



**Utrecht
University**

Department of Mathematics and Computer Science
Process Analytics

The Generation of viable Counterfactual Examples by finding minimal Edit Sequences using Event Data in Complex Processes

Master Thesis

Olusanmi A. Hundogan

Supervisors:

dr. ir. Xixi Lu
Yupei Du M. Sc.
August 10, 2022

Abstract

Contents

1	Introduction	4
1.1	Motivation	4
1.2	Problem Space	6
1.3	Related Literature	8
1.3.1	Generating Counterfactuals	8
1.3.2	Generating Counterfactual Sequences	8
1.3.3	Generating Counterfactual Time-Series	9
1.3.4	Generating Counterfactuals for Business Processes	10
1.4	Research Question	11
1.5	Outline	13
2	Background	15
2.1	Process Mining	15
2.1.1	A definition for Business Processes	15
2.1.2	What is Process Mining?	18
2.1.3	The Challenges of Process Mining	18
2.2	Multivariate Time-Series Modeling	19
2.2.1	What are Time Series Models?	20
2.2.2	The Challenges of Time Series Modelling	20
2.3	Counterfactuals	21
2.3.1	What are Counterfactuals?	21
2.3.2	The Challenges of Counterfactual Sequence Generation	23
2.4	Formal Definitions	24
2.4.1	Process Logs, Cases and Instance Sequences	24
2.4.2	State-Space Models	26
2.5	Representation	29
2.6	Long-Short-Term Memory Models	30
2.7	Damerau-Levenshtein	31
2.8	Evolutionary Algorithms	34

3	Methods	39
3.1	Methodological Framework	39
3.1.1	Architecture	39
3.1.2	Differences to DiCE4EL	40
3.2	Semi-Structured Damerau-Levenshtein	
	Distance	41
3.2.1	Semi-Structured Damerau Levenshtein	42
3.2.2	Discussion	43
3.3	Viability Measure	44
3.3.1	Similarity-Measure	45
3.3.2	Sparcity-Measure	45
3.3.3	Feasibility-Measure	46
3.3.4	Delta-Measure	47
3.3.5	Discussion	48
3.3.6	Differences to DiCE4EL	49
3.4	Prediction Model: LSTM	50
3.5	Counterfactual Generators	52
3.5.1	Baseline Model: Random Generator	52
3.5.2	Baseline Model: Sample-Based Generator	52
3.5.3	Baseline Model: Case-Based Generator	52
3.5.4	Generative Model: Evolutionary Algorithm	53
4	Evaluation	56
4.1	Datasets	56
4.2	Preprocessing	57
4.3	Experimental Setup	58
5	Results	61
5.1	Experiment 1: Model Selection	61
5.1.1	Model Configuration	61
5.1.2	Model Termination Point	64
5.1.3	Model Parameters	66
5.1.4	Model Candidates	68
5.2	Experiment 2: Model Comparison	69
5.2.1	Results	69
5.2.2	Analysis	71
5.3	Experiment 3: Evaluation under a different Viability Measure	72
5.3.1	Results	72
5.3.2	Analysis	74
5.4	Experiment 4: Qualitative Assessment	74
5.4.1	Results	74

5.4.2	Analysis	75
6	Discussion	77
6.1	Interpretation of Results	77
6.2	Limitations	78
6.3	Improvements	79
6.4	Future Work	80
7	Conclusion	81
	Appendices	82
A	Counterfactual Results	83

Chapter 1

Introduction

1.1 Motivation

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[51]. 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. A field that focuses on modelling processes is Process Mining (PM).

Research in the Process Mining discipline has shown that it is possible to predict the outcome of a particular process fairly well[28, 48]. For instance, in the medical domain, models have been shown to predict the outcome or trajectory of a patient's condition[34]. In the private sector, process models can be used to detect faults or outliers. The research discipline Deep Learning has shown promising results within domains that have been considered difficult for decades. The Moravex Paradox[1], which postulates that machines are capable of doing complex computations easily while failing in tasks that seem easy to humans such as object detection or language comprehension, does not hold anymore. Meaning that with enough data to learn, machines are capable of learning highly sophisticated tasks, better than any human. The same holds for predictive tasks. However, while many prediction models can predict certain outcomes, it remains a difficult challenge to understand their reasoning.

This difficulty arises from models, like neural networks, that are so-called *blackbox models*. Meaning, that their inference is incomprehensible, due to the vast amount of parameters involved. This lack of comprehension is undesirable for many fields like IT or finance. Not knowing why a loan was

given, makes it impossible to rule out possible biases. Knowing what will lead to a system failure, will help us knowing how to avoid it. In critical domains like medicine, the reasoning behind decisions become crucial. For instance, if we know that a treatment process of a patient reduces the chances for survival, we want to know which treatment step is the critical factor we ought to avoid. To summarise, knowing the outcome of a process often leads us to questions on how to change it. Formally, we want to change the outcome of a process instance, by making it maximally likely, with as little interventions as possible[39]. Figure 1.1 is a visual representation of the desired goal.



Figure 1.1: This figure illustrates a model, that predicts a certain trajectory of the process. However, we want to change the process steps in such a way, that it changes the outcome.

One-way to better understand the Machine Learning (ML) models lies within the eXplainable AI (XAI) discipline. XAI focuses the developments of theories, methods, and techniques that help explaining blackbox models to humans. Most of the discipline’s techniques produce explanations that guide our understanding. Explanations can come in various forms, such as IF-THEN rules[39, p.90] or feature importances[39, p.45]. but some are more comprehensible for humans than others.

A prominent and human-friendly approach are *counterfactuals*[39, p. 221]. Counterfactuals within the AI framework help us to answer hypothetical ”what-if” questions. Basically, if we know *what* would happen *if* we changed the execution of a process instance, we could change it for the better. In this thesis, we raise the question, how we can use counterfactuals to change the trajectory of a process models’ prediction towards a desired outcome. Knowing the answers not only increases the understanding of blackbox models, but also help us avoid or enforce certain outcomes.

1.2 Problem Space

In this thesis, we 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.

For process data, the literature often uses terms like structured and semi-structured, as they are related to the staticity and dynamicity. Both, structuredness and semi-structuredness, often relate to the data model, in which we structure the information at hand. As static data neither changes over time nor changes its structure, we can use structured data-formats such as tables to capture the information and each data point is an independent entity. We can take the MNIST dataset[16] or Iris dataset[3, 17] as examples for structured and static data. In both datasets, all data points are independent and have the same amount of attributes. In contrast, semi-structured data does not have to follow these strict characteristics. Here, data points often belong to a group of data points which constitutes the full entity. Furthermore, the attributes of each data point may vary. The grouping mechanism could take the form of associative links, class associations or temporal cause-effect relationships. Examples of these are Part-of-Speech datasets like Penn Treebank set[35]. Here, we often associate each data point with a sentence. However, the temporal relationship between words is debatable and hence whether the data is *dynamic* as well. Hence, not all semi-structured data sets are dynamic and vice versa. However, structured data will almost always be static, with the exception of time-series. Lastly, there is also unstructured data, which does not incorporate any specific data model. Corpora like the Brown dataset[19], for instance, are collections of text heavy unstructured information. In Figure 1.2, we show various examples of data.

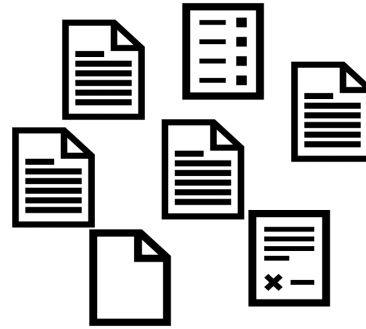
A major reason, why there has not been much research on counterfactuals for dynamic semi-structured data, emerges from a multitude of challenges, when dealing with counterfactuals and sequences. Three of these challenges are particularly important.

First, counterfactuals within AI attempt to explain outcomes which never occurred. A *what-if* questions often refer to hypothetical scenarios. Therefore, there is no evidential data, from which we 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.

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

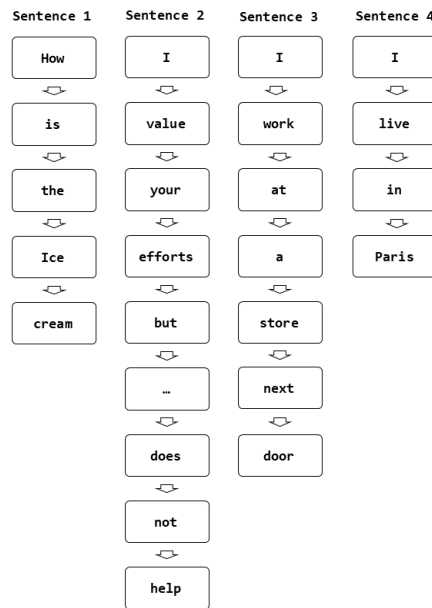
²With dynamic data, we refer to data that has a temporal relationship as a major component, which is also inherently sequential

s.length	s.width	p.length	p.width	variety
6.5	2.8	4.6	1.5	Versicolor
5.8	2.7	4.1	1.0	Versicolor
6.7	3.3	5.7	2.5	Virginica
4.6	3.4	1.4	0.3	Setosa
6.4	3.2	5.3	2.3	Virginica
5.9	3.0	4.2	1.5	Versicolor
7.4	2.8	6.1	1.9	Virginica
5.5	2.4	3.8	1.1	Versicolor
5.6	2.5	3.9	1.1	Versicolor
5.0	3.4	1.5	0.2	Setosa
6.9	3.1	5.4	2.1	Virginica
5.5	2.5	4.0	1.3	Versicolor
5.7	2.6	3.5	1.0	Versicolor
5.8	2.7	3.9	1.2	Versicolor
7.6	3.0	6.6	2.1	Virginica
6.7	3.3	5.7	2.1	Virginica
5.0	3.5	1.6	0.6	Setosa
7.7	2.8	6.7	2.0	Virginica
6.4	2.7	5.3	1.9	Virginica
7.7	3.8	6.7	2.2	Virginica
5.2	3.5	1.5	0.2	Setosa
5.7	3.8	1.7	0.3	Setosa



(a) An excerpt of the MNIST dataset. This is a structured dataset.

(b) A number of heterogenous documents. A dataset like this is unstructured.



(c) Multiple sequences of words. Each word forms a sentence of different lengths. Therefore, this data is semi-structured.

Figure 1.2: Schematic examples of static structured, dynamic semi-structured data and unstructured data.

Second, sequential data is highly variable in length, but process steps have complicated factors, too. The sequential nature of the data impedes the tractability of many problems due to the combinatorial explosion of possible sequences. 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.

Third, process data often requires knowledge of the causal structures that produced the data in the first place. However, these structures are often hidden and it is a NP-hard problem to elicit them[54].

