



**Utrecht
University**

Department of Mathematics and Computer Science
Process Analytics

The generation of interpretable counterfactual examples by finding minimal edit sequences using event data in complex processes

Master Thesis

Olusanmi Hundogan

Supervisors:

Xixi Lu

Yupei Du

March 1, 2022

Abstract

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Contents

List of terms	3
1 Introduction	6
1.1 Context of this Thesis	6
1.2 Problem Space	7
2 Background	8
2.1 Process Mining	8
2.1.1 A definition for Business Processes	8
2.1.2 What is Process Mining	10
2.1.3 The Challenges of Process Mining	10
2.2 Multivariate Time-Series Modelling	11
2.2.1 What are Time Series Models?	11
2.2.2 The Challenges of Time Series Modelling	12
2.2.3 A Process within the State-Space Framework	13
2.3 Counterfactuals	14
2.3.1 What are Counterfactuals?	14
2.3.2 The Challenges of Counterfactual Sequence Generation	15
2.4 Related Literature	16
2.4.1 The Generation of Counterfactual Explanations	16
2.4.2 Generating Sequential Counterfactual Explanations	16
2.4.3 Counterfactual Reasoning for Business Processes	17
2.4.4 The Construction of Counterfactual Time Series	18
2.5 Formal Definitions	19
2.6 Research Question	19
2.7 General Approach	19
2.8 What is Process Mining?	19
2.9 Challenges of Processing Process Data	19
3 Related Papers	20

4	Methods	21
4.1	Datasets	21
4.2	Preprocessing	21
4.3	Framework	21
5	Results	22
5.1	Evaluation	22
6	Discussion	23
7	Conclusion	24

List of terms

BI Business Intelligence. 3, *see*: Business Intelligence

BPM Business Process Management. 3, 10, *see*: Business Process Management

Business Intelligence XXX. 3

Business Process Management XXX. 3, 10

Causal Inference A discipline which seeks to incorporate causal relations to model phenomena in the real world.. 17

Comma Separated Values A structured data format to store information. Every line relates to a data point and every feature is separated by a separator. The separator is commonly a comma but other characters like tabs or semicolons are valid as well.. 3, 10

Continuous Process Improvement XXX. 3, 10

Corporate Performance Management XXX. 3

CPI Continuous Process Improvement. 3, 10, *see*: Continuous Process Improvement

CPM Corporate Performance Management. 3, *see*: Corporate Performance Management

CSV Comma Separated Values. 3, 10, *see*: Comma Separated Values

Data Mining XXX. 10

Deep Learning A sub-discipline of machine learning which focuses on neural networks as primary tool. The discipline emphasises the research and development of neural network architectures. Data preprocessing and feature engineering play a secondary role.. 16, 17

DKF Deep Kalman Filter. 17

DMM Deep Markov Model. 17

ELBO Evidence Lower-Bound. 17

Event Log A collection of event data, that's produced by the process. They are the main input of every process mining venture.. 10, 13

eXtensible Event Stream An XML-based data format to store event logs. The format was developed and adopted by the IEEE Task Force on Process Mining.. 5, 10

GAN Generative Adversarial Model. 4, 16, 17, *see*: Generative Adversarial Model

Generative Adversarial Model A model in which two neural network train simultaneously. The generative model tries to generate instance examples, while the discriminative models tries to distinguish real examples from generated ones. Both, the discriminator and the generator can be used in isolation afterwards.. 4, 16

HMM Hidden Markov Model. 17

Information System XXX. 4, 10

IS Information System. 4, 10, *see*: Information System

Markov Decision Process A probabilistic process which assumes that an agent can influence the outcome of the process by choosing decisions given a state and expecting a return for its performance.. 4, 18

Markov Process A probabilistic process whose outcomes depend on the process' state.. 4, 18

MDP Markov Decision Process. 4, 18, *see*: Markov Decision Process

ML Machine Learning. 6

MP Markov Process. 4, 18, *see*: Markov Process

Natural Language Processing A discipline that is mainly concerned with the analysis and modelling of natural language.. 5, 11

NLP Natural Language Processing. 5, 11, 16, 17, *see*: Natural Language Processing

PM Process Mining. 5, 17, *see*: Process Mining

Process Event Also called activities. A discrete step in the process.. 9

Process Instance Also called case. A collection of activities that belong to a common entity that is produced by the process.. 9–11, 13, 16

Process Mining A subdiscipline of Data Mining, which uses process logs to analyse and utilise data which was produced by processes.. 5, 17

Rashomon Effect Rashomon is a classic Japanese movie in which multiple witnesses tell a different equally valid story about the murder of a samurai. Although, each story acts as a valid explanation, they contradict each other. The same effect may apply to equally valid counterfactuals.. 15

Reinforcement Learning An are within Machine Learning, which seeks to allow an intelligent agent to choose the right actions by maximizing the cumulative rewards of a task setting.. 5, 18

