analyze the time complexity for each module, pinpointing the bottleneck in the whole pipeline for optimizations in Section 5.3.

### 5.2 Solution Pipeline and Complexity Analysis

The overall solution pipeline is depicted in Figure 7 with module (a) Preprocessing, (b) Cascading Analyst, and (c) K-Segmentation. In the following, we will dive into each module and analyze its computational complexity. The optimization modules with blue background in Figure 7 will be introduced in Section 5.3.

**Precomputation.** Module (a) is responsible for computing the diff score $\gamma(E)$ for each candidate explanation $E$ over every segment $[p_{j'}, p_j]$ where $1 \leq j' < j \leq n$. Given an explanation order threshold $\bar{\beta}$, we can enumerate all candidate explanations — each is in the form of $E = [A_1=a_1, ..., A_\beta=a_\beta]$ where $1 \leq \beta \leq \bar{\beta}$. Let $\epsilon$ be the total number of candidate explanations. By default, $\bar{\beta}$ is set as 3.

As illustrated in Figure 7, given an explanation $E$ and a segment $[p_{j'}, p_j]$, the difference score $\gamma(E)$, i.e., absolute-change in Definition 3.2, can be calculated by looking at the points at time $t_{j'}$ and $t_j$ of the aggregated time-series $ts(R)$ and $ts(R - \sigma_E R)$, i.e., the aggregated time-series from relation $R$ when excluding data with predicate $E$. As most aggregate function $f(M)$ is decomposable, e.g., SUM, AVG, Variance, we can derive $ts(R - \sigma_E R)$ by using $ts(R)$ and $ts(\sigma_E R)$. Since TSExplain is designed to integrate with existing interactive data analysis tools like PowerBI, data cube is typically maintained in memory and thus we can easily access $ts(R)$ and $ts(\sigma_E R)$ from the data cube.

*Time complexity:* for each segment $[p_{j'}, p_j]$ and explanation $E$, it takes $O(1)$ for computing the difference score $\gamma(E)$. There is in total $n^2$ segments and $\epsilon$ explanations, thus the total time complexity for module (a) is $O(\epsilon \cdot n^2)$.

**The Cascading Analysts (CA) Algorithm.** Module (b) is responsible for calculating top-explanation $\mathbb{E}_m^*$ for each segment $[p_{j'}, p_j]$. TSExplain employs [38] to identify top-explanation (Definition 3.5). In a nutshell, the CA algorithm simulates what a data analyst would perform during data analysis: recursively perform drill-down operations in dimensions and select data slices they are interested in. Each data slice is summarized by some conjunction predicate, i.e., an explanation in our context. Different from manual data analytics, here the drill-down dimensions and data slices are selected via dynamic programming to maximize the total diff score $\gamma(E)$ under the constraint that the number of data slices is $\leq m$.
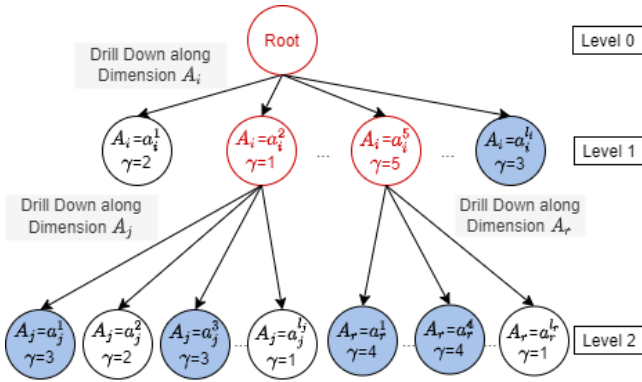


**Figure 8: Illustration of Cascading Analyst**

Given three explain-by attributes $\mathcal{A}=\{A_i, A_j, A_r\}$, Figure 8 illustrates the CA algorithm for identifying top-5 explanations. Each node in Figure 8 corresponds to an explanation $E$ and is associated with its diff score $\gamma(E)$ obtained from module (a). For instance, the left-most node at level one denotes explanation $E = (A_i=a_i^1)$ and the left-most node at level two denotes explanation $E = (A_i=a_i^2 \& A_j=a_j^1)$. To identify top-m non-overlapping explanations $\mathbb{E}_m^*$, the algorithm starts from the root node with $m$ quotas. It enumerates the first drill-down dimension, e.g., dimension $A_i$ as shown in Figure 8, as well as the quota assigned to each drill-down sub-tree, e.g., two out of five is assigned to the subtree rooted at node $(A_i=a_i^2)$ and $(A_i=a_i^5)$ respectively and another one is assigned to the subtree rooted at node $(A_i=a_i^{l_i})$. Again, for the sub-tree rooted at node $(A_i=a_i^2)$ with two quotas, the algorithm enumerates next dimension to drill down (i.e., $A_j$ in Figure 8), and assigns quota to each drill-down sub-tree (i.e., one to $(A_i=a_i^2 \& A_j=a_j^1)$ and one to $(A_i=a_i^2 \& A_j=a_j^3)$). This process is conducted recursively. The enumeration of drill-down dimension and quota assignment are performed via dynamic programming to maximize the total score $\sum_E \gamma(E)$ under the constraint that the total quota is $\leq m$. In Figure 8, the algorithm returns top-5 explanations (nodes in blue) with maximum $\sum_E \gamma(E)=3+3+4+4+3=17$. Please refer to [38] [3] for details.

---

[3][38] also works with the alternative Definition 3.5 in footnote 1.

*Time complexity:* the CA algorithm [38] takes $O(\epsilon \cdot |\mathcal{A}| \cdot m^2)$ per segment, where $\epsilon$ is the total number of candidate explanations, $|\mathcal{A}|$ is the number of explain-by attributes, and $m$ is a user-specified explanation number. By default, $m$ is set as 3. Since there are in total $n^2$ segments, the total time complexity is $O(\epsilon \cdot |\mathcal{A}| \cdot m^2 \cdot n^2)$

**K-Segmentation.** Module (c) is responsible for identifying the best K-segmentation scheme $\mathcal{P}_K^*$ that minimizes the total variance. As described in Section 4.1 (Eq. 6), we first compute $dist(P_k, P_x)$ based on the top-explanations $\mathbb{E}_m^*$ on centroid segment $P_k=[p_{j'}, p_j]$ and object $P_x=[p_x, p_{x+1}]$ obtained from module (b); next, $\text{var}(P_k)$ is calculated for every segment based on Eq. 7; lastly, DP is utilized for deriving $D(n, K)$ and the optimal segmentation scheme $\mathcal{P}_K^*$.

*Time complexity:* for each pair of $(P_k, P_x)$, calculating $dist(P_k, P_x)$ takes $O(m)$ in looking at top-$m$ explanations. There are $(n-l)$ centroid segments $P_k$ of length [4] $l$ where $1 \leq l < n$ and each $P_k$ of length $l$ contains $l$ objects $P_x$. In total, we have $\sum_{l=1}^{n-1} l(n-l) = O(n^3)$ pairs of $(P_k, P_x)$ and thus the time complexity for calculating all distance is $O(m \cdot n^3)$. Similarly, the total time complexity for calculating variance $\text{var}(P_k)$ of all $P_k$ is $O(n^3)$. With $\text{var}(P_k)$ available, DP takes $O(n \cdot nK) = O(n^2 K)$, as each step in Eq. 11 involves $O(n)$ for enumerating the last cut's position and there is in total $nK$ steps. Thus, the complexity of module (c) is $O(m \cdot n^3 + K \cdot n^2)$.