These challenges make the field, in which we can contribute a vast endeavor.

1.3 Related Literature

Many researchers have worked on counterfactuals and PM. Here, we combine the important concepts and discuss the various contributions to this thesis.

1.3.1 Generating Counterfactuals

The topic of counterfactual generation as explanation method was introduced by Wachter, Mittelstadt, and Russell in 2017[53]. The authors defined a loss function which incorporates the criteria to generate a counterfactual which maximizes the likelihood for a predefined outcome and minimizes the distance to the original instance. However, the solution of Wachter, Mittelstadt, and Russell 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, Molnar, Binder, and Bischl incorporates four main criteria for counterfactuals (see section 2.3) by applying a genetic algorithm with a multi-objective fitness function[13]. This approach strongly differs from gradient-based methods, as it does not require a differentiable objective function. However, their solution was only tested on static data.

1.3.2 Generating Counterfactual Sequences

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, Perer, and Ng are early examples of counterfactual NLP research[29, 36]. Their approach strongly focuses on

the manipulation of sentences to achieve the desired outcome. However, as Robeer, Bex, and Feelders puts it, their counterfactuals do not comply with *realisticness*[45].

Instead, Robeer, Bex, and Feelders showed that it is possible to generate realistic counterfactuals with a Generative Adversarial Model (GAN)[45]. 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[37, 55]. By modelling the latent space, we can simply sample from a distribution conditioned on the original instance. Bond-Taylor, Leach, Long, and Willcocks provides an overview of the strengths and weaknesses of common generative models.

Eventhough, a single latent vector model can theoretically produce multivariate sequences, it may still be too restrictive to capture the combinatorial space of multivariate sequences. Hence, most of the models within Natural Language Processing (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*[32]. Early models of this type of dynamic latent state-space models are the well-known *Kalman-Filter* for continuous 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, Shalit, and Sontag[29, 30]. Their Deep Kalman Filter (DKF) and subsequent Deep Markov Model (DMM) approximate the dynamic latent state-space by modelling the latent space given the data sequence and all previous latent vectors in the sequence. There are many variations[10, 18, 32] of Krishnan, Shalit, and Sontag’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 [21].

1.3.3 Generating Counterfactual Time-Series

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

First, Delaney, Greene, and Keane introduces a case-based reasoning to generate counterfactuals[15]. 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[39, p. 192]. In other words, if the model displays behaviour, which is not captured within the set of prototypical instances, most case-based techniques will fail to provide viable 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 continuous attributes can take values on a domain within infinite negative and positive bounds. These issues have not been explored in the paper of Delaney, Greene, and Keane, as it mainly deals with time series classification[15]. 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, Aksar, Leung, and Coskun also uses a case-based approach[5]. 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 also introduce a framework that allows for flexibility depending on the dataset.

1.3.4 Generating Counterfactuals for Business Processes

So far, none of the techniques have been applied to process data.

Within 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[6, 25, 46, 56]. Among these, Narendra, Agarwal, Gupta, and Dechu approach is one of the first to include counterfactual reasoning for process optimization[41]. Oberst and Sontag use counterfactuals to generate alternative solutions to treatments, which lead to a desired outcome[42]. Again, the authors do not attempt to provide an explanation of the models outcome and therefore, disregard multiple viability criterions for counterfactuals in XAI. Qafari and van der Aalst published the most recent paper on the counterfactual generation of explanations[44]. The authors, use a known Structural Causal Model (SCM), to guide the generation of their 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, De, and Gomez-Rodriguez develop the

first explanation method for process data[50]. However, their work closely resembles the work of Oberst and Sontag and treat the task as Markov Decision Process (MDP)[42]. 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. However, this often requires additional assumptions such as a given reward function and an action-space. For counterfactual sequence generation, there is no obvious choice for the reward function or the action-space.

Nonetheless, both Tsirtsis, De, and Gomez-Rodriguez and Oberst and Sontag contribute an important idea. The idea of incrementally generating the counterfactual instead of the full sequence. Hsieh, Moreira, and Ouyang has recently published an approach that builds on the same notion of incremental generation. Their approach has a very similar structure to our approach and appears to be the only one that we can compare our counterfactuals against.

For this reason, this thesis highlights some key differences and similarities. However, to understand the differences and similarities, we first have to establish some core concepts. In this section, we only discuss their approach, briefly.

The authors recognised that some processes have critical events, which govern the overall outcome. Hence, by simply avoiding the undesired outcome from critical event to critical event, it is possible to limit the search space and compute viable counterfactuals. They use an extension of DiCE[40] to generate counterfactuals. However, their approach requires concrete knowledge about these critical points. We propose a Framework that avoids this constraint.

To our knowledge, the authors are also the first authors that try to optimize their counterfactual process generation based on criteria that ensure their viability. However, in our approach, we use different operationalisations to quantify the criteria.

1.4 Research Question

As we seek to make data-driven process models interpretable, we have to understand the exact purpose of this thesis. Hence, we establish the open challenges and how this thesis attempts to solve them.

Having discussed the previous work on counterfactual sequence generation, a couple of challenges emerge. First, we need to generate on a set of criteria and therefore, require complex loss and evaluation metrics, that may

or may not be differentiable. Second, they cannot be logically impossible, given the data set. Hence, we have to restrict the space to counterfactuals of viable solutions, while being flexible enough to not just copy existing data instances. Third, using domain knowledge of the process significantly reduces the practicality of any solution. Therefore, we have to develop an approach, which requires only the given log as input while not relying on process specific domain knowledge. This begs the question, whether there is a method to generate sequential counterfactuals that are viable, without relying on process specific domain knowledge. In terms of specific research questions we try to answer:

RQ: How can an algorithm use existing counterfactual approaches for the generation of counterfactual sequences while incorporating structural differences between the factual sequence and the counterfactual sequence?

RQ1: How can we employ existing methods to compute viability, so that its optimization incorporates information about the structure of the sequence?

RQ2: To what extent can we generate counterfactuals that fulfill the criteria to be viable?

RQ3: How does an algorithm, which optimizes multiple viability quality metrics to perform against other approaches?

We approach these questions, by proposing a schematic framework which allows the exploration of several independent components. Figure 1.3 shows the conceptual framework of the base approach visually.

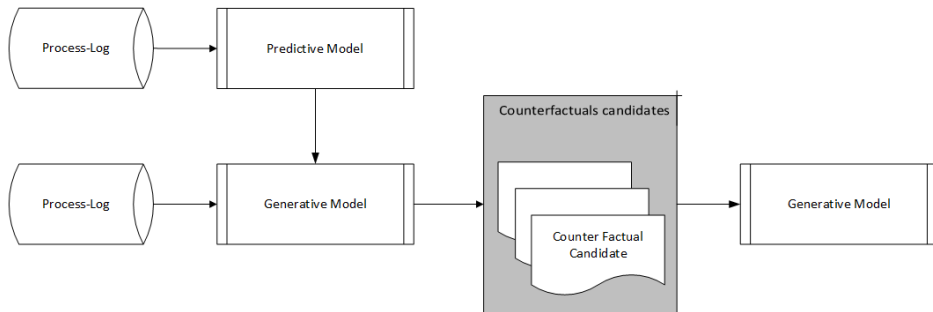


Figure 1.3: A simplified schematic representation of the framework which is explored in this thesis.

The framework contains three parts. First, we need a pretrained predictive component, which we aspire to explain. The component should *accurately*

predict the outcome of a process at any step. The accuracy-condition is favorable, but not necessary. If the component is accurately modelling the real world, we can draw real-world conclusions from the explanations generated. If the component is inaccurate, the counterfactuals only explain the prediction decisions and not the real world. The second part requires a generative component. The generative component needs to generate viable sequential counterfactuals which are logically *plausible*. A plausible counterfactual is one whose outcome can be predicted by the predictive component. If the predictive component cannot predict the counterfactual sequence, we can assume that the generative model is *unfaithful* to the predictive component that we want to explain. The third component is the evaluation metric upon which we decide the viability of the counterfactual candidates.

For the evaluation, we have to show the following:

- RQ1-H1: If we use a viability function which incorporates multiple criteria to determine counterfactuals, we consistently retrieve more viable counterfactuals, than choosing the at random.
- RQ1-H2: The generated counterfactuals consistently outperform the most viable counterfactuals among examples in the dataset.
- RQ2-H1: The results of the counterfactual are comparable to other existing literature.

1.5 Outline

The remainder of the thesis is outlined as follows: In chapter 2, we introduce all of the important concepts that are crucial to this thesis. Most importantly, we introduce the main research discipline PM and the subject of our research: *Counterfactuals*. Furthermore we cover some necessary background required to understand the methods, we employ.

The chapter 3, introduces our methodological framework in further detail. The chapter explains all the important components and methods, we apply, to answer the research question. Among these methods, we introduce the the methodological architecture, a modified version of the Damerau-Levenshtein distance.

chapter 4 covers the main approach behind our experimental setup. We discuss how we attempt to answer our research questions and introduce the datasets we are using and how we conduct the preprocessing.

In chapter 5 we report on the results and insights we gain from executing our research approach.

All the results are summarised in [**ch:discussion**]. Here, we summarize and interpret our results. We discuss limitations and possible improvements. We also discuss implications for future research endeavors.

The chapter 7 summarizes the thesis and the implications for the PM research field.

Chapter 2

Background

This chapter explores the most important concepts for this work. Hence, we focus on the problem domain, starting with an overview about PM. Afterwards, we discuss the nature of the data, we handle in this thesis by discussing *Multivariate Discrete Time-Series*. Next, we introduce counterfactuals and establish how we characterise *viable* counterfactuals.

2.1 Process Mining

This thesis focuses on processes and the modelling of process generated data. Hence, it is important to establish a common understanding for this field.

2.1.1 A definition for Business Processes

Before elaborating on Process Mining, we have to establish the meaning of the term *process*. The term is widely-used and therefore has a rich semantic volume. A process generally refers to something that advances and changes over time[14]. Despite, legal or biological processes being valid interpretations, too, we focus on *business processes*.