RL Reinforcement Learning. 5, 18, *see*: Reinforcement Learning

SCM Structural Causal Model. 5, 18, *see*: Structural Causal Model

Structural Causal Model Also called *Structural Equation Model (SEM)*. An SCM is a set of variables and equations. The equations' the relationship of outputs from other equations and the variables. The relationship between functions and variables inform the causal relationship. SCMs can be represented as directed graphs in which nodes describe variables and equations and the vertices dependencies. Furthermore, a SCM can produce a dataset containing all the values thhat where generated.. 5, 18

Total Quality Management XXX. 5

TQM Total Quality Management. 5, *see*: Total Quality Management

XAI eXplanable AI. 6, 12, 14, 15, 18

XES eXtensible Event Stream. 5, 10, *see*: eXtensible Event Stream

Chapter 1

Introduction

1.1 Context of this Thesis

Many processes, often medical, economical, or administrative in nature, are governed by sequential events and their contextual environment. Many of these events and their order of appearance play a crucial part in the determination of every possible outcome. With the rise of AI and the increased abundance of data in recent years several techniques emerged that help to predict the outcomes of complex processes in the real world. **[Expand the domain application.]**

For instance, research in the Process Mining discipline has shown that is possible to predict the outcome of a particular process fairly well **CITE** . **[However, while many prediction models can easily certain outcomes, it remains a difficult challenge to understand what led to a particular outcome. This obstacle is undesirable, as knowing the main factors to an outcome can help understand how to steer a process to a desired outcome with minimal effort.]** In other words, we want to change the outcome of a particular event, by making it maximally likely, with as little interventions as possible **CITE TEST** .

One-way to better understand the Machine Learning (ML) models lies within the eXplainable AI (XAI) discipline. XAI dedicates its research to the **research and** development of so-called *black-box models* that are difficult to interpret. Most of the discipline's techniques produce explanations that guide our understanding.

A prominent and human-friendly approach uses the generation of counterfactuals as primary explanation tool. Counterfactuals within the AI framework help us to answer hypothetical "what-if" questions. In this thesis, we will raise the question, how we can use counterfactuals to change the trajec-

tory of a models' prediction towards a desired outcome. Knowing the answers will help us further understand what to do to avoid or enforce the outcome of a process. **[WHY]**

1.2 Problem Space

In this paper, we will approach the problem of generating counterfactuals for processes. The literature has provided a multitude of techniques to generate counterfactuals for AI models, that are derived from static data¹. However, little research has focussed on counterfactuals for dynamic data². A major reason, emerges from a **[multitude – better #]** of challenges, when dealing with counterfactuals and sequences. First, counterfactuals within AI attempt to explain outcomes, that did not happen. Therefore, there is no evidence data, from which one can infer predictions. Subsequently, this lack of evidence further complicates the evaluation of generated counterfactuals. In other words, you cannot validate the correctness of a theoretical outcome that has never occurred. Second, sequential data is not only has a highly variable form, too **CITE**. The sequential nature of the data impedes the tractability of many problems due to the combinatorial explosion of possible sequences which depends on the length of the sequence. Third, process data of requires knowledge of the underlying and often hidden causal structures that produce the data in the first place. However, these structures are often hidden and it is a NP-hard problem to elicit them **CITE Check process discovery literature**. Furthermore, the data generated is seldomly one-dimensional or discrete. Henceforth, each dimension's contribution can vary in dependance of its context, the time and magnitude. Hence, the field in which we can contribute to this open challenge is vast. As a result, we have to restrict the solution space by imposing limitations and assumptions. Therefore, the result of this paper will describe a framework that will only apply to a subset of problems. In the following sections, we will explore these restrictions by describing the most important concepts in chapter 2.

¹With static data, we refer to data that does not change over a time dimension.

²With dynamic data, we refer to data that has time as a major component, which is also inherently sequential

Chapter 2

Background

This chapter will explore the most important concepts for this work. Most of the concepts can have several meanings depending on the varying context in which they are applied. For this purpose, we will provide an intuitive understanding, the ensuing challenges, a concrete definition for this work and lastly and a mathematically formal description. The concepts we will cover encompass [sequence modelling](#), process mining and counterfactual explanations.

2.1 Process Mining

2.1.1 A definition for Business Processes

Before elaborating on Process Mining, we have to establish the meaning of the term *process* in the context of this paper. The term is broadly used in many contexts and therefore has a rich semantic volume. A process generally refers to something that advances and changes over time[6]. Although, legal or biological processes may be valid understandings, we focus on processes *business processes*.