**Takeaway.** To sum up, the total time complexity is $O(\epsilon \cdot n^2 + \epsilon \cdot |\mathcal{A}| \cdot m^2 \cdot n^2 + m \cdot n^3 + K \cdot n^2)$ — it grows linearly to $\epsilon$ and $|\mathcal{A}|$, quadratic to $m$, and cubed to $n$. In general, $m$, $K$, and $|\mathcal{A}|$ are small due to user's limited perception. Thus, the runtime depends mostly on the number of candidate explanations $\epsilon$ and the size of the aggregated time series $n$. Next, we focus on reducing $\epsilon$ and $n$ for speedup.

## 5.3 Optimizations

As the takeaway above describes, the runtime largely depends on the number of candidate explanations $\epsilon$ and the time series length $n$. For Liquor dataset used in Section 7, $\epsilon$ is around 5000 even when only two explain-by attributes are considered, i.e., $|\mathcal{A}|=2$, and $n$ is around 300. Next, we introduce two optimizations for speedup: *(I.)* guess-and-verify to reduce $\epsilon$ and *(II.)* sketching to reduce $n$.

*5.3.1 Guess-and-verify.* As discussed in Section 5.2, the CA algorithm is one of the bottlenecks in our pipeline. Guess-and-verify is designed to reduce the input explanation number ($\epsilon$) in CA.

**High-level Idea.** To ensure non-overlapping, CA algorithm recursively drills down dimensions and selects data slices (explanations) as shown in Figure 8. The intuition behind guess-and-verify is that explanation $E$ with higher diff score $\gamma(E)$ is more likely to be in the top-$m$ non-overlapping explanations $\mathbb{E}_m^*$. Thus, instead of using all candidate explanations as the input of CA algorithm, we can take a *guess* and limit the input to only include top explanations with the highest diff score. A smaller input size can dramatically reduce the runtime of the CA algorithm, but the returned results $\overline{\mathbb{E}_m^*}$ may not be the optimal top-$m$ non-overlapping explanations $\mathbb{E}_m^*$. To mitigate this, we *verify* whether the returned result is optimal. This process is repeated until the result is verified to be optimal.

Specifically, we first sort all explanations in descending order of $\gamma(E)$, denoted as $\chi$. Next, we go through the following two phases

---

[4]A segment $[p_i, p_j]$ has length $(j - i)$.

iteratively: *(1) guess* and *(2) verify*. In each iteration, first take a guess on the input size $\bar{m}$; then run the CA algorithm with the first $\bar{m}$ explanations in $\chi$ and obtain the result $\overline{\mathbb{E}}$; verify if $\overline{\mathbb{E}}$ is optimal.
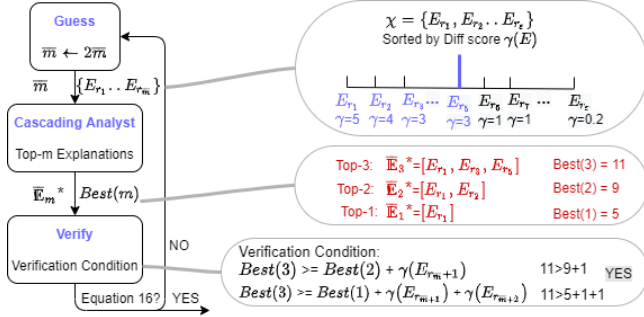


**Figure 9: Guess-and-Verify**

**Guess.** Naively selecting the first $m$ explanations in $\chi$ might result in explanations that overlap with each other, though the total difference score is the largest. Alternatively, we hypothesize that the answer of the CA algorithm comes from the top $\bar{m}$ candidate explanations in $\chi$. If the returned result passes the verification condition, terminate; otherwise, increase $\bar{m}$ to $2\bar{m}$ as shown in Figure 9.

**Verify.** Taking the first $\bar{m}$ explanations in $\chi$ as the input, the result $\overline{\mathbb{E}}_m^*$ returned by CA algorithm is not guaranteed to be optimal since some explanation $E \in \mathbb{E}_m^*$ might rank after $\bar{m}$ in $\chi$. To ensure the optimality, we design a sufficient condition such that once this condition is satisfied, it is guaranteed that the returned result is optimal. Let $\chi=[E_{r_1}, E_{r_2}.., E_{r_\epsilon}]$ be the ordered explanation list ranked by $\gamma(E)$. With $[E_{r_1}, E_{r_2}.., E_{r_{\bar{m}}}]$ as input, let $\overline{\mathbb{E}}_m^*$ be the top-$m$ non-overlapping explanations returned by CA algorithm and $Best[m]$ be the corresponding total difference score. Given $m$, the CA algorithm not only returns $Best[m]$, but also $Best[m']$ for every $1\leq m'<m$ as a side product of dynamic programming. $Best[0] = 0$. With these notations, we can now present the verification condition in Eq. 12. The high-level idea is that each explanation can be categorized into the following two classes: (1) with rank $\leq \bar{m}$, and (2) with rank $>\bar{m}$. Thus, we can upper bound the total difference score of each candidate $m$ *non-overlapping explanations* by the right-hand side of Eq. 12, where the first term corresponds to the score of class (1) and the second term upper bounds the score of class (2). Hence, if we ensure that the current best solution (left-hand side of Eq. 12) has higher total score than all candidate solutions (right-hand side of Equation 12), we can safely terminate and output the optimal top-explanation. Proofs are omitted due to space limitation.

$$Best[m] \geq Best[m'] + \sum_{1 \leq j \leq m-m'} \gamma(E_{r_{\bar{m}+j}}) \ \ \forall 0 \leq m' < m \quad (12)$$

*Time complexity:* `guess-and-verify` decreases the complexity from $O(\epsilon \cdot |\mathcal{A}| \cdot m^2 \cdot n^2)$ to $O(\bar{m} \cdot |\mathcal{A}| \cdot m^2 \cdot n^2)$ **at the best case**. Empirically, when $m = 3$, we initialize $\bar{m} = 30$.

*5.3.2 `Sketching`.* As discussed in Section 5.2, the time complexity in each module is at least quadratic to the time series size $n$. This is because each point is treated as a candidate cutting position in K-Segmentation and thus the total number of segments involved in each module is $O(n^2)$. Hence, reducing the number of candidate cutting positions can dramatically improve the efficiency. `Sketching` is designed for this purpose, as depicted in Figure 7.

**High-level Idea.** Given a time series with $n$ points, K-Segmentation aims to select $(K$-1) cutting points out of $(n$-2) non-endpoints. Our intuition is that some points are worse-suitable to be used as the cutting points and can be eliminated in a more cost-effective manner; next, since the remaining points is of a much smaller size, it is affordable to input them in our solution pipeline (Section 5.2). We call the remaining points `sketch`, as its role is to represent the original $n$ points in the given time series. In particular, `Sketching` consists of two phases: *(I.)* `sketch` selection and *(II.)* pipeline with `sketch`. We propose to select sketch using our proposed pipeline in Section 5.2, but with the constraint that each segment's length to be within $L$, where $L << N$.

**Sketch Selection.** There are two main requirements for sketch selection. First, the process should be efficient; second, the selected `sketch` should contain promising points for small-variance K-segmentation scheme (Eq. 7). Strawman approaches like random sampling or evenly spaced sampling are fast, but does not meet requirement two. To satisfy both, we propose to utilize our proposed pipeline in Section 5.2, but with some constraint. As we will detail below, the constraint is used to reduce the pipeline runtime (requirement one); and the pipeline is used to identify promising cutting points for minimizing the total variance in Eq. 7 (requirement two).

As discussed, the number of segments considered in our solution pipeline (Section 5.2) is $O(n^2)$. To alleviate this bottleneck, we can restrict each segment's length to be within $L$, where $L << N$. In this way, we only need to compute the diff score $\gamma(E)$ (module a), top-explanations $\mathbb{E}_m^*$ (module b), distance $dist(P_k, P_x)$ and variance $\text{var}P_k$ (module c) for segments with length $\leq L$. This reduces the total number of considered segments from $O(n^2)$ to $O(Ln)$. More specifically, to derive a `sketch` of size $|S|$, we set $K = |S|$ in our K-segmentation pipeline and set the maximum segment length threshold as $L$. Empirically, we set $L=\min(0.05N, 20)$ and $|S|=\frac{3n}{L}$.