An example is a loan application process in which an applicant may request a loan. The case would then be assessed and reviewed by multiple examiners and end in a final decision. The loan might end up in an approval or denial. The *business* part is misleading as these processes are not confined to commercial settings alone. For instance, a medical business process may cover a patients 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 (HCI) 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 receival 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 work, we mainly follow the understanding outlined in van der Aalst, Adriansyah, de Medeiros, Arcieri, Baier, Blickle, Bose, van den Brand, Brandtjen, Buijs, Burattin, Carmona, Castellanos, Claes, Cook, Costantini, Curbera, Damiani, de Leoni, Delias, van Dongen, Dumas, Dustdar, Fahland, Ferreira, Gaaloul, van Geffen, Goel, Günther, Guzzo, Harmon, ter Hofstede, Hoogland, Ingvaldsen, Kato, Kuhn, Kumar, La Rosa, Maggi, Malerba, Mans, Manuel, McCreesh, Mello, Mendling, Montali, Motahari-Nezhad, zur Muehlen, Munoz-Gama, Pontieri, Ribeiro, Rozinat, Seguel Pérez, Seguel Pérez, Sepúlveda, Sinur, Soffer, Song, Sperduti, Stilo, Stoel, Swenson, Talamo, Tan, Turner, Vanthienen, Varvaressos, Verbeek, Verdonk, Vigo, Wang, Weber, Weidlich, Weijters, Wen, Westergaard, and Wynn[51]. Each step, including start- and end-points, is a process event which was caused by an *activity*. Often, both terms, *event* and *activity*, are used interchangeably. However, there are subtle differences. We understand an event as something that happens at a specific point in time. The driving question is *when* the event happens. In contrast, an activity is related to the content of an event. Here, we ask *what* happens at a point in time. For instance, if we apply for a loan that requires an approval by one person and afterwards a second approval, we can call both activities **APPROVAL**. Although both activities are fundamentally the *same*, they happen at different points in time. Henceforth, both events remain *different*. Mainly, because one can argue that both events have varying time dependent contexts. For instance, an approval at daytime might be caused by different reasons, than an event caused at night-time.

Each process event may contain additional information in the form of event attributes. If a collection of events *sequentially* relate to a single run through a process, we call them *process instance* or *trace*. These instances do not have to be completed. Meaning, the trace might end prematurely. 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 location, a patients age or the buyers budget. In its entirety, a business process can be summarised as a *graph*, a *flowchart* or another kind of visual representation. Figure 2.1's graphical representation is an example of such a *process map*[51].

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

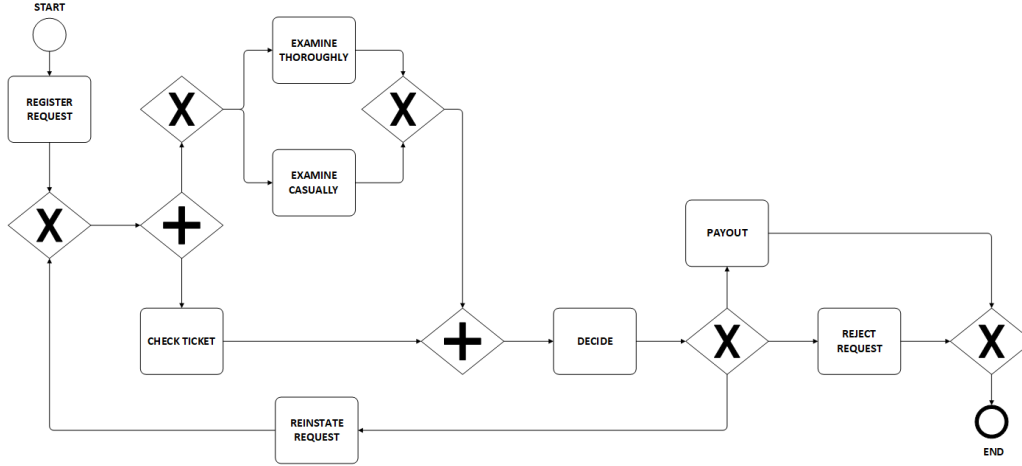


Figure 2.1: This graph shows an example of a Business Process Modell Notation (BPMN) process map.

A finite series of discrete events with one or more starting points, intermediary steps and end points. Each intermediate step has at least one precedent and at least one antecedent step.

However, we have to address a number of issues with this definition.

First, it excludes infinite processes like solar system movements or continuous processes such as weather changes. 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 within an 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 starting from the customers first visit until the successful shipment, an Amazon vendor might only be interested in the delivery process of a product.

Third, the example provided in Figure 2.1 may not relate to the *real* underlying data generating process. As process *models* are inherently simplified, they may or may not be accurate. The *true* process is often unknown. Therefore, we distinguish between the *true process* and a *modelled process*. The *true process* is a hypothetical concept whose *true* structure remains unknown. In, contrast, a process *model* simplifies and approximates the characteristics of the *true process*.

2.1.2 What is Process Mining?

Having established a definition for a process, we next discuss *Process Mining*. This young discipline has many connections to other fields that focus on the modelling and analysis of processes such as Continuous Process Improvement (CPI) or Business Process Management (BPM)[51]. However, its data-centric approaches originate in Data Mining. The authors van der Aalst, Adriansyah, de Medeiros, Arcieri, Baier, Blickle, Bose, van den Brand, Brandtjen, Buijs, Burattin, Carmona, Castellanos, Claes, Cook, Costantini, Curbera, Damiani, de Leoni, Delias, van Dongen, Dumas, Dustdar, Fahland, Ferreira, Gaaloul, van Geffen, Goel, Günther, Guzzo, Harmon, ter Hofstede, Hoogland, Ingvaldsen, Kato, Kuhn, Kumar, La Rosa, Maggi, Malerba, Mans, Manuel, McCreesh, Mello, Mendling, Montali, Motahari-Nezhad, zur Muehlen, Munoz-Gama, Pontieri, Ribeiro, Rozinat, Seguel Pérez, Seguel Pérez, Sepúlveda, Sinur, Soffer, Song, Sperduti, Stilo, Stoel, Swenson, Talamo, Tan, Turner, Vanthienen, Varvaressos, Verbeek, Verdonk, Vigo, Wang, Weber, Weidlich, Weijters, Wen, Westergaard, and Wynn 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”[51]. The discipline revolves around the analysis of event logs. A event log is a collection of process instances, which are retrieved from various sources like an Information System (IS) or database. Logs are often stored in data formats such as Comma Separated Values (CSV) or eXtensible Event Stream (XES)[51].

2.1.3 The Challenges of Process Mining

As mentioned in chapter 1, process data modelling and analysis is a challenging task. van der Aalst, Adriansyah, de Medeiros, Arcieri, Baier, Blickle, Bose, van den Brand, Brandtjen, Buijs, Burattin, Carmona, Castellanos, Claes, Cook, Costantini, Curbera, Damiani, de Leoni, Delias, van Dongen, Dumas, Dustdar, Fahland, Ferreira, Gaaloul, van Geffen, Goel, Günther, Guzzo, Harmon, ter Hofstede, Hoogland, Ingvaldsen, Kato, Kuhn, Kumar, La Rosa, Maggi, Malerba, Mans, Manuel, McCreesh, Mello, Mendling, Montali, Motahari-Nezhad, zur Muehlen, Munoz-Gama, Pontieri, Ribeiro, Rozinat, Seguel Pérez, Seguel Pérez, Sepúlveda, Sinur, Soffer, Song, Sperduti, Stilo, Stoel, Swenson, Talamo, Tan, Turner, Vanthienen, Varvaressos, Verbeek, Verdonk, Vigo, Wang, Weber, Weidlich, Weijters, Wen, Westergaard, and Wynn mentions a number of issues that arise from processes[51].

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 for information mining purposes. The information is often incomplete, due to a lack of context information, the omission of logged process steps, or wrong levels of granularity[51].

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. 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[51].

A third issue which contributes to the datasets' incompleteness and complexity is a phenomenon called *concept drift*[51]. This phenomenon relates to the possibility of changes 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[51]. 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[28, 48].

2.2 Multivariate Time-Series Modeling

The temporal and multivariate nature of process instance often turns PM into a Multivariate Time-Series Modeling problem. Therefore, it is necessary to establish an understanding for this type of data structure.

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 and 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 NLP, that the words in a sentences for many languages do not seem to only depend on prior words but also on future words[20]. Hence, we can assume that a human has formulated his sentence in the brain before expressing it in a sequence of words. In contrast to sequences, time series cannot depend on future values. The general understanding of *time* is causal and forward directed. 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[32]. 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 introduces a combinatorial set of possible realisations (often called *productions*). For instance, a set of two objects $\{A, B\}$ produces 7 theoretical combinations ($\{\emptyset\}$, $\{A\}$, $\{B\}$, $\{A, B\}$, $\{B, A\}$, $\{A, A\}$, $\{B, B\}$). Just by adding C and then D to the object set increases the number of combinations to 40 and 341, respectively. Second, sequential data may contain cyclical patterns which increase the number of possible productions to infinity[54]. Both, the combinatorial increase and cycles, yield a set of a countable infinite number of possible productions. However, as processes may also contain additional information a third obstacle arises. Including additional information extends the set to an uncountable number of possible productions. 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[2]. This idea postulates the predictive capability of a sequence given its preceding sequence. In other words, if we know that C must be followed by D, then 341 the number of possible combinations reduces to 156. All of these possible 156 combinations are now temporally-related and hence, we speak of a *time-series*.

However, the prediction of sequences recontextualises the issue to 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? We often use data-driven AI-methods like Hidden-Markov-Models or Deep Learning to solve the first question. However, the second question is more subtle. At first glance, it is easy to believe that both questions are quite similar, because we could assume that the precedent causes the antecedent. Meaning, that we can use the data available to elicit sequential correlative patterns. In reality, the latter question is much more difficult as data often does not include any information about the inter-relationships. To illustrate this difficulty, we could say that the presence of C causes D. But if D also appears to be valid in a sequence 'AABD', it cannot be caused by the presence of C alone.

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.

2.3 Counterfactuals

Counterfactuals are an important explanatory tool to understand a models' cause for decisions. Generating counterfactuals is main focus of this thesis. Hence, we 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*”[11]. A simpler definition from Starr states that 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 in the causal inference[23] literature. A prominent figure within the causal inference discipline is Pearl, Glymour, and Jewell, who postulates that a “*kind of statement – an 'if' statement in which the 'if' portion is untrue or unrealized – is known as a counterfactual*”[43]. What binds all of these definitions is the notion of causality within *what-if* scenarios.

For this paper, we 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 pre-defined output*” according to Molnar[39, p. 212]. Note that XAI mainly

concerns itself with the explanation of *models*, which are always subject to inductive biases and therefore, inherently subjective. The idea behind counterfactuals as explanatory tool¹ is simple. We understand the outcome of a model, if we know *what* outcome would occur *if* we changed its input.

Let us assume, a student is approaching an important deadline, which she desires to meet. Every day, she has a multitude of options to choose from. Either, continue with the report (option A), focus on learning more about the topic (option B), pursue her hobby as a break (option C), meet up with friends (option D), or procrastinate (option E). Furthermore, we assume, there are 7 days left and she can either miss (0) the deadline or meet it (1). The approach she follows is *ABABDEA* and she misses the deadline. Let us refer to this sequence of actions as the factual *sequence 1*. Then, a counterfactual *ABABDBA* that meets the deadline tells us that **E** (probably) caused missing the deadline. In other words, if the student had not procrastinated two days before the deadline she could have made it on time.