An example is a loan application process in which an applicant may request a loan at a specific point in time. The case is then assessed and reviewed by multiple approvers and ends in a final decision. The loan may be granted or denied. The *business* part may be misleading as these processes are not confined to commercial settings. For instance, a medical business process may cover a patient's admission to a hospital, followed by a series of diagnostics and treatments and ending with the recovery or death of a patient. Another example from a human-computer-interaction [\[Add to glossary\]](#) perspective would be an order process for an online retail service like Amazon. The buyer

might start the process by adding articles to the shopping cart and proceeding with specifying their bank account details. This order process would end with the submission or receipt of the order.

All of these examples have a number of common characteristics. They have a clear starting point which is followed by numerous intermediary steps and end in one of the possible sets of outcomes. For this paper we will refer to each step, including start and end points, as Process Event. Each Process Event may contain additional information in the form of event attributes. A collection of these Process Events refer to a Process Instance, if they all relate to a single run of a process. In line with the aforementioned examples, these Process Instances could be understood as a single loan application, a medical case or a buy order. We can also attach Process Instance related information to each instance. Examples would be the applicants race, a patients age or the buyers budget. In its entirety, a business process can be summarised as a *graph* or *flowchart*, in which every node represents an event and each arc the path to another event. This graphical representation is referred to as *process map*. Figure 2.1 shows an example of such a representation.

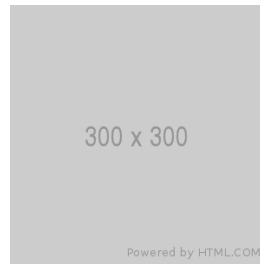


Figure 2.1: This graph shows an example of various process maps.

In conclusion, in this thesis a *business process* refers to

A finite series of discrete events with one or more starting points, intermediary steps and end points.

However, we have to address a number of issues with this definition. First, this definition excludes infinite processes like [XXX] or continuous processes such as [XXX]. There may be valid arguments to include processes with these characteristics, but they are not relevant for this thesis. Second, in each example we deliberately used words that accentuate modality such as *may*, *can* or *would*. It is important to understand that each process anchors its definition in its application context. Hence, what defines a business process is indisputably subjective. For instance, while an online marketplace like Amazon might be interested in the process from the customers first click

to the successful shipment, an Amazon vendor might be interested in the delivery process of a product only. Third, the example provided in Figure 2.1 may not relate to the reality of a data generating process. In line with the second point, these examples subjective models of a process. They may or may not be accurate. The *true* process is often unknown to every actor. Therefore, we will distinguish between the *true process model* and a *process model*. The *true process model* is a hypothetical concept whose *true* structure remains unknown.

2.1.2 What is Process Mining

Having established our understanding of a process, we can turn towards *Process Mining*. This young discipline has many connections to other fields that focus on the modeling and analysis of processes such as Continuous Process Improvement (CPI) or Business Process Management (BPM). However, its data-centric approaches originate in Data Mining. The authors W. van der Aalst et al. describe this field as a discipline “to discover, monitor and improve real processes (i.e., not assumed processes) by extracting knowledge from event logs readily available in today’s (information) systems” [24]. The discipline revolves around the analysis of Event Logs. An Event Log is a collection of Process Instances. These logs are retrievable from various sources like an Information Systems (ISs) or database. Those logs are often stored in data formats such as Comma Separated Values (CSV) or eXtensible Event Stream (XES).

2.1.3 The Challenges of Process Mining

As mentioned in chapter 1, process data modelling and analysis is a challenging task. W. van der Aalst et al. mentions a number of issues that arise from processes [24].

The first issue arises from the quality of the data set. Process logs are seldomly collected with the primary goal of mining information and hence, often appear to be of subpar quality. The information is often incomplete due to a lack of context information, the omission of logged process steps or wrong levels of granularity.

This issue is exacerbated by the second major issue with process data. Mainly, its complexity. Not only does a process logs complexity arise from the variety of data sources and differing levels of complexity, but also from the data’s characteristics. The data can often be viewed as multivariate sequence with discrete and continuous features and variable length. This characteristic alone creates problems explored in section 2.2. **[Also refer to**

variability in sequence section.] However, the data is also just a *sample* of the process. Hence, it may not reflect the real process in its entirety. In fact, mining techniques need to incorporate the *open world assumption* as the original process may generate unseen Process Instances.

A third issue which contributes to the datasets incompleteness and complexity is a phenomenon called *concept drift*. This phenomenon relates possibility of a change in the *true* process. The change may occur suddenly or gradually and can appear in isolation or periodically. An expression of such a drift may be a sudden inclusion of a new process step or domain changes of certain features. These changes are not uncommon and their likelihood increases with the temporal coverage and level of granularity of the dataset **CITE** . In other words, the more *time* the dataset covers and the higher its detail, the more likely a change might have occurred over the time.

All three issues relate to the *representativeness* of the data with regards to the unknown *true* process that generated the data. However, they also represent open challenges that require research on their own. For our purpose, we have to assume that the data is representative and its underlying process is static. These assumptions are widely applied in the body of process mining literature **CITE** .