*Time complexity:* In phase I (`sketch` selection), *module (a):* the time of computing the diff scores turns to $O(\epsilon \cdot L \cdot n)$; *module (b):* the time of CA algorithm turns to $O(\epsilon \cdot |\mathcal{A}| \cdot m^2 \cdot Ln)$; *module (c):* the time of computing distance and variance turns to $O(m \cdot L^2 \cdot n)$, and the time of DP turns to $O(L \cdot nK)$. **The total time complexity is reduced by at least $\frac{L}{n}$, compared to the pipeline without constraint.** In phase II (pipeline with `Sketch`), *module (a):* the time of computing the diff scores turns to $O((\epsilon \cdot |S|^2)$; *module (b):* the time of cascading analyst turns to $O(\epsilon \cdot |\mathcal{A}| \cdot m^2 \cdot |S|^2)$; *Module (c):* the time of computing distance and variance scores turns to $O(m \cdot |S|^2 \cdot n)$, and the time of DP turns to $O(K \cdot |S|^2)$. **The total time complexity is reduced by $(\frac{|S|}{n})^2$, compared to the pipeline without `sketch`.**

## 6 OPTIMAL SELECTION OF $K$

In real-world datasets, it is hard for users to specify the number of segments $K$ in advance. By varying segment number $K$, TSExplain outputs segmentation schemes with different variance scores, generating a $K$-$Variance$ curve. The left-hand side of Figure 11 illustrates

an example. Intuitively, $K$-$Variance$ curves decrease monotonically as the increase of $K$. At an extreme, when $K=n$-1, the total variance reaches a minimum score of zero. Furthermore, the total variance score drops quickly when $K$ is small and slows down when $K$ grows larger. This indicates that the marginal improvement of increasing $K$ becomes smaller when $K$ is large. Also, a larger $K$ brings about too many segments, which would exceed user perception limitations. Thus, our goal is to identify the optimal $K$ with relatively low variance and keep the segmentation scheme concise.

Such a task is well-studied in the machine learning community [23]. We borrow the idea of a well-known method named the "Elbow method" [23, 32] which picks the "elbow point" of the $K$-$Variance$ curve as the optimal $K$. We use a task-agnostic algorithm [40] to automatically determine the "elbow point" of our $K$-$Variance$ curve. This algorithm first normalizes the curve to be from $(0, 0)$ to $(1, 1)$. Then, it picks $K^*=\arg\max_K [\texttt{total\_var}(K) - K]$ as the "elbow point" where $\texttt{total\_var}(K)$ denotes the normalized total variance when the segment number is $K$.

In our implementation, we collect the dynamic programming results $D(n, K)$ varying $K$ from 1 to 20, plot the $K$-$Variance$ curve, and then choose the elbow point. We note that compared to calculating $D(n, K = 20)$, collecting $D(n, K)$ with varying $K$ from 1 to 20 does not add extra cost, since $D(n, K)$ gets generated for $1 \leq K < 20$ during the dynamic programming process of $D(n, K = 20)$. We constrain $K$ to be at most 20 due to user perception limitation: when $K$ is too large, e.g., $K \geq 20$, it would be hard for users to interpret the explanation results. We admit that when the time series is long, i.e., $n$ is large, restricting $K$ under 20 might not return the explanations at the finest granularity, but we argue that explanations at coarse grain with $K \leq 20$ is a better choice considering user perception limit. Empirically, we observe that TSExplain chooses 6 or 7 segments in most cases in our real-world experiments (Section 7.4).

## 7 EXPERIMENTS

In this section, we answer two questions: (1) how effective is TSExplain in identifying the evolving explanations; (Section 7.3 and 7.4) (2) how fast is TSExplain (Section 7.5).

The current TSExplain is implemented in C++. All experiments below are run single-threaded on a Macbook Air 2020 with Apple M1 chip 8-core CPU and 16GB memory.

### 7.1 Datasets

We introduce the datasets used in the experiments.

*7.1.1 Synthetic datasets.* We synthesize 20 different datasets with 7 different levels of $SNR_{dB}$ (section 4.2.1). In total, we have 140 different datasets. The aggregated time series is `SELECT T, count(sales) FROM R GROUP BY T` and the explain-by attribute is `category`.

*7.1.2 Real-world datasets.* We use the three real-word datasets in the motivating examples of Section 1. Here, we briefly go through the aggregated time series, explain-by attributes. In practice, we expect users to provide explain-by attributes based on their domain knowledge.

***Covid.*** [7] records the daily/total confirmed cases of 58 states. Naturally, there are two aggregated time series: ① the covid total confirmed cases trend, `SELECT date, SUM(total-confirmed-cases)` `FROM Covid GROUP BY date;` ② the covid daily confirmed cases trend, `SELECT date, SUM(daily-confirmed-cases) FROM Covid GROUP BY date`. We choose `state` as our explain-by attribute to answer "which `states` are the main contributors to the rises or drops?".

***S&P 500.*** contains 503 company[5] stock price (`price`) and free-float shares (`share`) from 2020-1-1 to 2020-10-1. Based on the S&P 500 index formula[5], we derive the S&P 500 index's time serie as `SELECT date,` $\frac{\texttt{SUM(price*share)}}{\texttt{divisor}}$ `AS SP500-index FROM Sp500 GROUP BY date`, where divisor is a constant. We try to explain the S&P 500 index's crashes and rebound using the hierarchical explain-by features - `category`, `subcatagory`, `stock`.

***Liquor.*** contains liquor purchase transactions in Iowa from 2020-1-2 to 2020-6-30. The time series is `SELECT date, SUM(Bottles Sold)` `FROM Liquor GROUP BY date`. We pick four attributes out of 24 attributes as our explain by features: `Bottle Volume (ml)` – the size of each bottle in a purchase (e.g., 750ml); `Pack` – the number of bottles per pack (e.g., 6); `Category Name` – the category of the purchased liquor (e.g., American Flavored Vodka); `Vendor Name` – the vendor of the purchased liquor (e.g. Phillips Beverage). Below, we use BV, P, CN, and VN to represent them respectively for short.

### 7.2 Baseline

TSExplain is the first explanation-aware segmentation to surface the evolving explanations for time series. The closest segmentation works are the bottom up algorithm (*Bottom-Up*) which performs best overall in piecewise linear approximation [21] and recent semantic segmentation algorithm - *FLUSS* [9] and *NNSegment* [42]. All these methods are explanation-agnostic, partition time series solely based on the visual shapes, and require segment number as input. For a fair comparison, the $K$ for baselines is either given or borrowed from TSExplain's results. Implementation wise, we reproduce *Bottom-Up* based on the pseudo-code in Keogh et al [21]. *FLUSS* is implemented using Stump library [26] and *NNSegment* is implemented using the authors' code [42].

### 7.3 Explanations of Synthetic Datasets

We evaluate TSExplain's effectiveness on synthetic datasets and perform quantitative comparisons between TSExplain and the three baselines. For fair comparison, we adopt the oracle segment number $K$ of the ground truth, and we run TSExplain and baselines with known $K$.

**Metric** We propose a metric to compare the effectiveness of these methods by quantifying the distance between these methods' output and ground truth. We calculate the edit distance between outputs and ground truth. Since different datasets have different segment number $K$ and time series lengths $n$, we normalize our edit distance by $K$ and $n$. The lower the metric is, the more effective the method is. We term this metric *distance percent*($\%$).