As counterfactuals only address explanations of one model result and not the model as a whole, they are *local* explanations[39, p. 212]. According to Molnar *Valid* counterfactuals satisfy **four** criteria[39, p. 212]:

- Similarity: A counterfactual should be similar to the original instance. If a successful counterfactual to sequence 1 was *ABABEEA*, we would already have difficulties to discern whether meeting with friends *D*, procrastinating *E* or both caused the outcome of missing the deadline 0. Hence, we want to be able to easily compare the counterfactual with the original. We can achieve this by minimizing their mutual distance.
- Sparcity: In line with the notion of similarity, we want to change the original instance only minimally. If the sequence had many changes, it would similarly impede the understanding of causal relationships in sequence 1.
- Feasibility: Each counterfactual should be feasible. In other words, impossible values are not allowed. As an example, if the student followed a strict *AAAAAAEA* would not be feasible if we consider students could burn-out. Typically, we can use data to ensure this property. However, the *open-world assumption* impedes this solution. With *open-world*, we mean that processes may change and introduce behaviour that has not been measured before. A student might only attempt a Bachelor's

¹There are other explanatory techniques in XAI like *feature importances* but counterfactuals are considered the most human-understandable

thesis once. Especially, for long and cyclical sequences, we have to expect previously unseen sequences.

Likelihood: A counterfactual should produce the desired outcome if possible. This characteristic is ingrained in Molnar’s definition. However, as the model might not be persuaded to change its prediction, we relax this condition. We say that we want to increase the likelihood of the outcome as much as possible. If the counterfactual *ABABDXA* hinges on *X* as in an earthquake occurring that postpones the deadline, the sequence would be highly unrealistic. Hence, we cannot be certain of our conclusion for sequence 1. Therefore, we want the counterfactual’s likelihood to be at least more likely than the factual outcome.

All four criteria allow us to assess the viability of each generated counterfactual and thus, help us to define an evaluation metric for each individual counterfactual. However, we also seek to optimise certain qualities on the population level of the counterfactual candidates.

Diversity: 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 *E* causes an outcome of 0, but what if outcome 0 is by more than *E*? If we are able to find counterfactuals all counterfactuals that involve *E* and that lead to missing the deadline, we get a better understanding of what caused outcome 0.

Realism: For a real world application, we still have to evaluate their *reasonability* within the applied domain. This is a characteristic that can only be evaluated by a domain expert.

We refer to both sets of viability criterions as *individual viability* and *population viability*. However, to remain concise, we use *viability* to refer to the individual criterions only. We explicitly mention *population viability* if we refer to criterions that concern the population.

2.3.2 The Challenges of Counterfactual Sequence Generation

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

The most important disadvantage of counterfactuals is the Rashomon Effect[39, ch.9.3]. If all of the counterfactuals are viable, but contradict each other, we have to decide which of the *truths* are worth considering.

This decision reveals the next challenge of evaluation. Although, the criteria can support us with the decision, it remains an open research question *how* to evaluate counterfactuals according to Carvalho, Pereira, and Cardoso. So far, no one was able to establish a standardized evaluation protocol[26]. Every automated measure comes with implicit assumptions and they cannot guarantee a realistic explanations. Furthermore, we attempt to explain something with – in simple terms – *experiences* that never actually occurred. We still need domain experts to assess their *plausibility*.

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. 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.

2.4 Formal Definitions

Before diving into the rest of this thesis, we have to establish preliminary definitions, we use in this work. With this definitions, we share a common formal understanding of mathematical descriptions of every concept used within this thesis.

2.4.1 Process Logs, Cases and Instance Sequences

We start by formalising the event log and its elements. We use a medical process as an example to provide a better semantic understanding. An event log is denoted as L . Here, L could be as database which logs the medical histories of all patients in a hospital.

We assume the database logs all interactions, be it therapeutic or diagnostic and store them as an event with a unique identifier. Let \mathcal{E} be the universe of these event identifiers and $E \subseteq \mathcal{E}$ a set of events. The set E could consist, for instance, of a patients first session with a medical professional, then a diagnostic scan, followed by therapie sessions, surgery and more.

All of these interactions with one patient make up a case, which has a unique identifier, too. Let C be a set of case identifiers and $\pi_\sigma : E \mapsto C$ a surjective function that links every element in E to a case $c \in C$ in which c signifies a specific case. The function allows us to associate every event within the database to a single patient. The function’s surjective property ensures for each case there exists at least one event.

For a set of events $E \subseteq \mathcal{E}$, we use a shorthand s^c being a particular sequence $s^c = \langle e_1, e_2, \dots, e_t \rangle$ with c as case identifier and a length of t . Each s is a trace of the process log $s \in L$. To understand the difference between c and s , we can say, that c is the ID for the case of patient X. Henceforth, s^c reflects all interactions that the database has logged for patient X.

These events are ordered in the sequence, in which they occurred for patient X. Therefore, let \mathcal{T} be the time domain and $\pi_t : E \mapsto \mathcal{T}$ a non-surjective linking function which strictly orders a set of events. In other words, every event in the database maps to one point in time. If the database logs every event on a daily basis, then all possible dates in history constitute \mathcal{T} . However, not every day has to be linked to a case as π_t is non-surjective.

Let \mathcal{A} be a universe of attribute identifiers, in which each identifier maps to a set of attribute values $\bar{a}_i \in \mathcal{A}$. An attribute identifier describes everything the database might store for a patient, such as heart-rate or blood sugar level. If the database logs the heart-rate, then heart-rates of -42 beats-per-minute are not possible. Hence, \bar{a}_i can per definition only map to positive integers.

Let \bar{a}_i correspond to a set of possible attribute values by using a surjective mapping function $\pi_A : \mathcal{A} \mapsto A$. Then, each event e_t consists of a set $e_t = \{a_1 \in A_1, a_2 \in A_2, \dots, a_I \in A_I\}$ with the size $I = |\mathcal{A}|$, in which each a_i refers to a value within its respective set of possible attribute values. In other words, every event consists of a set of values. If the event was recorded after a physio therapeutic session, then a_1 might be the specific degree to which you can move your ligaments and a_2 a description for the type of activity. If the event was recorded after a breast-cancer scan, the a_1 , a_2 and a_3 might relate to the specific diameter, the threat-level and again an indicator for the activity type. Conversely, we define a mapping from an attribute value to its respective attribute identifier $\pi_{\bar{a}} : A \mapsto \mathcal{A}$. Hence, we can map every event attribute value back to its attribute identifier.

The following part is not necessarily connected with *what* is stored within the database symbolically, but rather *how* it is represented in the database or during processing.

We require a set of functions F to map every attribute value to a representation which can be processed by a machine. Let $\pi_d : A_i \mapsto \mathbb{N}$ be a surjective function, which determines the dimensionality of a_i and also F be a set of size I containing a representation function for every named attribute set. We denote each function $f_i \in F$ as a mapper to a vector space $f_i : a_i \mapsto \mathbb{R}^d$, in which d represents the dimensionality of an attribute value $d = \pi_d(A_i)$. For instance categorical variables will map to a one-hot-encoded vector. Numerical values like heart-beat might be recorded in scalar form.

With these definitions, we denote any event $e_t \in s^c$ of a specific case c

as a vector, which concatenates every attribute representation f_i as $\mathbf{e}_t^c = [f_1; f_2; \dots; f_I]$. Therefore, \mathbf{e}_t^c is embedded in a vector space of size D which is the sum of each individual attribute dimension $D = \sum_i \pi_d(A_i)$. In other words, we concatenate all representations, whether they are scalar or vectors to one final vector representing the event. Furthermore, if we refer to a specific named attribute set A_i , we use the shorthand \bar{a}_i .

Figure 2.4 shows a schematic representation of a log L , a case c and an event e .

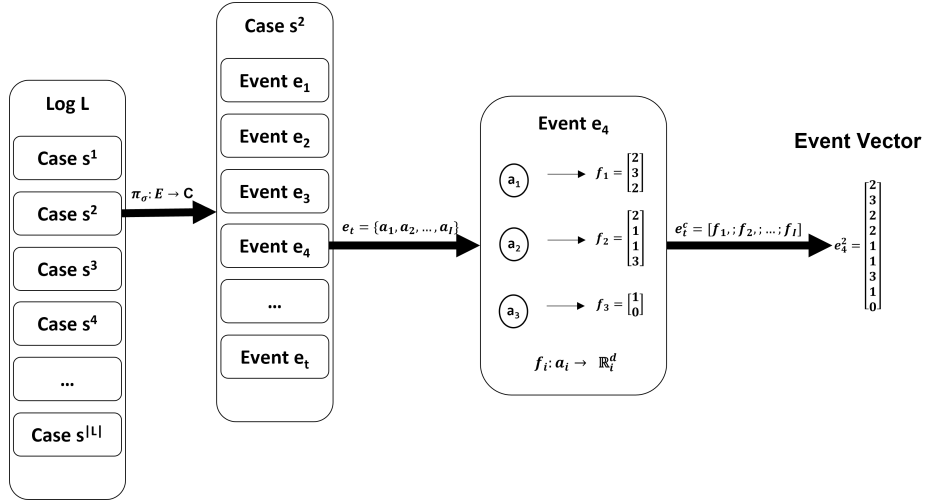


Figure 2.2: This figure shows the representation of a log L which contains anumber of cases s . Case s^2 contains a number of events e_t . Each events has attribute values a_i , which are mapped to vector spaces of varying dimensions. At last, all of the vectors are concatenated.

2.4.2 State-Space Models

Generally speaking, every time-series can be represented as a state-space model[27]. 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.

$$\begin{aligned} \mathbf{z}_{t+1} &= h(t, \mathbf{z}_t, \mathbf{u}_t) \\ \mathbf{e}_t &= g(t, \mathbf{z}_t, \mathbf{u}_t) \\ \mathbf{z}_{t+1} &:= \frac{d}{dt} \mathbf{z}_t \end{aligned} \tag{2.1}$$

Here, \mathbf{u}_t represents the input, \mathbf{z}_t the state at time t . The function h maps t , \mathbf{z}_t and \mathbf{u}_t to the next state \mathbf{z}_{t+1} . The event \mathbf{e}_t acts as an output computed by

function g which takes the same input as h . The variables \mathbf{z}_t , \mathbf{u}_t and \mathbf{e}_t are vectors with discrete or continuous features. The distinction of \mathbf{z}_{t+1} and \mathbf{e}_t decouples *hidden*² states, from *observable* system outputs. Figure 2.3 shows a graphical representation of these equations.