2.2 Multivariate Time-Series Modelling

The data which is mined in Process Mining is typically a multivariate time-series. It is important to establish the characteristics of time-series.

2.2.1 What are Time Series Models?

A time series can be understood as a series of observable values, that depend on previous values. The causal dependence turns time-series into a special case of sequence models. Sequences do not *have to* depend on previous values. They might depend on previous and future values or not be interdependent at all. An example of a sequence model would be a language model. Results in Natural Language Processing (NLP), that the words in a sentences for many languages do not seem to only depend on prior words but also on future words **CITE** . Hence, we can assume that a human has formulated his sentence in the brain before expressing it in a sequence of words **CITE** . In contrast to sequences, time series cannot depend on future values. The general understanding of *time* is linear and forward directed **CITE** . The notion of time relates to our understanding of *cause and effect*. Hence, we can decompose any time series in a precedent (causal) and an antecedent (effect)

part (CHECK [13]). A time series model attempts to capture the relationship between precedent and antecedent.

2.2.2 The Challenges of Time Series Modelling

The analysis of unrestricted sequential opens up a myriad of challenges. First, sequential data introduce a combinatorial set of possible realisations. For instance, a set of two objects $\{A, B\}$ yields 7 theoretical combinations ($\{\emptyset\}$, $\{A\}$, $\{B\}$, $\{A, B\}$, $\{B, A\}$, $\{A, A\}$, $\{B, B\}$). Just by adding C and D to the object set increases the number of combinations to 40 and then 341. Second, sequential data may contain cycles which increases the number of possible productions to infinity CITE . Both, the combinatorial increase and cycles, contain a set of a countable infinite number of possibilities for discrete sets. However, as processes may also contain additional information a third obstacle arises. Including additional information increases the set to an uncountable number of possible values. With these obstacles in mind, it often becomes intractable to compute an exact model.

Hence, we have to include restrictive assumptions to reduce the solution space to a tractable number. A common way to counter this combinatorial explosion is the inclusion of the *Granger Causality* assumption. This idea postulates the predictive capability of a sequence given its preceding sequence. In other words, if we know that D can only follow after C, then 341 possible combinations reduce to 170. All of these possible 170 combinations are now temporally-related and hence, we speak of a *time-series*.

However, the prediction of sequences raises two new questions. First, if we know the precedence of a time-series, what is the antecedent? And second, if we can predict the antecedent accurately, what caused it? At first glance, it is easy to believe that both questions are quite similar, because we could assume that the precedent causes the antecedent. However, the first question is often solved using predictive AI models that rely on data, like Hidden-Markov-Models or Deep Learning. However, the latter question is much more difficult as data cannot help explain causes of sequences that never occurred. To illustrate this impossibility, if we never encountered 'D' in our dataset consisting of A, B and C, we cannot reliably say what would cause D. Answering this question requires additional tools within the XAI framework. One such method is the focus of this thesis and is further explored in section 2.3.

2.2.3 A Process within the State-Space Framework

Generally speaking, every time-series can be represented as a state-space model[10]. Within this framework the system consists of *input states* for *subsequent states* and *subsequent outputs*. A mathematical form of such a system is shown in Equation 2.1.

$$\begin{aligned} \mathbf{z}(t+1) &= \mathbf{h}(t, \mathbf{z}(t), \mathbf{u}(t)) \\ \mathbf{x}(t) &= \mathbf{f}(t, \mathbf{z}(t), \mathbf{u}(t)) \end{aligned} \quad (2.1)$$

Here, \mathbf{u} represents the input, \mathbf{x} the state, t the time. The function f maps t , $\mathbf{z}(t)$ and $\mathbf{u}(t)$ to the next state $\mathbf{z}(t+1)$. \mathbf{x} acts as an output computed by function f which takes the same input as h . **The variables \mathbf{z} , \mathbf{u} , t and \mathbf{y} are vectors with discrete or continuous features.** The distinction of $\mathbf{z}(t+1)$ and $\mathbf{x}(t)$ decouples *hidden*¹ states, from *observable* system outputs. Figure 2.2 shows a graphical representation of these equations.

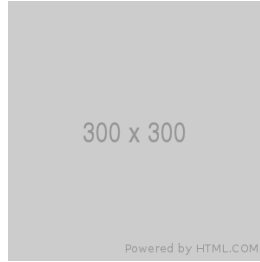


Figure 2.2: This figure shows a simplified graphical representation of a state-space model. Each arrow represents the flow of information.