Figure 10 shows the comparison between TSExplain and baselines. As *FLUSS* and *NNSegment* both involves parameter, i.e. period and window size, we try mutiple parameters and report the best overall results we found. The x-axis is SNR, and the y-axis

---
[5]S&P 500 has 505 components and the components are adjusted along the time. We select the 503 companies that are in the component list during the whole period.

is the *distance percent(%)*. We report the average distance percent for each SNR level. As we can see, TSExplain always has the best performance than all three baselines and *Bottom-Up* is the most comparable baseline among all three. When SNR > 35, the *distance percent(%)* of TSExplain is close to 0, indicating for cleaner datasets, TSExplain's output is almost the same as ground truth. However, the baselines are incapable to detect ground truth even with clean datasets.
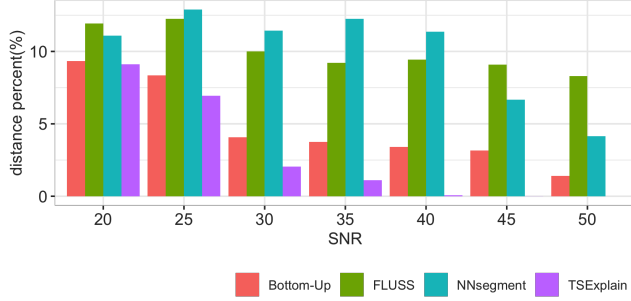


**Figure 10: Distance percentage of** TSExplain **and baselines**

**Takeaway.** For synthetic datasets, TSExplain performs much more effectively than all baselines. *Bottom-Up* is the most comparable baseline. For less noisy dataset, TSExplain can accurately detect the ground truth, while baselines can not.

## 7.4 Case Study on Real-World Datasets

This subsection demonstrates TSExplain's effectiveness on three different real-world datasets: Covid, S&P500, and Liquor, and gives an illustrative comparison between TSExplain and baselines. For fair comparison, TSExplain recommends the $K$ and our baseline uses the same segment number $K$. In this experiments, we focus on the top three explanations and set each explanation's order as 3. For very fuzzy datasets, we apply a moving average to smooth it before explaining it. The moving average window can be customized in the interface's panel as shown in the demo [6].

*7.4.1 Covid.* We explain the two time series in Covid dataset, `total-confirmed-cases` and `daily-confirmed-cases` separately.

`Total-confirmed-cases`. TSExplain identifies that the optimal $K$ equals 6 based on the optimal selection of K in Section 6. Figure 11 shows TSExplain's output segmentation scheme with the top three explanations' trend and three baselines' output. Please refer to the legend of Figure 2 for the top-3 explanations of TSExplain. In TSExplain's output, the first segment is from 1/22 to 3/14, where the increase of `total-confirmed-cases` is due to WA, NY, CA's increase. From 3/14 to 5/4, NY increases the most, followed by NJ and MA. Then, between 5/4 and 5/29, IL, CA, NY slowly increase. Later on, the confirmed cases surge mainly because of the sharp increase in CA, TX, FL, and IL, especially IL increases quickly

Contrarily, in the baselines, some neighboring segments' explanations are exactly the same, i.e., 6/16–7/31 and 7/31–11/2 in *Bottom-Up*, 6/28–8/9 and 8/9–10/3 in *NNSegment*, as shown in Figure 11c and Figure 11e; while *FLUSS* segmentation (Figure 11d) segments
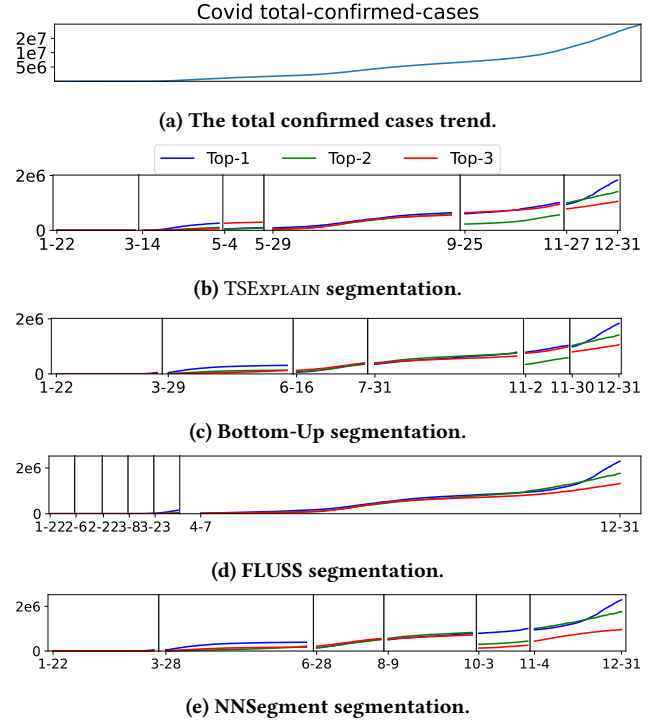


**Figure 11: Segmentation of total-confirmed-cases.**

the early time into a lot of small segments which is hard to interpret. In addition, none of the baselines detect the top-explanations' changes from NY during 3/14-5/4 to IL during 5/4-5/29 as reported in news[50].

| Segment | Top-1 Expl | Top-2 Expl | Top-3 Expl |
|---|---|---|---|
| 1/22 ~3/7 | Washington + | New York + | California + |
| 3/7 ~4/7 | New York + | New Jersey + | Massachusetts + |
| 4/7 ~5/25 | New York - | New Jersey - | California + |
| 5/25 ~7/16 | Florida + | Texas + | California + |
| 7/16 ~9/9 | Florida - | Texas - | California - |
| 9/9 ~11/10 | Illinois + | Texas + | Wisconsin + |
| 11/10 ~12/31 | California + | New York + | Illinois - |

**Table 3: Expl. in Fig. 12(middle). +/- denote change effect.**

`Daily-confirmed-cases`. In Figure 12 and table 3, TSExplain segments this time series into seven periods. Specifically, from 3/7 to 4/6, NY, NJ, and MA's rises contribute to the overall rise in the US. From 4/7 to 5/25, NY and NJ decline dramatically, and TSExplain captures an interesting pattern that CA starts to rise. During the holiday season in 2020, namely the last segment, CA and NY surge again while IL declines. Compared to TSExplain, the *Bottom-Up* segmentation does not detect the changes during 3/7 - 5/25. The *FLUSS* performs very poorly without revealing the explanation during 2/28-9/10 at all. The *NNSegment* segmentation is similar to our TSExplain explanation. However, it segments the up and down trend during 6/23 -8/26 to a whole segment which makes it users hard to interpret the explanations.

(a) The daily confirmed cases trend.

(b) TSExplain **segmentation.**

(c) Bottom-Up segmentation.

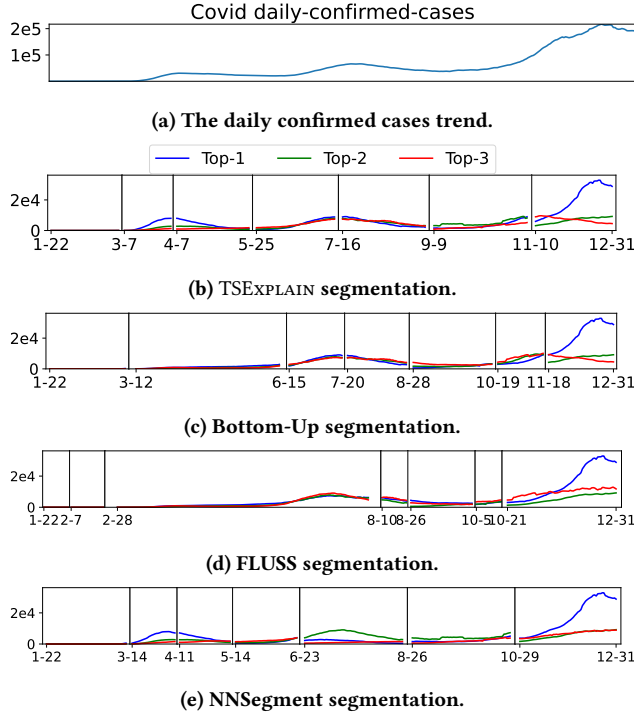(d) FLUSS segmentation.