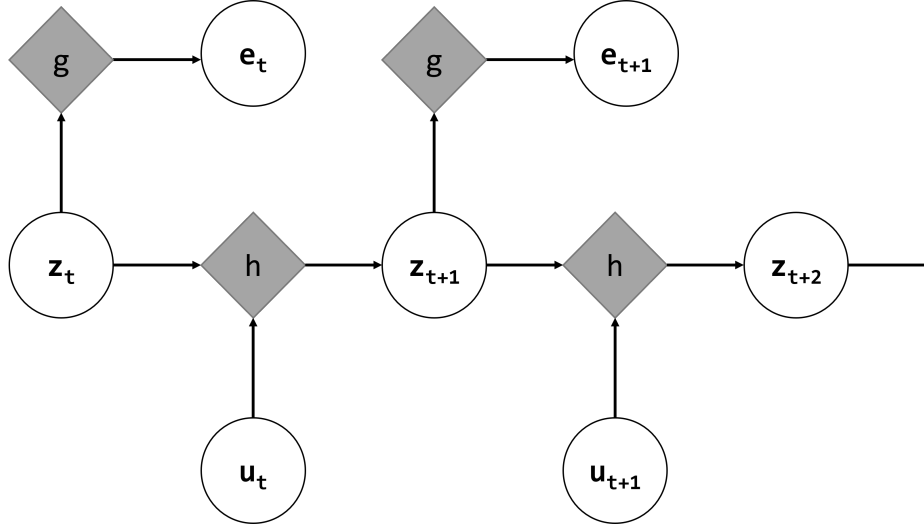


Figure 2.3: This figure shows a simplified graphical representation of a state-space model. Each arrow represents the flow of information.

The body of literature for state-space models is too vast to discuss them in detail. However, for process mining, we can use this representation to discuss the necessary assumptions for process mining. In line with the process-definition in section 2.1, we can understand the event log as a collection of the 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 treat t as a discrete scalar value to denote discrete sequential time steps. Hence, if we have $\sigma = \{a, b, b, c\}$, then t , describes the index of each element in σ . The input \mathbf{u}_t represents all context information of the process. Here, \mathbf{u}_t subsumes observable information such as the starting point and process instance-related features. The functions h and g determine the transition of a process' state to another state and its output over time. Note, that this formulation disregards any effects of future timesteps on the current timestep. Meaning, that the state transitions are causal and therefore, ignorant of the future. As we

²A 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. For an introduction to state-space models see:[22]

establish in section 2.1, we can assume that a process is a discrete sequence, whose transitions are time-variant. In this framework, we try to identify the parameters of the functions h and g . Knowing the functions, it becomes simple to infer viable counterfactuals. However, the function parameters are often unknown and therefore, we require probabilistic approaches.

We can formulate Equation 2.1 probabilistically as shown in Equation 2.2.

$$\mathbb{E}[p(z_{t+1} \mid t, z_{1:T}, u_{1:T}, x_{1:T}, \theta_h)] = \int z_{t+1} \cdot p(z_{t+1} \mid t, z_{1:T}, u_{1:T}, x_{1:T}, \theta_h) \quad (2.2)$$

$$\mathbb{E}[p(x_t \mid t, z_{1:T}, u_{1:T}, \theta_g)] = \int x_t \cdot p(x_t \mid t, z_{1:T}, u_{1:T}, \theta_g)$$

Note, that h and g are substituted with probability density functions parametrized with θ_h and θ_g . T signifies the full sequence including future timesteps. Both expectations are intractable as they require integrating over n -dimensional vectors. To solve the intractability, we characterize the system as a *Hidden Markov Process* and Probabilistic Graphical Model (PGM). This framework allows us to leverage simplifying assumptions such as the independence from future values and *d-separation*.

These characteristics change the probabilities in Equation 2.2 to Equation 2.3:

$$p(z_{t+1} \mid z_{1:t}, u_{1:t}, \theta_h) = \prod_{1}^t p(z_t \mid z_{1:t}, u_t, \theta_h) \quad (2.3)$$

$$p(x_t \mid z_{1:t}, \theta_g) = \prod_{1}^t p(x_{t-1} \mid z_{1:t}, \theta_g) \quad (2.4)$$

For $p(z_{t+1} \mid t, z_{1:T}, u_{1:T}, x_{1:T}, \theta_h)$, we ignore future timesteps, as T changes into t . *d-separation* allows us to ignore all \mathbf{e}_t of previous timesteps. The graphical form also decomposes the probability into a product of probabilities that each depend on all previous states and its current inputs. Previous \mathbf{e}_t are ignored due to *d-separation*. $p(x_t \mid t, z_{1:T}, u_{1:T}, \theta_g)$ only depends on its current state, which is in line with HMMs. Note, that we deliberately not assume a *strong Markov Property*, as the Deep Learning-Framework allows us to take all previous states into account. The *strong Markov Property* would assume that only the previous state suffices. At last, we assume that we do not model automatic or any other process whose state changes without a change in the input or previous states. Hence, we remove the dependency

on the independent t variable. Only the previous states $z_{1:T}$ and the input information \mathbf{u}_t remain time-dependent.

In this probabilistic setting, the generation of counterfactuals, amounts to drawing samples from the likelihood of Equation 2.3. We then use the samples to reconstruct the most-likely a counterfactual $e_{1:t}^*$. Hence, our goal is to maximize both likelihoods.

2.5 Representation

To process the data in subsequent processing steps, we have to discuss the way we encode the data. There are a multitude of ways to represent a log. We introduce four ways and the reason we choose the *hybrid-vector-representation*. Figure 2.4 shows schematically, how we can represent process data.

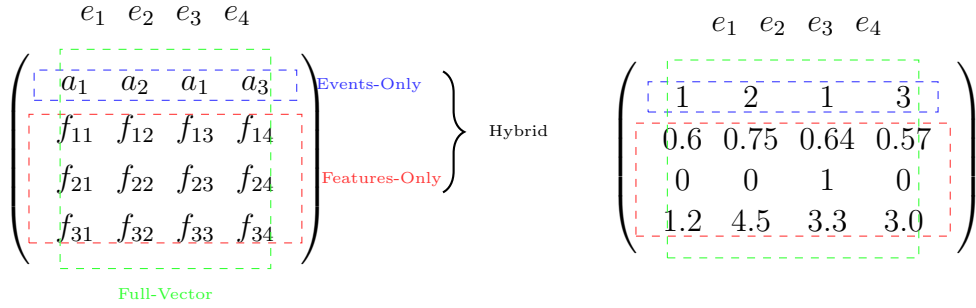


Figure 2.4: All four possible representations on an exemplary process instance.

First, we can choose to concentrate on *event-only-representation* and ignore feature attributes entirely. However, feature attributes hold significant amount of information. Especially in the context of using counterfactuals for explaining models as the path of a process instance might strongly depend on the event attributes. Similar holds for a *feature-only-representation*.

The first is a *single-vector-representation* with this representation we can simply concatenate each individual representation of every original column. This results in a matrix with dimensions (case-index, max-sequence-length, feature-attributes). The advantage of having one vector is the simplicity with which it can be constructed and used for many common frameworks. Here, the entire log can be represented as one large matrix. However, eventhough, it is simple to construct, it is quite complicated to reconstruct the former values. It is possible to do so by keeping a dictionary which holds the mapping between original state and transformed state. However, that requires every subsequent procedure to be aware of this mapping. Furthermore, we use

methods, that treat events and their associated features (event attributes) separately. For instance, if we want to sample from a markov model with transition probabilities and emission probabilities, then it is much easier to first sample the event trajectory and then, the conditional feature attributes. Or, if we attempt to compute an edit distance between two sequences, it is easier to compute those, if we keep events and event attributes separate.

Therefore, we decide to keep the original sequence structure of events as a separate matrix and complementary to the remaining event attributes. If required, we turn the label encoded activities ad-hoc to one-hot encoded vectors. Thus, this *hybrid-vector-representation* grants us greater flexibility. However, we now need to process two matrices. The first matrix has the dimensions (case-index, max-sequence-length) and the latter (case-index, max-sequence-length, feature-attributes).

2.6 Long-Short-Term Memory Models

In order to explain the decisions of a prediction we have to introduce a predictive model, which needs to be explained. Any sequence model suffices. Additionally, the model's prediction do not have to be accurate. However, the more accurate the model can capture the dynamics of the process, the better the counterfactual functions as an explanation of these dynamics. This becomes particularly important if the counterfactuals are assessed by a domain expert.

In this thesis, the predictive model is an Long Short-Term Memory (LSTM) model. LSTMs are well-known models within Deep Learning, that use their structure to process sequences of variable lengths[24]. LSTMs are an extension of Recurrent Neural Networks (RNNs). We choose this model as it is simple to implement and can handle long-term dependencies well.

Generally, RNNs are Neural Networks (NNs) that maintain a state h_{t+1} . The state is computed and then propagated to act as an additional input alongside the next sequential input of the instance x_{t+1} . The hidden state h is also used to compute the prediction o_t for the current step. The formulas attached to this model are shown in

$$h_{t+1} = \sigma(Vh_t + Ux_t + b) \quad (2.5)$$

$$o_t = \sigma(Wh_t + b) \quad (2.6)$$

Here, W , U and V are weight matrices that are multiplied with their respective input vectors h_t , x_t . b is a bias vector and σ is a nonlinearity

function. LSTM fundamentally work similarly, but have a more complex structure that allows to handle long-term dependencies better. They manage this behaviour by introducing additional state vectors, that are also propagated to the following step. We omit discussing these specifics in detail, as their explanation is not further relevant for this thesis. For our understanding it is enough to know that h_t holds all the necessary state information. Figure 2.5 shows a schematic representation of an RNN.

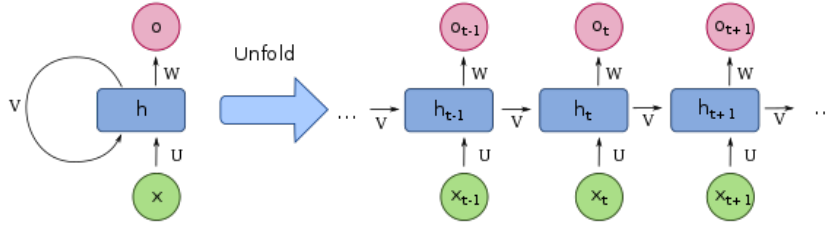


Figure 2.5: A schematic representation of an RNN viewed in compact and unfolded form.