The body of literature for state-space models is too vast to discuss it in detail². However, for process mining we can use this representation to discuss the necessary assumptions with regards to process mining. In accordance to the process-definition in section 2.1, we can understand the Event Log as a collection of observable outputs of a state-space model. The state of the process is hidden as the *true* process which generated the data cannot be observed as well. The time t is a step within the process. Hence, we will treat t as a discrete skalar value to denote discrete sequential time steps. The input \mathbf{u} represents all context information of the process. Here, \mathbf{u} subsumes observable information such as the starting point and Process Instance-related features. The functions f and h determine the transition of

¹A state does not have to be hidden. Especially, if we know the process and the transition rules. However, those are often inaccessible if we only use log data. Instead, many techniques try to approximate the hidden state given the data instead.

²For an introduction to state-space models see: XXX

a process' state to another state and its output over time. As we establish in section 2.1, we can assume that a process is a discrete sequence, whose transitions are time-variant. **[I might not need the linear formalisation as neural networks model nonlinear functions.]** *OPTIONAL:* These assumptions simplify the representation to a linear system of equations as depicted in Equation 2.2.

$$\begin{aligned} \mathbf{z}(k+1) &= \mathbf{A}(k)\mathbf{z}(k) + \mathbf{B}(k)\mathbf{u}(k) \\ \mathbf{y}(k) &= \mathbf{C}(k)\mathbf{z}(k) + \mathbf{D}(k)\mathbf{u}(k) \end{aligned} \tag{2.2}$$

This representation extends Equation 2.1 by including time dependent transition matrices A, B, C and D. We refer to this model as *explicit discrete time-variant* system. The transition matrices are time-dependant, as a process' transition may change at every time step.

[Note that k replaces t to express the *discreteness* of each time step. Also note the change of $z'(t)$ to $z(k+1)$ as it introduces a causal dependence on the previous state $z(k)$.]

[A number of AI techniques were developed to model this representation bla bla bla (HMM, Kalman, etc – Has further formalisation).] [ALSO MENTION additional assumptions or reductions.]

2.3 Counterfactuals

Counterfactuals are an important explanatory tool to understand a models' cause for decisions. Generating counterfactuals is main focus of this thesis. Hence, we will establish the most important characteristics of counterfactuals in this section.

2.3.1 What are Counterfactuals?

Counterfactuals have various definitions. However, their semantic meaning refers to “a conditional whose antecedent is false”[4]. A simpler definition from Starr states, counterfactual modality concerns itself with *what is not, but could or would have been*. Both definitions are related to linguistics and philosophy. Within AI and the mathematical framework various formal definitions can be found within causal inference[8]. Here, citeauthor describes a counterfactual as Causal inference definition. What binds all of these definitions is the notion of causality within “what if” scenarios.

However, for this paper, we will use the understanding established within the XAI context. Within XAI, counterfactuals act as a prediction which “describes the smallest change to the feature values that changes the prediction

to a predefined output” [16]. Note that XAI mainly concerns the explanation of models, which are always subject to inductive biases of the model itself and therefore inherently subjective. The idea behind counterfactuals as explanations³ is that we understand the output of a model, if we know what change caused would cause a different outcome. For instance, let’s denote a sequence 1 as *ABCDEF \mathbf{G}* , then a counterfactual *ABCDEX \mathbf{Z}* would tell us that **F** (probably) caused **G** in sequence 1. As counterfactuals only address explanations of single model instances and not the model as a whole, they are called *local* explanation.

Valid counterfactuals satisfy four criteria. First, a counterfactual should be minimally different from the true instance. If the counterfactual to sequence 1 was *AACDEX \mathbf{Z}* we would already have difficulties to discern whether B or F or both caused G at the end of sequence 1. Second, a counterfactual should produce a predefined outcome as closely as possible. This characteristic is ingrained in Molnars definition. If the counterfactual *ABCDEX \mathbf{Z}* ends with Z but this sequence is highly unrealistic, then cannot be certain of our conclusion for sequence 1. Third, we typically desire multiple diverse counterfactuals. One counterfactual might not be enough to understand the causal relationships in a sequence. In the example above we might have a clue that F causes G but what if G is not only caused by F? If we are able to find counterfactuals *VBCDEF \mathbf{H}* and *ABCDEX \mathbf{Z}* but all other configurations lead to G, then we know positions 1 and 6 cause G. As last criterion, each counterfactual should be possible. A sequence *ABCDE1 \mathbf{G}* would not be possible if numerals are not allowed. All four criteria allow us to assess the validity of each generated counterfactual and thus, help us to define an evaluation metric.

2.3.2 The Challenges of Counterfactual Sequence Generation

The current literature surrounding counterfactuals expose a number of challenges when dealing with counterfactuals.

The most important disadvantage of counterfactuals is the Rashomon Effect [16, ch. 9.3]. If all of the counterfactuals are valid, but contradict each other, we have to decide which of the *truths* are worth considering.