(e) NNSegment segmentation.

**Figure 12: Segmentation of daily-confirmed-cases. Left is the K-Variance curve. Right shows aggregated time series(top), TSExplain results(middle), baseline results(bottom).**



(a) The S&P500 trend.

(b) TSExplain **segmentation.**

(c) Bottom-Up segmentation.

(d) FLUSS segmentation.
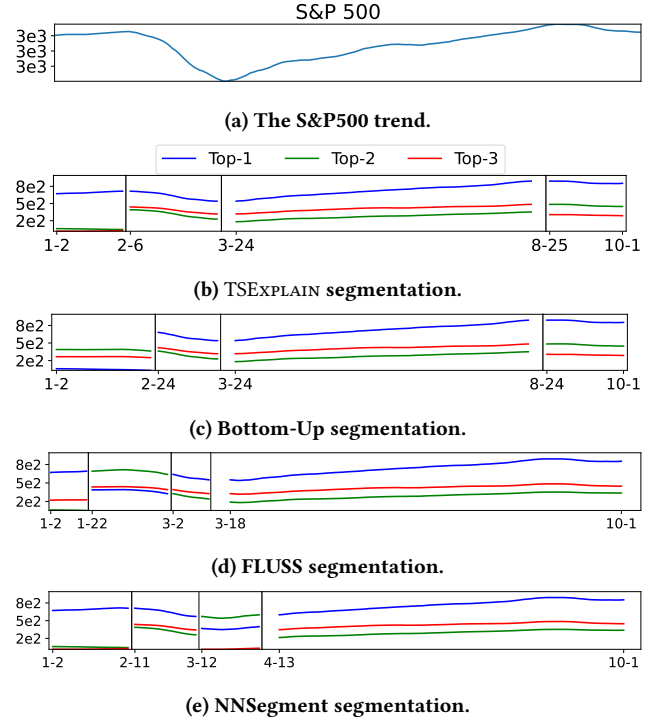
(e) NNSegment segmentation.

**Figure 13: Segmentation of S&P 500. Left shows the K-Variance Curve. Right shows aggregated time series(top), TSExplain results(middle), baseline results(bottom).**

| Segment | Top-1 Expl | Top-2 Expl | Top-3 Expl |
|---|---|---|---|
| 1/2 ~2/6 | technology + | energy - | *internet retail +* |
| 2/6 ~3/24 | technology - | financial - | communication - |
| 3/24 ~8/25 | technology + | consumer cyclical + | communication + |
| 8/25 ~10/1 | technology - | communication - | financial - |

**Table 4: TSExplain Explanations of S&P 500 in Figure 13(middle). All the explanations are related to attribute** `category`, **except** *internet retail +* **is related to** `subcategory`. **We omit the attribute name for short.**

*7.4.2 S&P 500.* Figure 13 and Table 4 show TSExplain explanation results. TSExplain finds the elbow point at four and segments the time series into four small segments. Before the market crash from 1/2 to 2/6, the S&P 500 rises mainly due to the rises of `category` technology and `subcategory` internet retail; meanwhile, `category` energy slightly drops. Also, TSExplain recognizes the market crash at 3/24. TSExplain explains that stocks belonging to `category` technology, financial, and communication contribute the most to the crash. After that, the `category` technology contributes most to the recovery during 3/24 ~8/25 and the drop during 8/25 ~10/1. We can also discover an interesting fact that financial `category` drops dramatically from 2/6 to 3/24, but does not bounce back a lot as technology and communication `category` do during the stock market recovery. In the *Bottom-Up* approach, although it detects the market crashing point on 3/24, however, the decreasing influencers

– technology, financial, communication service categories are detected starting from 2/24 which is much later than TSExplain. And *FLUSS* and *NNSegment* do not detect the crash point well and also can not recognize the drop from Aug to Oct.

*7.4.3 Liquor.* Figure 14 shows that TSExplain segments the time series into seven segments and Table 5 shows each segment's explanations. TSExplain recognizes that overall, from 1/20 to 4/21, the large pack liquor increases a lot, i.e., pack = 12, 24, 48. These indicate that people favor large-pack liquor at the beginning of the pandemic. We also remark that the underlying top explanations can be a combination of up and down trends. For example, the main reasons for the overall increase from 3/6 to 3/31 are `BV=1000(-)`, `BV=1750&P=6(+)` and `BV=750&P=12(+)`. TSExplain explains this in such way that ① `BV=1000` decreases sharply, otherwise the overall bottles sold can increase much more; ② `BV=1750&P=6(+)` and `BV=750&P=12(+)` directly contribute to the increase. Moreover, with some background knowledge, we find that liquor with `BV=1000` is mainly sold in independent stores. In March, Iowa's close down proclamation[17] requires restaurants and bars to shut down, and most indepen dent retailers rely on selling liquor to bars and restaurants. As a result, their business significantly declined. In late April, Iowa Governor issued a proclamation reopening restaurants, and the business of independent liquor stores gradually recovered. We can see TSExplain recognizes that from 4/21 to 5/8, `BV=1000ml&Pack=12` increases a lot and from 5/8 to 6/10, `BV=1000` becomes the top increasing explanation indicating that independent stores benefited from the

reopening policy. In comparison, the explanations of the *Bottom-Up* results look flat, as shown at the bottom of Figure 14(c), indicating the detected explanations change subtly. This is similar with *FLUSS* and *NNSegment*. Moreover, the top-2 and top-3 explanations during 3/25-4/17 returned by *FLUSS* have only subtle changes in bottles sold. For *NNSegment*, the explanations during 1/23-2/7 and 2/7-2/25 are exactly the same, and the changes made by top-2 and top-3 explanations of the last segment are very small as well. Also, we remark that although we specify four explain-by attributes, the results are only about BV and P. This indicates that TSExplain is able to identify interesting attributes and ignore the less interesting ones, i.e. VN and CN.



(a) The bottles sold trend.



(b) TSExplain **segmentation.**



(c) Bottom-Up Segmentation.



(d) FLUSS segmentation.



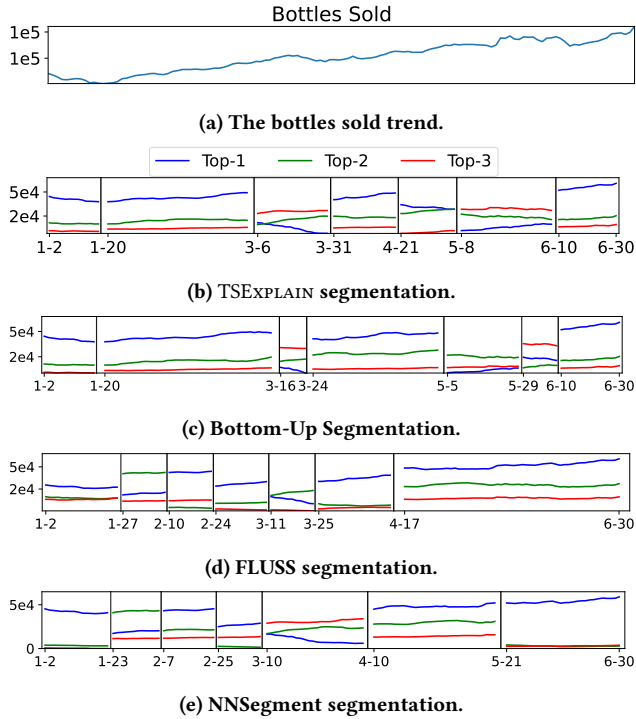(e) NNSegment segmentation.