2.7 Damerau-Levenshtein

The Damerau-Levenshtein distance function is a modified version of the Levenshtein distance[33], which is a widely used to compute the edit-distance of two discrete sequences[4, 38]. The most important applications are within the NLP discipline and the Biomedical Sciences. Within these areas, we often use the Levenshtein distance to compute the edit-distance between two words, two sentences or two DNA sequences. Note, that the elements of these sequences are often atomic symbols instead of multidimensional vectors. Generally, the distance accounts for inserts, deletions and substitutions of elements between the two sequences. Damerau modified the distance function to allow for transposition operations. For Process Mining, transpositions are important as one event can transition into two events that are processed in parallel and may have varying processing times. In Figure 2.6, we schematically show two sequences and their distance.

Equation 2.7 depicts the recursive formulation of the distance. The distance computes the costs of transforming the sequence a to b , by computing the

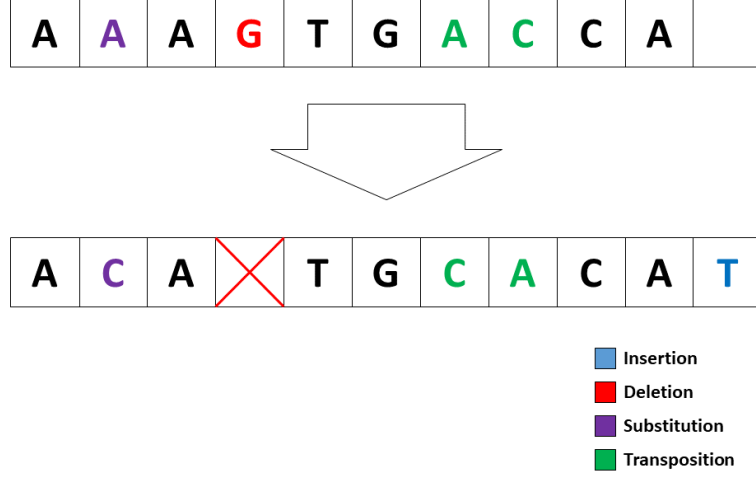


Figure 2.6: Two arbitrary sequences and their edit difference according to Damerau. The edit distance is the sum of each operation necessary to transform the sequence to another sequence. Blue shows an insert, red a deletion, purple a substitution and green a transposition. Therefore the edit distance is 4.

minimum of five separate terms.

$$d_{a,b}(i, j) = \min \begin{cases} d_{a,b}(i-1, j) + 1 & \text{if } i > 0 \\ d_{a,b}(i, j-1) + 1 & \text{if } j > 0 \\ d_{a,b}(i-1, j-1) + 1 & \text{if } i, j > 0 \\ d_{a,b}(i-2, j-2) + 1 & \text{if } i, j > 1 \wedge a_i = b_{j-1} \wedge a_{i-1} = b_j \\ 0 & \text{if } i = j = 0 \end{cases} \quad (2.7)$$

The recursive form $d_{a,b}(i, j)$ for sequences a and b with respective elements i and j takes the minimum of each of each allowed edit operation. In particular, no change, deletion, insertion, substitution and transposition. For each operation, the algorithm adds an edit cost of 1.

We cannot use the Damerau-Levenshtein distance for process mining, if the process carries additional information about event attributes. Mainly, because two events may be emitted by the same activity, but they may still carry different event attributes.

To illustrate the issue, we explore a couple of examples. Lets assume, we have two strings $s^1 = aaba$ and $s^2 = acba$. Using the Damerau-Levenshtein distance, the edit distance between both sequences is 1, as we can recognise a substitution at the second position in both strings. However, this representation is insufficient for process instances. Therefore, we now characterise the two sequences as process events rather than strings in Equation 2.8.

$$s^1 = \{a, a, b, a\} \quad (2.8)$$

$$s^2 = \{a, a^*, b, a\} \quad (2.9)$$

$$s^3 = \{a, c, b, a\} \quad (2.10)$$

$$s^4 = \{a, a, b\} \quad a, b, c \in \mathbb{R}^3 \quad (2.11)$$

$$a = \begin{bmatrix} 2 \\ 1 \\ 4 \end{bmatrix} \quad a^* = \begin{bmatrix} 3 \\ 3 \\ 4 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad c = \begin{bmatrix} 3 \\ 1 \\ 4 \end{bmatrix} \quad (2.12)$$

If we do not consider attribute values, it becomes clear that s^2 , s^3 and s^4 have an edit-distance to s^1 of 0, 1 and 1. However, with attribute values in mind, s^1 and s^2 display clear differences. Similarly, s^1 and s^3 not only differ in terms of activity but also attribute value. Lastly, s^1 and s^4 are the same in attribute values, but one element still misses entirely. It appears unintuitive that each of these differences are associated with the same cost. The examples show that we can neither disregard attribute values nor events, while computing the edit distance of two process instances.

Instead, we have to define a cost function which takes attribute variables into account. Therefore we modify the Damerau-Levenshtein distance by introducing a cost function instead of a static cost. Here, the cost of each edit-type is determined by a distance-function, which considers the difference between event-attributes. Therefore, we propose an edit-function, which captures structural sequence differences, as well as, content related differences. Going back to our example, if assume our cost function to only count differences in attributes, then the difference between s^1 and s^2 shall be 2 as their activities are the same, but the first two event attributes are different. To illustrate the structural elements, the difference between s^1 and s^3 shall be 3 instead of 2. Even if both a and c have two common event attributes, the activities they represent are still different. For instance, if both s^1 and s^3 were medical processes and a and c represented taking a cancer drug or a placebo, anyone would understand both activities are different even if the patient took the same dosage.

2.8 Evolutionary Algorithms

Many of our generative models are based on Evolutionary Algorithms. This section provides a small overview about this optimization technique.

All evolutionary algorithms use ideas that resemble the process of evolution. There are four broad categories: A Genetic Algorithm (GA) uses bit-string representations of genes, while Genetic Programming (GP) uses binary codes to represent programs or instruction sets. Evolutionary Strategy (ES) require the use of vectors. Lastly, Evolutionary Programming (EP), which closely resembles ES, without imposing a specific data structure type[31, 52]. Our approach falls into the category of GA. We refer to the literature review of Vikhar for more insights into the field. The most vital concept in this category is the *gene* representation. For our purposes, the gene of a sequence consists of the sequence of events within a process instance. Hence, if an offspring inherits one gene of a parent, it inherits the activity associated with the event and its event attributes.

Our goal is to generate candidates by evaluating the sequence based on our viability measure. Our measure acts as a fitness function. The candidates that are deemed fit enough are subsequently selected to reproduce offspring. The offspring is subject to mutations. Then, we evaluate the new population repeat the procedure until a termination condition is reached. It differs from gradient-based methods such as Deep Learning, because it does not require us to use differentiable functions. Hence, we can directly optimise the viability measure established in section 3.3.

For the algorithm, we follow a rigid structure of the operations as outlined in 1. As 1 shows, we define 5 fundamental operations. Initiation, Selection, Crossover, Mutation and Recombination.

Algorithm 1 The basic structure of an evolutionary algorithm.

Require: Hyperparameters

Ensure: The result is the final population

```
population  $\leftarrow$  INIT population;  
while not termination do  
    parents  $\leftarrow$  SELECT population;  
    offspring  $\leftarrow$  CROSSOVER parents;  
    mutants  $\leftarrow$  MUTATE offspring;  
    survivors  $\leftarrow$  RECOMBINE population  $\cup$  mutants;  
    termination  $\leftarrow$  DETERMINE termination  
    population  $\leftarrow$  survivors  
end while
```

Initiation

The initiation process refers to the creation of the initial set of candidates for the selection process in the first iteration of the algorithm. Often, this amounts to the random generation of individuals. In this thesis, we call this method the *Random-Initiation*. However, choosing among a subset of the search space can allow for a faster convergence. We chose to implement three different subspaces as a starting point. First, by sampling from the data distribution of the Log (*Sampling-Based-Initiation*). Second, by picking individuals from a subset of the Log (*Case-Based-Initiation*).

Selection

The selection process chooses a set of individuals among the population according to a selection procedure. These individuals will go on to act as material to generate new individuals. Again, there are multiple ways to accomplish this. In this thesis, we explore three methods. First, the *Roulette-Wheel-Selection*. Here, we compute the fitness of each individual in the population and choose a random sample proportionate to their fitness values. Next, the *Tournament-Selection*, which randomly selects pairs of population individuals and uses the individual with the higher fitness value to succeed. Last, we select individuals based on the elitism criterion. In other words, only a top-k amount of individuals are selected for the next operation (*Elitism-Selection*). This approach is deterministic and therefore subject to getting stuck in local minima.

Crossover

Within the crossover procedure, we select random pairing of individuals to pass on their characteristics. Again allowing a multitude of possible procedures. We can uniformly choose a fraction of genes of one individual (*Parent 1*) and overwrite the respective genes of another individual (*Parent 2*). The result is a new individual. We call that (*Uniform-Crossover*). Figure 2.7 shows a simple schematic example. By repeating this process towards the opposite direction, we create two new offsprings, which share characteristics of both individuals. The amount of inherited genes can be adjusted using a rate-factor. The higher the crossover-rate, the higher the risk of disrupting possible sequences. If we turn to Figure 2.7 again, we see how the second child has 2 repeating genes at the end. If a process does not allow the transition from *activity 8* to another *activity 8*, then the entire process instance becomes infeasible.

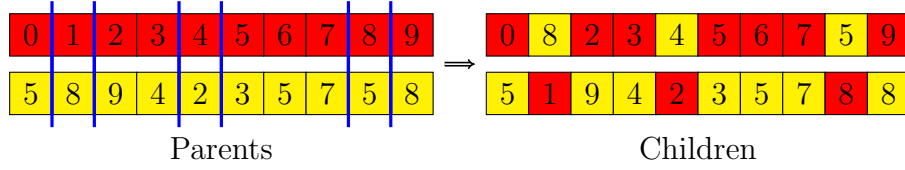


Figure 2.7: A crossing process of uniformly applying characteristics of one sequence to another.

The second approach is suitable for sequential data of same lengths. We can choose a point in the sequence and pass on genes of *Parent 1* onto the *Parent 2* from that point onwards and backwards (*One-Point-Crossover*). Thus, creating two new offsprings again as depicted in Figure 2.8.

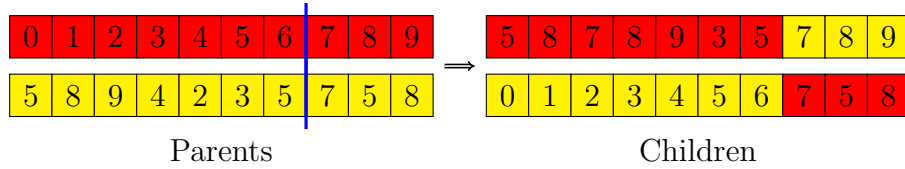


Figure 2.8: A One-Point example of applying characteristics of one sequence to another using one split point

The last option is called *Two-Point-Crossover* and resembles its single-point counterpart. However, this time, we choose two points in the sequence and pass on the overlap and the disjoints to generate two new offsprings. Again, Figure 2.9 describes the procedure visually.