This decision reveals the next challenge of evaluation **CITE**. Although, the criteria can support us with the decision, it remains a question *how* to evaluate counterfactuals. Every automated measure comes with implicit

³There are other explanatory techniques in XAI like *feature importances* but counterfactuals are considered the most human-understandable

assumptions and often do not guarantee a realistic explanation. We still need domain experts to assess their validity.

The generation of counterfactual sequences contribute to both former challenges, due to the combinatorial expansion of the solution space. This problem is common for counterfactual sentence generation and has been addressed within the NLP [CITE](#). However, as process mining data not only consist of discrete objects like *words*, but also event and case features, the problem remains a daunting task. So far, little work has gone into the generation of multivariate counterfactual sequences like Process Instances [CITE](#).

2.4 Related Literature

2.4.1 The Generation of Counterfactual Explanations

The topic of counterfactual generation as explanation method was introduced by Wachter et al. in 2017 [CITE](#). The authors defined a loss function which incorporates the criteria to generate a counterfactual maximizes the likelihood for a predefined outcome and minimizes the distance to the original instance. However, the solution of Wachter et al. did not account for the minimalisation of feature changes and does not penalize unrealistic features. Furthermore, their solution cannot incorporate categorical variables.

A newer approach by Dandl et al. incorporates all four main criteria for counterfactuals (see section 2.3) by applying a genetic algorithm with a fitness function that incorporates all of the main criteria [CITE](#). This approach strongly differs from gradient-based methods, as it does not require a differentiable function which requires optimization. However, their solution only works with structured data.

2.4.2 Generating Sequential Counterfactual Explanations

When it comes to sequential data most researchers work on ways to generate counterfactuals for natural language. This often entails generating univariate discrete counterfactuals with the use of Deep Learning techniques. Martens and Provost and later Krause et al. are early examples of counterfactual NLP research. Their approach strongly focuses on the manipulation of sentences to achieve the desired outcome. However, as Robeer et al. puts it, their counterfactuals do not comply with *realisticness*.

Instead, Robeer et al. showed that it is possible to generate realistic counterfactuals with a Generative Adversarial Model (GAN). They use the model

to implicitly capture a latent state space and sample counterfactuals from it. Apart from implicitly modelling the latent space with GANs, it is possible to sample data from an explicit latent space. Examples of these approaches often use an encoder-decoder pattern in which the encoder encodes a data instance into a latent vector, which will be perturbed and then decoded into a similar instance[15][26]. By modelling the latent space, we can simply sample from a distribution conditioned on the original instance. Bond-Taylor et al. provides an overview of the strengths and weaknesses of common generative models.

Eventhough, a latent vector model can theoretically produce multivariate sequences, it a single latent-vector may be too restrictive to capture the combinatorial space of multivariate sequences. Hence, most of the models within NLP were not used to produce a sequence of vectors, but a sequence of discrete symbols. For process instances, we can assume a causal relation between state vectors in a sequential latent space. We call models that capture a sequential latent state-space which has causal relations *dynamic*[13]. Early models of this type of dynamic latent state-space models are the well-known *Kalman-Filter* for continous states and Hidden Markov Model (HMM) for discrete states. In recent literature, many techniques use Deep Learning to model complex state-spaces. The first models of this type were developed by Krishnan et al. Their Deep Kalman Filter (DKF) and subsequent Deep Markov Model (DMM) approximate the dynamic latent state-space by modeling the latent space given the data sequence and all previous latent vectors in the sequence. There are many variations of Krishnan et al.’s model, but most use Evidence Lower-Bound (ELBO) of the posterior for the current Z_t given all previous $\{Z_{t-1}, \dots, Z_1\}$ and X_t .

2.4.3 Counterfactual Reasoning for Business Processes

So far, none of the models have been applied to process data which can be treated as multivariate time-series.

Within Process Mining (PM), Causal Inference has long been used to analyse and model business processes. Mainly, due to the causal relationships underlying each process. However, early work has often attempted to incorporate domain-knowledge about the causality of processes in order to improve the process model itself[2, 9, 21, 27]. Among these, Narendra et al. approach is one of the first to include counterfactual reasoning for process optimization. Oberst and Sontag use counterfactuals to generate alternative solutions to treatments, which lead to a desired outcome. Again, the authors do not attempt to provide an explanation of the models outcome and therefore, disregard multiple **[criteria for viable counterfactuals]** in

XAI. [19] published the most recent paper on the counterfactual generation of explanations. **[Explain their approach.] The authors, use a known Structural Causal Model (SCM), to guide the generations of counterfactuals. However, this approach requires a process model which is as close as possible to the *true* process model. For our approach, we assume that no knowledge about the dependencies are known.**

Within the XAI context, Tsirtsis et al. develop the first explanation method. However, their work closely resembles the work of Oberst and Sonntag and treat the task as Markov Decision Process (MDP). This extension of a regular Markov Process (MP) assumes that an actor influences the outcome of a process given the state. This formalisation allows the use of Reinforcement Learning (RL) methods like Q-learning or SARSA **CITE**. However, this often requires additional assumptions such as a given reward function a discrete action-space. **However, within the context of this thesis, there is no obvious reward function nor discrete action-space available.**

2.4.4 The Construction of Counterfactual Time Series

Within the *multivariate time-series* literature two recent approaches yield ideas worth discussin.