**Figure 14: Segmentation of Liquor. Right shows the K-Variance Curve. Left shows aggregated time series(top), TSExplain results(middle), baseline results(bottom).**

| Segment | Top-1 Expl | Top-2 Expl | Top-3 Expl |
|---|---|---|---|
| 1/2 ~1/20 | P=12 - | P=6 - | BV=375&P=24 - |
| 1/20 ~3/6 | P=12 + | P=6 + | P=48 + |
| 3/6 ~3/31 | BV=1000 - | BV=1750&P=6 + | BV=750&P=12 + |
| 3/31 ~4/21 | P=12 + | BV=1750&P=6 - | P=24 + |
| 4/21 ~5/8 | BV=1750&P=12 - | P=6 + | BV=1000&P=12 + |
| 5/8 ~6/10 | BV=1000 + | BV=1750&P=6 - | BV=750&P=12 - |
| 6/10 ~6/30 | P=12 + | BV=1750&P=6 + | P=24 + |

**Table 5: TSExplain Explanations of Liquor in Fig. 12(middle).**

| dataset | $\epsilon$ | filtered $\epsilon$ | n |
|---|---|---|---|
| total-confirmed-cases | 58 | 54 | 345 |
| daily-confirmed-cases | 58 | 55 | 345 |
| S&P 500 | 610 | 329 | 151 |
| Liquor | 8197 | 1812 | 128 |

**Table 6: Real-world Dataset Statistics.**

*7.4.4 Sensitivity to K.* TSExplain adopts the elbow method to detect the optimal $K$. We observe that a slight change of the optimal $K$ will only bring up a slight shift in the results, e.g., remove or add one cutting point if K minuses/adds 1.

**Takeaway.** TSExplain effectively explains real-world datasets while the baseline defects in that:

- Less explanation diversity: the neighboring segments have the same explanations.
- Less effectiveness: the underlying key explanations can not be detected, and explanations are detected relatively late.
- Less significance: The explanations detected change subtly.

## 7.5 Efficiency Evaluation

We conduct three experiments to evaluate the efficiency. In the first experiment, we report the breakdown latency and study the impact of our optimization strategies. In the second experiment, we compare the end-to-end runtime between TSExplain and all three baselines. In the third experiment, we study the scalability of TSExplain.

*7.5.1 Latency Breakdown and Quality.*

**Methods.** Besides the two optimizations in Section 5.3, we introduce another straight-forward optimization - `filter`. The filtering protocol works as follows: given an explanation $E$, if each point in its aggregated time series has value smaller than a *ratio* of the corresponding value in the overall aggregated time series, we filter this explanation $E$ as its support is low and thus insignificant. We set *ratio* as 0.001 by default.

We study TSExplain with different optimizations: VanillaTSExplain (short as `Vanilla`) is the plain version without any optimization; `w filter` filters out the explanations below the default *ratio*; `O1` represents the algorithm which applies `guess-and-verify` after filtering; likewise, `O2` applies `sketching` and `O1+O2` applies both optimizations. Table 6 summarizes the statistics of the datasets. $\epsilon$ is the total candidate explanation number, $n$ is the time series length, and #record is the dataset record number. We also report the filtered $\epsilon$. We remark that in this experiment, $K$ is unspecified, and the time of selecting the optimal $K$ via the Elbow Method (Section 6) is included in the latency (Figure 15).

**Latency.** Figure 15 illustrates the breakdown latency of TSExplain. The overall latency breaks down into three parts – precomputation(blue), the cascading analysts algorithm(orange), K-segmentation(green) corresponding to three modules in section 5.2.

For Covid `total` and `daily-confirmed-cases` dataset, the filtering strategy only slightly improves since the predicate number $\epsilon$ does not change a lot before and after filtering. In contrast, the `sketching` in `O2` significantly reduces the latency of the cascading
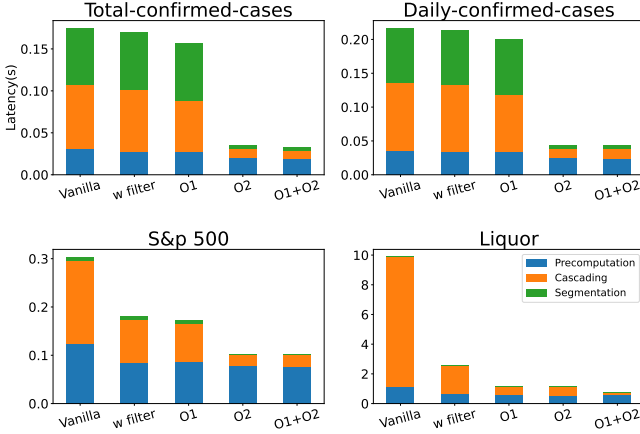
Figure 15: Latency of TSExplain.

| dataset | Variance(Vanilla) | Variance(O1+O2) |
|---------|-------------------|-----------------|
| total-confirmed-cases | 22.602 | 22.744 |
| daily-confirmed-cases | 91.619 | 91.994 |
| S&P 500 | 5.002 | 5.002 |
| Liquor | 33.6533 | 33.6533 |

Table 7: Quality of optimization strategies.

analysts algorithm and K-segmentation. Overall, O1+O2 reduces the latency of `total-confirmed-cases` from 175 ms to 33ms, and the latency of `daily-confirmed-cases` from 217ms to 43ms. For S&P 500 dataset, the filtering strategy reduces the latency to half. `Guess-and-verify` slightly reduces the latency, and together with `sketching` optimization, the overall latency is only 102 ms. For the Liquor dataset, the predicate number $\epsilon$ is very large even after filtering. The cascading analyst algorithm is the bottleneck. The `Vanilla` version takes 9.888s. After filtering, it still takes 2.59s. Since the predicate number $\epsilon$ is large, O1 (`Guess-and-verify`) takes big effect in this case. Each optimization O1 or O2 alone can significantly shrink the runtime to around 1.1s. The two optimizations together reduce the runtime to 756 ms. Although the precomputation is inevitable, the time of the cascading analyst algorithm has been reduced from `vanilla` — 8.754s to O1+O2 — around 200ms.

***Quality.*** Except `guess-and-verify`, `filter` and `sketching` both approximate the results without formal guarantee. Thus, we study the impact of optimizations in terms of the result quality.

Table 7 illustrates the segmentation scheme's total variance of O1+O2 compared with the `Vanilla` version. The variances and output segmentation of both algorithms are exactly the same for S&P 500 and Liquor datasets. For the Covid datasets, the difference is less than 1% and in the output segmentation, only two cutting points are slightly different — the corresponding distances are less than four days. Such a small discrepancy demonstrates that our optimization's effect on result quality is neglectable.

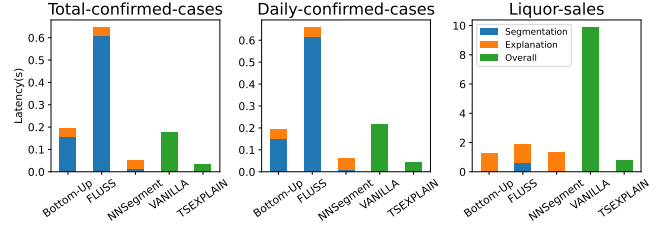### 7.5.2 End-to-end Efficiency Comparison with Baselines.



Figure 16: Efficiency comparison with baselines.

***Methods.*** The three baselines in Section 7.2 solely focus on visual shape based segmentation without providing any explanations and require the segmentation number as an input. To make them comparable, first, after segmenting using each baseline, we add the explanation module using the CA algorithm in Section 5.2; second, we reuse the optimal $K$ TSExplain finds in Section 7.4, and then run all baselines and TSExplain with this given optimal $K$. We also remark that the latency of determining the optimal K is very low, around 2ms in our experiments.