Obviously, we can increase the number of crossover points even further. However, this increase comes at the risk of disrupting sequential dependencies.

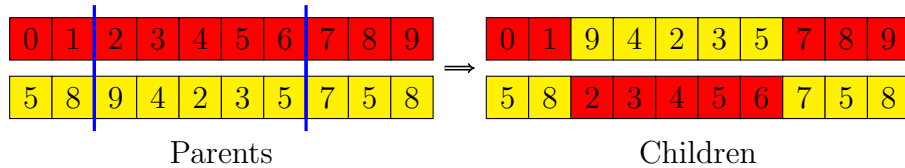


Figure 2.9: The process of applying characteristics of one sequence to another using two split points.

Mutation

Mutations introduce random perturbations to the offsprings. Here, we apply only one major operation. However, the extend in which these mutations are applicable can still vary.

Before elaborating on the details, we have to briefly discuss four modification types that we can apply to sequences of data. Reminiscent of edit distances, which were introduced earlier in this thesis, we can either insert, delete or change a step. These edit-types are the fundamental edits we use to modify sequences. For a visual explanation of each edit-type we refer to Figure 2.6 in section 2.7.

However, we can change the rate to which each operation is applied over the sequence. We call these parameters *mutation-rates*. In other words, if the delete-rate equals 1 every individual experiences a modification which results in the deletion of a step. Same applies to other edit types. There are still three noteworthy topics to discuss.

First, these edit-types are disputable. One can argue, that change and transpose are just restricted versions of delete-insert compositions. For instance, if we want to change the activity *Buy-Order* with *Postpone-Order* at timestep 4, we can first, delete *Buy-Order* and insert *Postpone-Order* at the same place. Similar holds for transpositions, albeit more complex. Hence, these operations would naturally occur over repeated iterations in an evolutionary algorithm.

However, these operations follow the structure of established edit-distances like the Damerau-Levenshtein distance. Furthermore, they allow for efficient restrictions with respect to the chosen data encoding. For instance, we can restrict delete operations to steps that are not padding steps. In contrast insert operations can be restricted to padding steps only.

Second, we could introduce different edit-rates for each edit-type. However, this adds additional complexity and needlessly increases the search space for hyperparameters.

Third, as we chose the hybrid encoding scheme, we have to define what an insert or a change means for the data. Aside from changing the activity, we also have to choose reasonable data attributes. This necessity requires to define two ways to produce them. We can either choose the features randomly, or choose to sample from a distribution which depends on the previous activities. We name the former approach *Default-Mutation*. We can simplify the latter approach by invoking the markov assumption and sample the feature attributes given the activity in question (*Sample-Based-Mutation*).

Recombination

This operation decides which individuals remain in the population for the next iteration¹. Here, we introduce three variations.

We name the strict selection of the best individuals among the offsprings

and the previous population *Fittest-Survivor-Recombination*. This recombiner strictly optimizes the population and is susceptible to getting stuck in local minima. In contrast, we name the addition of the top-k best offsprings to the initial population *Best-of-Breed-Recombination*. The former will guarantee, that the population size remains the same across all iterations but is prone to local optima. The latter only removes individuals after a population threshold was reached. Afterwards, the worst individuals are removed to make way for new individuals. Furthermore, we propose one additional recombination operator. The operator selects the new population in a different way than the former recombination operators. Instead of using the viability directly, we sort each individual by every viability component, separately. This approach allows us to select individuals regardless of the scales of every individual viability measure. We refer to this method as *Ranked-Recombination*.

¹We have to point out that in the literature, recombination is often synonymous with crossover. Both steps are similar in their filtering purpose. However, the selector filters potential parents while the recombiner filters the population. However, in this thesis recombination refers to the update process which generates the next population.

Chapter 3

Methods

In this chapter, we describe details of our framework **Name it** and discuss advantages and limitations. Therefore, we provide a more detailed overview and additionally describe all components. As the framework resembles the work of Hsieh, Moreira, and Ouyang, we also discuss differences and similarities between both solutions.

3.1 Methodological Framework

3.1.1 Architecture

To generate counterfactuals, we need to establish a conceptual framework, which consists of three main components. The three components are shown in Figure 3.1.

The first component is a predictive model. As we attempt to explain model decisions with counterfactuals, the model needs to be pretrained. We can use any model that can predict the probability of a sequence. This condition holds for models trained for process outcome classification and next-activity prediction. The model used in this thesis is a simple LSTM model using the process log as an input. The model is trained to predict the next action given a sequence.

The second component is a generative model. The generative model produces counterfactuals given a factual sequence. In our approach, each generative model should be able to generate a set of counterfactual candidates given one factual sequence. Specifically, we compare an evolutionary approach against 3 different generative baseline approaches. The baselines do not iteratively optimise towards viability criterions. All approaches allow us to use a factual sequence as a starting point for the generative production

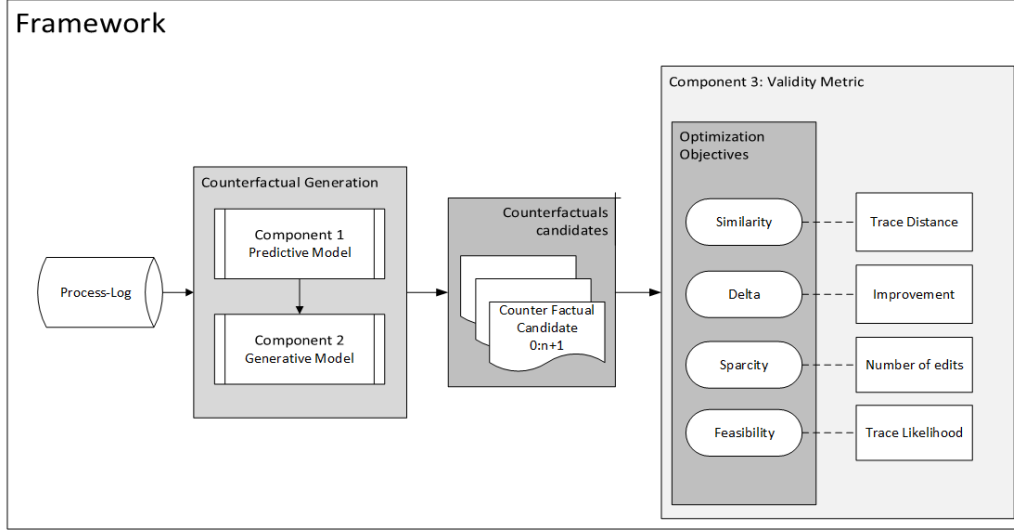


Figure 3.1: The methodological framework of this thesis. The input is the process log. The log will be used to train a predictive model (Component 1) and the generative model (Component 2). This process produces a set of candidates which are subject to evaluation via the validity metric (Component 3).

of counterfactuals. Furthermore, they also generate multiple variations of the final solution.

The generated candidates are subject to the third major component’s scrutiny. To select the most *viable* counterfactual candidate, we evaluate their viability score using a custom metric. The metric incorporates all main criterions for viable counterfactuals mentioned in section 2.3. We measure the *similarity* between two sequences using a multivariate sequence distance metric. The *delta* between the likelihood of the factual and the counterfactual. For this purpose, we require the predictive model, as it computes a predictions score that reflects the likelihood. We measure *sparsity* by counting the number of changes in the features and computing the edit distance. Lastly, we need to determine the *feasibility* of a counterfactual. This requires splitting the feasibility into two parts. First, the likelihood of the sequence of each event and second, the likelihood of the features given the event that occurred.

3.1.2 Differences to DiCE4EL

Hsieh, Moreira, and Ouyang has recently published a paper on the counterfactual generation of process data. They call their framework DiCE4EL and shares many ideas with our framework. Therefore, we want to highlight the key differences and similarities.

In their approach they attempt to solve various issues that we have also

highlighted in section 1.2. Furthermore, they do so by incrementally generating the model in a sequential order. However, unlike Oberst and Sontag, whose solution creates counterfactuals for every step in the sequence, Hsieh, Moreira, and Ouyang focus on critical decision points they call milestones.

To gain a better understanding, it is important to briefly outline the event log the authors use. It was taken from a Dutch bank which processes loan applications in which customers request a certain amount of money. The activities relate to either application states or manual work activities. The application states consist of tasks generated by a machine and manual work activities produced by humans. Hence, the manual work items occur in reference to the application state. For instance, if the loan application is in a *pre accepted* state, then the next events are often produced by humans who are reviewing the state. Those events are essentially, sub-events of the application state. The human activities do not have to happen sequentially. They can occur in parallel. The moment all manual work items conclude, marks the decision for the next application state. For instance, from *pre accepted* to *accepted*. Now, to understand why the milestone approach works, requires to know that an application loan process will change to a rejection state, for instance, if all manual work items are completed. There will not be applications that suddenly switch to another state although the work items of a previous state have not concluded, yet. Thus, one can split the entire sequence into subsequences or ignore the sub events entirely, which reduces the search space significantly.

One issue with this approach is the fact that one would first have to identify these milestones. Hence, a crucial distinction to our proposed framework, is their dependence on knowledge about the true process as displayed in this section. Our framework does not leverage structural information about the true process model in question. We believe this is the core contribution in contrast to their approach.

However, similarities between both frameworks do exist. Mainly, our approach also relies on prediction scores of the model we attempt to explain. Similar to Hsieh, Moreira, and Ouyang, we incorporate these scores into our quality measure.

3.2 Semi-Structured Damerau-Levenshtein Distance

Before discussing the viability function, we have to introduce an edit-distance for sequences. An edit-distance is used to compute distance between two

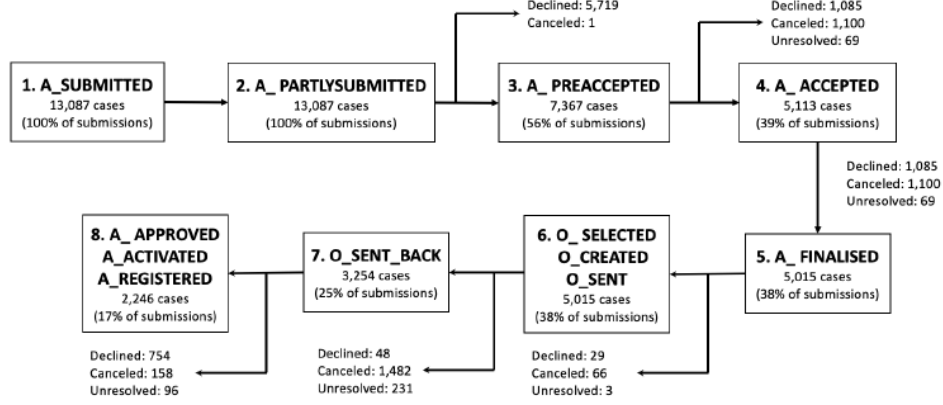


Figure 3.2: Milestones of loan application process captured in BPIC2012 as identified in [7]

sequences. Two distances take the *structural* characteristics of sequences into account.