First, Delaney et al. introduces a case-based reasoning to generate counterfactuals. Their method uses existing counterfactual instances, or *prototypes*, in the dataset. Therefore, it ensures, that the proposed counterfactuals are *realistic*. However, case-based approaches strongly depend on the *representativeness* of the prototypes **CITE**. In other words, if the model displays behaviour, which is not capture within the set of prototypical instances, most case-based techniques will fail to provide valid counterfactuals. The likelihood of such a break-down increases due to the combinatorial explosion of possible behaviours if the *true* process model has cycles or continuous event attributes. Cycles may cause infinite possible sequences and continous attributes can take values on a domain within infinite negative and positive bounds. However, despite these shortcomings, case-based approaches may act as a valuable baseline against other sophisticated approaches.

The second paper within the multivariate time series field by Ates et al. also uses a case-based approach. However, it contrasts from other approaches, as it does not specify a particular model but proposes a general framework instead. Hence, within this framework, individual components could be substituted by better performing components. Describing a framework, rather than specifying a particular model, allows to adapt the framework, due to the heterogeneous process dataset landscape. In this paper, we will also introduce a framework that allows for flexibility depending on the dataset. **The**

framework will be evaluated in two steps. The first step aims to compare various model types against each other based on the counterfactual validity. The second step scrutinizes the best framework configurations from step one, by presenting its results to a domain expert.

2.5 Formal Definitions

To formalise the log and a process model, we will use the formalisation established by citeauthor **CITE**.

2.6 Research Question

2.7 General Approach

2.8 What is Process Mining?

2.9 Challenges of Processing Process Data

Chapter 3

Related Papers

Chapter 4

Methods

4.1 Datasets

4.2 Preprocessing

4.3 Framework

Chapter 5

Results

5.1 Evaluation

Chapter 6

Discussion

Chapter 7

Conclusion

Bibliography

- Ates, E., Aksar, B., Leung, V. J., & Coskun, A. K. (2021, May 19). Counterfactual Explanations for Multivariate Time Series. In *2021 International Conference on Applied Artificial Intelligence (ICAPAI)* (pp. 1–8). 2021 International Conference on Applied Artificial Intelligence (ICAPAI). doi:10.1109/ICAPAI49758.2021.9462056
- Baker, J., Song, J., & Jones, D. R. (2017). Closing the Loop: An Empirical Investigation of Causality in IT Business Value. *undefined*. Retrieved March 1, 2022, from <https://www.semanticscholar.org/paper/Closing-the-Loop%3A-An-Empirical-Investigation-of-in-Baker-Song/df210060211bdc598f2d3382c68c615319287f71>
- Bond-Taylor, S., Leach, A., Long, Y., & Willcocks, C. G. (2021, April 14). Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models. arXiv: 2103.04922 [cs, stat]. Retrieved October 1, 2021, from <http://arxiv.org/abs/2103.04922>
- Counterfactual. (n.d.). doi:10.1093/oi/authority.20110803095642948
- Dandl, S., Molnar, C., Binder, M., & Bischl, B. (2020). Multi-Objective Counterfactual Explanations. In T. Bäck, M. Preuss, A. Deutz, H. Wang, C. Doerr, M. Emmerich, & H. Trautmann (Eds.), *Parallel Problem Solving from Nature – PPSN XVI* (pp. 448–469). doi:10.1007/978-3-030-58112-1_31
- Definition of PROCESS. (n.d.). Retrieved February 17, 2022, from <https://www.merriam-webster.com/dictionary/process>
- Delaney, E., Greene, D., & Keane, M. T. (2021). Instance-Based Counterfactual Explanations for Time Series Classification. In A. A. Sánchez-Ruiz & M. W. Floyd (Eds.), *Case-Based Reasoning Research and Development* (pp. 32–47). doi:10.1007/978-3-030-86957-1_3
- Hitchcock, C. (2020). Causal Models. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Summer 2020). Metaphysics Research Lab, Stanford University. Retrieved February 10, 2022, from <https://plato.stanford.edu/archives/sum2020/entries/causal-models/>