***Results.*** Figure 16 reports the end-to-end efficiency comparison. For each baseline, we report the segmentation and explanation time separately, while for TSExplain, we show the overall time since our segmentation module interleaves with the explanation module. To illustrate the effectiveness of our proposed optimizations, we also report VanillaTSExplain (VANILLA). We can tell that for different datasets, *FLUSS* is always the slowest, *NNSegment* and *Bottom-Up* rank in the middle. VanillaTSExplain is similar to the *Bottom-Up* on the *COVID-19* datasets and becomes slow when the predicate number goes up in the *Liquor-sales* dataset. Yet, combined with all the proposed optimizations, TSExplain is the fastest compared to all baselines on all datasets.

### 7.5.3 Scalability.

In the scalability experiment, we synthesized new time series following the procedure in Sec VII.A of different lengths = 100, 200, 400, 800, 1600, 3200, 6400. For each length, we synthesize five different time series. We explain these time series using VanillaTSExplain (without any optimizations) and TSExplain with all optimizations mentioned in the paper. Figure 17 reports the average latency of different time series lengths. We terminate when the latency is greater than 100s. VanillaTSExplain's latency increases exponentially. With optimizations, TSExplain's latency increases much slower when the time series becomes longer. In particular, TSExplain can interactively explain the synthetic time series with length = 3200 in 982 ms.

**Takeaway.** TSExplain can explain these three time series within 800ms, and our optimizations have accelerated the running time up to 13× with neglectable effects on quality. What's more, TSExplain is faster than all the baselines.

## 8 DISCUSSION

**Time-varying Attribute.** Different from non-temporal attributes, time-varying attributes are those whose values may change over
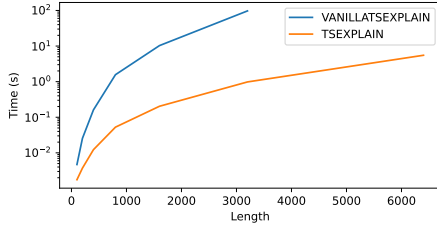
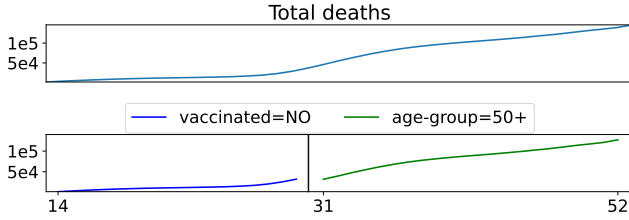**Figure 17: Different latency for time series of different lengths.**



**Figure 18: Segmentation of weekly total deaths. Overall trend(top).** TSExplain **results(bottom).**

time [18]. Augmenting dataset with time-varying attributes can sometimes add insights in explaining trends. Below we demonstrate an example of using TSExplain with time-varying attributes.

The covid death dataset [4] records weekly deaths of different age groups and different vaccinated status from week 14 to 52 in 2021. The `vaccinated` attribute is a time-varying attribute since a person of `vaccinated=NO` can shift to `vaccinated=YES` in the near future. However, `age-group` is a static attribute because one person belonging to a specific age group will not change within a year. We explain the total death trend in Figure 18(top) using the attributes `age-group` and `vaccinated`. Figure 18(bottom) shows the result that before week 31, the main contributor is the unvaccinated people and after week 32, the main contributor shifts to elder people with `age-group=50+`. Combined with human knowledge that the vaccinated population is gradually becoming large along the time, we can get the insight that at the beginning, unvaccinated people (including unvaccinated elders) are the major factor of total deaths since unvaccinated young people also face high risk. Later on, elder people no matter vaccinated or not are the major reason as more young people get protected from vaccines while elder people do not get protected that well even vaccinated.

**Real-time Time Series.** We briefly discuss how TSExplain can be extended to support real-time time series explanation. TSExplain first gives users the segmentation results of existing time series and meanwhile, caches all unit segments' top explanations. When new data arrives, it incrementally computes the top explanations for the new time series, runs the segmentation algorithm based on the existing time series' cutting point and newly arrived data points, and updates the segmentation results.

**Seasonal Datasets** For seasonality datasets, TSExplain can explain the seasonality dataset directly and detect the repeated pattern of evolving explanations which indicates the periodicity property. Users can also first decompose the seasonal datasets [15] and explain the seasonality and trend separately.

## 9 CONCLUSION

This work introduces TSExplain, the first explanation engine that identifies the evolving explanations for aggregated time series. We formulate the problem for deriving evolving explanations as a *K-Segmentation* problem, aiming to partition the input time series into K smaller segments such that each period has consistent top explanations. We propose a novel variance metric to quantify the consistency in each segment and develop a dynamic programming algorithm for identifying the optimal *K-Segmentation* scheme. TSExplain can automatically identify the optimal *K* using the "elbow method". In the experiments, we show TSExplain can effectively discover the evolving explanation on synthetic and real-world datasets. We propose optimizations that enable TSExplain to answer all our queries interactively within one second. Several future work directions include extending the difference metric library, recommending explain-by attributes, adding hints for segments with higher variance for further inspection.

# REFERENCES

[1] Firas Abuzaid, P. Kraft, Sahaana Suri, Edward Gan, E. Xu, Atul Shenoy, Asvin Anathanaraya, John Sheu, E. Meijer, Xi Wu, J. Naughton, Peter Bailis, and M. Zaharia. 2018. DIFF: A Relational Interface for Large-Scale Data Explanation. *Proc. VLDB Endow.* 12 (2018), 419–432.

[2] Peter Bailis, Edward Gan, S. Madden, D. Narayanan, Kexin Rong, and Sahaana Suri. 2017. MacroBase: Prioritizing Attention in Fast Data. *Proceedings of the 2017 ACM International Conference on Management of Data* (2017).

[3] Berit Hoffmann. 2021. Unveiling Sisu's vision for Decision Intelligence. https://sisudata.com/blog/unveiling-vision-decision-intelligence.

[4] CDC COVID-19 Response, Epidemiology Task Force. 2022. Rates of COVID-19 Deaths by Age Group and Vaccination Status. https://data.cdc.gov/Public-Health-Surveillance/Rates-of-COVID-19-Cases-or-Deaths-by-Age-Group-and/3rge-nu2a.

[5] CHAD LANGAGER. 2021. How Is the Value of the S&P 500 Calculated? https://www.investopedia.com/ask/answers/05/sp500calculation.asp. [Online; accessed 5-October-2021].

[6] Yiru Chen and Silu Huang. 2021. TSExplain: Surfacing Evolving Explanations for Time Series. In *Proceedings of the 2021 International Conference on Management of Data.* 2686–2690.

[7] Ensheng Dong, Hongru Du, and Lauren Gardner. 2020. An interactive web-based dashboard to track COVID-19 in real time. *The Lancet infectious diseases* 20, 5 (2020), 533–534.

[8] David H Douglas and Thomas K Peucker. 1973. Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *Cartographica: the international journal for geographic information and geovisualization* 10, 2 (1973), 112–122.

[9] Shaghayegh Gharghabi, Yifei Ding, Chin-Chia Michael Yeh, Kaveh Kamgar, Liudmila Ulanova, and Eamonn Keogh. 2017. Matrix profile VIII: domain agnostic online semantic segmentation at superhuman performance levels. In *2017 IEEE international conference on data mining (ICDM).* IEEE, 117–126.

[10] Shaghayegh Gharghabi, Chin-Chia Michael Yeh, Yifei Ding, Wei Ding, Paul Hibbing, Samuel LaMunion, Andrew Kaplan, Scott E Crouter, and Eamonn Keogh. 2019. Domain agnostic online semantic segmentation for multi-dimensional time series. *Data mining and knowledge discovery* 33, 1 (2019), 96–130.

[11] Google Trend. 2020. Taylor Swift compared with Kim Kardashian. https://trends.google.com/trends/explore?q=%2Fm%2F0dl567,%2Fm%2F0261x8t&date=now%207-d&geo=US.

[12] Jim Gray, Surajit Chaudhuri, Adam Bosworth, Andrew Layman, Don Reichart, Murali Venkatrao, Frank Pellow, and Hamid Pirahesh. 1997. Data cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals. *Data mining and knowledge discovery* 1, 1 (1997), 29–53.

[13] John A Hartigan and Manchek A Wong. 1979. Algorithm AS 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)* 28, 1 (1979), 100–108.

[14] Jim Hunter and Neil McIntosh. 1999. Knowledge-based event detection in complex time series data. In *Joint European Conference on Artificial Intelligence in Medicine and Medical Decision Making.* Springer, 271–280.

[15] Rob J Hyndman and George Athanasopoulos. 2018. *Forecasting: principles and practice.* OTexts.

[16] Imply. 2021. Imply Explain Feature. https://docs.imply.io/latest/explain/.

[17] Iowa Government. 2020. Gov. Reynolds issues a State of Public Health Disaster Emergency. https://idph.iowa.gov/News/ArtMID/646/ArticleID/158309/Gov-Reynolds-issues-a-State-of-Public-Health-Disaster-Emergency-31720?fbclid=IwAR1gjXuRQH37snyHWqL7-GlI4DOsRLAliqd94DD7m_sJaJr38lSOC4I7pwE. [Online; accessed 5-October-2021].

[18] Christian S Jensen and Richard T Snodgrass. 1996. Semantics of time-varying information. *Information Systems* 21, 4 (1996), 311–352.

[19] Manas Joglekar, Hector Garcia-Molina, and Aditya Parameswaran. 2017. Interactive data exploration with smart drill-down. *IEEE Transactions on Knowledge and Data Engineering* 31, 1 (2017), 46–60.

[20] Johns Hopkins University. 2021. COVID-19 Data Repository. https://github.com/CSSEGISandData/COVID-19.

[21] Eamonn Keogh, Selina Chu, David Hart, and Michael Pazzani. 2004. Segmenting time series: A survey and novel approach. In *Data mining in time series databases.* World Scientific, 1–21.

[22] Eamonn J Keogh, Padhraic Smyth, et al. 1997. A probabilistic approach to fast pattern matching in time series databases.. In *Kdd*, Vol. 1997. 24–30.

[23] Trupti M Kodinariya and Prashant R Makwana. 2013. Review on determining number of Cluster in K-Means Clustering. *International Journal* 1, 6 (2013), 90–95.

[24] Antti Koski, Martti Juhola, and Merik Meriste. 1995. Syntactic recognition of ECG signals by attributed finite automata. *Pattern Recognition* 28, 12 (1995), 1927–1940.

[25] Jokin Labaien, Ekhi Zugasti, and Xabier De Carlos. 2020. Contrastive explanations for a deep learning model on time-series data. In *International Conference on Big Data Analytics and Knowledge Discovery.* Springer, 235–244.

[26] Sean M Law. 2019. STUMPY: A powerful and scalable Python library for time series data mining. *Journal of Open Source Software* 4, 39 (2019), 1504.

[27] Chenjie Li, Zhengjie Miao, Qitian Zeng, Boris Glavic, and Sudeepa Roy. 2021. Putting Things into Context: Rich Explanations for Query Answers using Join Graphs. In *Proceedings of the 2021 International Conference on Management of Data.* 1051–1063.

[28] Zhicheng Liu and Jeffrey Heer. 2014. The Effects of Interactive Latency on Exploratory Visual Analysis. *IEEE Trans. Vis. Comput. Graph.* 20, 12 (2014), 2122–2131. https://doi.org/10.1109/TVCG.2014.2346452

[29] Yasuko Matsubara, Yasushi Sakurai, and Christos Faloutsos. 2014. Autoplait: Automatic mining of co-evolving time sequences. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data.* 193–204.

[30] Zhengjie Miao, Andrew Lee, and Sudeepa Roy. 2019. LensXPlain: Visualizing and explaining contributing subsets for aggregate query answers. *Proceedings of the VLDB Endowment* 12, 12 (2019), 1898–1901.

[31] Microsoft. 2021. PowerBI. https://powerbi.microsoft.com/en-us/.

[32] Andrew Ng. 2012. Clustering with the k-means algorithm. *Machine Learning* (2012).

[33] PowerBI. 2021. Create Key Influencers Visualizations. https://docs.microsoft.com/en-us/power-bi/visuals/power-bi-visualization-influencers.

[34] PowerBI. 2021. Key Influencer Over Time. https://ideas.powerbi.com/ideas/idea/?ideaid=57440365-96af-4362-9b8f-5d096bd92788.

[35] Urs Ramer. 1972. An iterative procedure for the polygonal approximation of plane curves. *Computer graphics and image processing* 1, 3 (1972), 244–256.

[36] Clayton Rooke, Jonathan Smith, Kin Kwan Leung, Maksims Volkovs, and Saba Zuberi. 2021. Temporal Dependencies in Feature Importance for Time Series Predictions. *arXiv preprint arXiv:2107.14317* (2021).

[37] Sudeepa Roy and Dan Suciu. 2014. A formal approach to finding explanations for database queries. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data.* 1579–1590.

[38] M. Ruhl, M. Sundararajan, and Qiqi Yan. 2018. The Cascading Analysts Algorithm. *Proceedings of the 2018 International Conference on Management of Data* (2018).

[39] Sunita Sarawagi. 2001. idiff: Informative summarization of differences in multidimensional aggregates. *Data Mining and Knowledge Discovery* 5, 4 (2001), 255–276.

[40] Ville Satopaa, Jeannie Albrecht, David Irwin, and Barath Raghavan. 2011. Finding a" kneedle" in a haystack: Detecting knee points in system behavior. In *2011 31st international conference on distributed computing systems workshops.* IEEE, 166–171.

[41] Sisu Data. 2021. Augmented Analytics for Strategic Data Teams. https://sisudata.com/.

[42] Torty Sivill and Peter Flach. 2022. LIMESegment: Meaningful, Realistic Time Series Explanations. In *International Conference on Artificial Intelligence and Statistics.* PMLR, 3418–3433.

[43] Tableau. 2021. BI and Analytics Software. https://www.tableau.com/.

[44] Tableau. 2021. Explain Data. www.tableau.com/products/new-features/explain-data.

[45] Sana Tonekaboni, Shalmali Joshi, Kieran Campbell, David K Duvenaud, and Anna Goldenberg. 2020. What went wrong and when? Instance-wise feature importance for time-series black-box models. *Advances in Neural Information Processing Systems* 33 (2020), 799–809.

[46] HJLM Vullings, MHG Verhaegen, and Henk B Verbruggen. 1997. ECG segmentation using time-warping. In *International Symposium on Intelligent Data Analysis.* Springer, 275–285.

[47] Xiaolan Wang, X. Dong, and A. Meliou. 2015. Data X-Ray: A Diagnostic Tool for Data Errors. In *SIGMOD '15.*

[48] Wikipedia, the free encyclopedia. 2021. Signal-to-noise ratio. https://otexts.com/fpp2/classical-decomposition.html.

[49] Wikipedia, the free encyclopedia. 2021. S&P500 Index. https://en.wikipedia.org/wiki/S%26P_500.

[50] WTTW. 2020. Illinois Seeing More and More COVID-19 Cases as Testing Continues to Increase. https://news.wttw.com/2020/05/02/illinois-seeing-more-and-more-covid-19-cases-testing-continues-increase.

[51] E. Wu and S. Madden. 2013. Scorpion: Explaining Away Outliers in Aggregate Queries. *Proc. VLDB Endow.* 6 (2013), 553–564.