- Hompes, B. F. A., Maaradji, A., La Rosa, M., Dumas, M., Buijs, J. C. A. M., & van der Aalst, W. M. P. (2017). Discovering Causal Factors Explaining Business Process Performance Variation. In E. Dubois & K. Pohl (Eds.), *Advanced Information Systems Engineering* (pp. 177–192). doi:10.1007/978-3-319-59536-8_12
- Kalman, R. E. (1960). A New Approach to Linear Filtering and Prediction Problems. *Transactions of the ASME—Journal of Basic Engineering*, 82, 35–45.
- Krause, J., Perer, A., & Ng, K. (2016, May 7). Interacting with Predictions: Visual Inspection of Black-box Machine Learning Models. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems* (pp. 5686–5697). doi:10.1145/2858036.2858529
- Krishnan, R., Shalit, U., & Sontag, D. (2017). Structured Inference Networks for Nonlinear State Space Models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 31(1). Retrieved February 22, 2022, from <https://ojs.aaai.org/index.php/AAAI/article/view/10779>
- Leglaive, S., Alameda-Pineda, X., Girin, L., & Horaud, R. (2020, February 10). A Recurrent Variational Autoencoder for Speech Enhancement. arXiv: 1910.10942 [cs, eess]. Retrieved February 7, 2022, from <http://arxiv.org/abs/1910.10942>
- Martens, D., & Provost, F. (2014). Explaining data-driven document classifications. *MIS Quarterly*, 38(1), 73–100. doi:10.25300/MISQ/2014/38.1.04
- Melnyk, I., Santos, C. N. dos, Wadhawan, K., Padhi, I., & Kumar, A. (2017, December 4). Improved Neural Text Attribute Transfer with Non-parallel Data. arXiv: 1711.09395 [cs]. Retrieved February 28, 2022, from <http://arxiv.org/abs/1711.09395>
- Molnar, C. (2019). *Interpretable machine learning. A Guide for Making Black Box Models Explainable*. Retrieved from <https://christophm.github.io/interpretable-ml-book/>
- Narendra, T., Agarwal, P., Gupta, M., & Dechu, S. (2019). Counterfactual Reasoning for Process Optimization Using Structural Causal Models. In T. Hildebrandt, B. F. van Dongen, M. Röglinger, & J. Mendling (Eds.), *Business Process Management Forum* (pp. 91–106). doi:10.1007/978-3-030-26643-1_6
- Oberst, M., & Sontag, D. (2019, June 6). Counterfactual Off-Policy Evaluation with Gumbel-Max Structural Causal Models. arXiv: 1905.05824 [cs, stat]. Retrieved September 22, 2021, from <http://arxiv.org/abs/1905.05824>
- Qafari, M. S., & van der Aalst, W. M. P. (2021). Case Level Counterfactual Reasoning in Process Mining. In S. Nurcan & A. Korthaus (Eds.),

- Intelligent Information Systems* (pp. 55–63). doi:10.1007/978-3-030-79108-7_7
- Robeer, M., Bex, F., & Feelders, A. (2021, November). Generating Realistic Natural Language Counterfactuals. In *Findings of the Association for Computational Linguistics: EMNLP 2021* (pp. 3611–3625). EMNLP-Findings 2021. doi:10.18653/v1/2021.findings-emnlp.306
- Shook, C. L., Ketchen Jr., D. J., Hult, G. T. M., & Kacmar, K. M. (2004). An assessment of the use of structural equation modeling in strategic management research. *Strategic Management Journal*, 25(4), 397–404. doi:10.1002/smj.385
- Starr, W. (2021). Counterfactuals. In E. N. Zalta (Ed.), *The Stanford Encyclopedia of Philosophy* (Summer 2021). Metaphysics Research Lab, Stanford University. Retrieved February 9, 2022, from <https://plato.stanford.edu/archives/sum2021/entries/counterfactuals/>
- Tsirtsis, S., De, A., & Gomez-Rodriguez, M. (2021, July 6). Counterfactual Explanations in Sequential Decision Making Under Uncertainty. arXiv: 2107.02776 [cs, stat]. Retrieved September 9, 2021, from <http://arxiv.org/abs/2107.02776>
- van der Aalst, W., Adriansyah, A., de Medeiros, A. K. A., Arcieri, F., Baier, T., Blickle, T., . . . Wynn, M. (2012). Process Mining Manifesto. In F. Daniel, K. Barkaoui, & S. Dustdar (Eds.), *Business Process Management Workshops* (pp. 169–194). doi:10.1007/978-3-642-28108-2_19
- Wachter, S., Mittelstadt, B., & Russell, C. (2017). Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR. *ArXiv*. doi:10.2139/ssrn.3063289
- Wang, K., Hua, H., & Wan, X. (2019, December 12). Controllable Unsupervised Text Attribute Transfer via Editing Entangled Latent Representation. arXiv: 1905.12926 [cs]. Retrieved November 9, 2021, from <http://arxiv.org/abs/1905.12926>
- Wang, Z., Zhang, J., Xu, H., Chen, X., Zhang, Y., Zhao, W. X., & Wen, J.-R. (2021, July 11). Counterfactual Data-Augmented Sequential Recommendation. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (pp. 347–356). doi:10.1145/3404835.3462855