3.2.1 Semi-Structured Damerau Levenshtein

In order to reflect these differences in attribute values, we introduce a modified version of the Damerau-Levenshtein distance, that not only reflects the difference between two process instances, but also the attribute values. We achieve this by introducing a cost function $cost_{a_i, b_j}$, which applies to a normed vector-space¹. Concretely, we formulate the modified Damerau-Levenshtein distance as shown in Equation 3.1. For the remainder, we refer to this edit-

distance as Semi-structured Damerau-Levenshtein distance (SSDLD).

$$d_{a,b}(i, j) = \min \begin{cases} d_{a,b}(i-1, j) + \text{cost}(\mathbf{0}, b_j) & \text{if } i > 0 \\ d_{a,b}(i, j-1) + \text{cost}(a_i, \mathbf{0}) & \text{if } j > 0 \\ d_{a,b}(i-1, j-1) + \text{cost}(a_i, b_j) & \text{if } i, j > 0 \\ & \& \bar{a}_i = \bar{b}_j \\ d_{a,b}(i-1, j-1) + \text{cost}(a_i, \mathbf{0}) + \text{cost}(\mathbf{0}, b_j) & \text{if } i, j > 0 \\ & \& \bar{a}_i \neq \bar{b}_j \\ d_{a,b}(i-2, j-2) + \text{cost}(a_i, b_{j-1}) + \text{cost}(a_{i-1}, b_j) & \text{if } i, j > 1 \\ & \& \bar{a}_i = \bar{b}_{j-1} \\ & \& \bar{a}_{i-1} = \bar{b}_j \\ 0 & \& i = j = 0 \end{cases} \quad (3.1)$$

Here, $d_{a,b}(i, j)$ is the recursive form of the Damerau-Levenshtein-Distance. a and b are sequences and i and j specific elements of the sequence. $\text{cost}(a, b)$ is a cost function which takes the attribute values of a and b into account. The first two terms correspond to a deletion and an insertion from a to b . The idea is to compute the maximal cost for that the wrongfully deleted or inserted event. The third term adds the difference between two events with identical activities \bar{a}_i and \bar{b}_j . As mentioned earlier, two events that refer to the same activity can still be different due to event attributes. The distance between the event attributes determines *how* different these events are. The fourth term handles the substitution of two events. Here, we compute the substitution cost as the sum of an insertion and a deletion. The fifth term computes the cost after transposing both events. This cost is similar to term 3 only that we now consider the differences between both events after they were aligned. The last term relates to the stopping criterion of the recursive formulation of the Damerau-Levenshtein distance.

3.2.2 Discussion

There are two important points to discuss, as they might incite disagreements with the validity of our viability measure.

If we assess the first two terms, we use $\text{cost}(x, 0)$ to denote the maximal distance of inserting and deleting x . $\text{cost}(x, 0)$ can be read as cost between x

¹A normed vector-space is a vector space, in which all vectors have the same dimensionality. For instance, if all vectors have three dimensions, we can call the vector-space *normed*.

and a null-vector of the same size. However, it is noteworthy to state that this interpretation does not hold for any arbitrary cost-function. For instance, the cosine-distance does not work with a null vector, as it is impossible to compute the angle between x and a null vector. Here, the maximum distance would just amount to 1. In contrast, the family of Minkowsky distance works well with this notion, because they compute a distance between two points and not two directions.

Furthermore, the intuition behind most of the terms, requires an established notion between events and their event attribute. Generally, we can have two notions of this relationship.

For the first relationship, we consider the event and its attributes as separate entities. This notion is reasonable, as some attributes remain static throughout the whole process run. If we take a loan application process as an example, an applicant's ethnic background does not change regardless of the event. This characteristic can be considered a case attribute, which remains static throughout the process run. This understanding would require us to modify the viability measures, as they treat the activity independently from its attribute values. In other words, if the activities of two events are \bar{a} and \bar{b} , but their attribute values are $(\frac{2}{3})$ and $(\frac{2}{3})$, these events may be seen as more similar than two \bar{a} and \bar{a} with attribute values $(\frac{2}{3})$ and $(\frac{5}{0})$.

In contrast, a second notion would treat each event as an independent and atomic point in time. Hence, a and b would be considered completely different even if their event attributes are the same. This understanding is also a valid proposition, as you could argue that an event which occurs at nighttime is not the same event as a daytime event. Here, the time domain is the main driver of distinction and the content remains a secondary actor.

All the terms described in the SSDLD follow the second notion. There are two reasons for this decision. First, treating event activities and event attributes separately would further complicate the SSDLD, as we would have to expand the cost structure to account for unchangable event attributes. Second, the unmodified Damerau-Levenshtein distance applies to discrete sequences, such as textual data with atomic words or characters. By treating each event as an discrete sequence element, we remain faithful to the original function.

3.3 Viability Measure

Earlier, in section 2.3, we have discussed what determines *good* counterfactuals. However, we have not introduced our approach to operationalize the notion of *viability*. To recall, a counterfactual is hardly useful, if it is vastly

different from the factual example or, if it requires changes that are logically implausible. For instance, if patients are required to vastly change their behavior in many aspects of their life or change their race these counterfactuals are hardly useful for the patient or a medical professional. We are more interested in what we have to change *at least*. Also, if the counterfactual is, per se, unrealistic or bears no change in outcome, we lack any interest in those counterfactuals, as well. For processes, these issues become even more complicated as they are semi-structured and often multivariate. How we operationalize these criteria is explained in the following.

3.3.1 Similarity-Measure

We use a function to compute the distance between the factual sequence and the counterfactual candidates. Here, a low distance corresponds to a small change. For reasons explained earlier (section 3.2), we want to take the structural distance and the feature distance into account. Henceforth, we use the previously established SSDLD. The similarity distance uses a cost function as specified in Equation 3.2.

$$\begin{aligned} cost(a_i, b_j) &= L2(a_i, b_j) \\ a_i, b_j &\in \mathbb{R}^d \end{aligned} \tag{3.2}$$

Here, $dist(x, y)$ is an arbitrary distance metric. i and j are the indices of the sequence elements a and b , respectively.

3.3.2 Sparsity-Measure

Sparsity refers to the number of changes between the factual and counterfactual sequence. We typically want to minimize the number of changes. However, sparsity is hard to measure, as we cannot easily count the changes. There are two reasons, why this is the case: First, the sequences that are compared can have varying lengths. Second, even if they were the same length, the events might not line up in such a way, that we can simply count the changes to a feature. Hence, to solve this issue, we use the previously established SSDLD. The sparsity distance uses a cost function as specified in Equation 3.3.

$$\begin{aligned} cost(a_i, b_j) &= \sum_d \mathbb{I}(a_{id} = b_{jd}) \\ a_i, b_j &\in \mathbb{R}^d \end{aligned} \tag{3.3}$$

Here, $\sum_d \mathbb{I}(a_{id} = b_{jd})$ is an indicator function, that is used to count the number of changes in a vector.

3.3.3 Feasibility-Measure

To determine the feasibility of a counterfactual trace, it is important to recognise two components.

First, we have to compute the probability of the sequence of event transitions. This is a difficult task, given the *Open World assumption*. In theory, we cannot know whether any event *can* follow after another event or not. However, if the data is representative of the process dynamics, we can make simplifying assumptions. For instance, we can compute the first-order transition probability by counting each transition. However, the issue remains that longer sequences tend to have a zero probability if they have never been seen in the data.

Second, we have to compute the feasibility of the individual feature values given the sequence. We can relax the computation of this probability using the *Markov Assumption*. In other words, we assume that each event vector depends on the current activity, but none of the previous events and features. Meaning, we can model density estimators for every event and use them to determine the likelihood of a set of features.

There are many ways to estimate the density of a data set. For our purposes, we incorporate the sequential structure of the log data and make simplifying assumptions. First, we consider every activity as a state in the case. Second, each state is only dependent on its immediate predecessor and neither on future nor on any any states prior to its immediate predecessor. Third, the collection of attributes within an event depend on the activity which emits it. The second assumption is commonly known as *Markov Assumption*. With these assumptions in place, we can model the distribution by knowing the state transition probability and the density to emit a collection of event attributes given the activity.

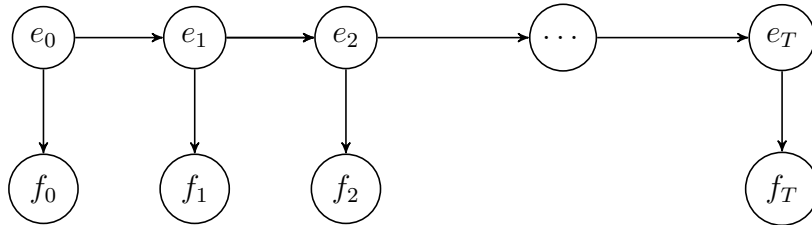


Figure 3.3: The feasibility model in graphical form. e_t represents an event and f_t the features it emits.

Here, e_t represents the transition from one event state to another. Likewise,

f represent the emission of the feature attributes. Hence, the probability of a particular sequence is the product of the transition probability multiplied with state emission probability for each step. Note, that this is the same as the feasibility measure as in Equation 3.4.

$$p(e_{0:T}, f_{0:T}) = p(e_0) p(f_0 | e_0) \prod_1^T p(e_t | e_{t-1}) p(f_t | e_t) \quad (3.4)$$

To conclude this section, we have to stress again, that there are many ways to define feasibility. We chose a probabilistic approach. There is an issue with this approach. Shorter sequences naturally have higher probabilities. Hence, we introduce a bias into our viability measure towards short sequences. This bias can be beneficial or detrimental depending on the use case. For instance, a medical process model might favor shorter counterfactual explanations. Mainly, because we want to understand how we can effectively reduce the time of illness. However, if we want to explain a highly standardised manufacturing process that went wrong in one instance; then, we would rather keep the counterfactual as close as possible to the factual.

3.3.4 Delta-Measure

For this measure, we evaluate the likelihood of a counterfactual trace by determining whether a counterfactual leads to the desired outcome or not. For this purpose, we use the predictive model, which returns a prediction for each counterfactual sequence. As we are predicting process outcomes, we typically predict a class. However, it is often difficult to force a deterministic model to produce a different class prediction. Therefore, we can relax the condition by maximising the prediction score of the desired counterfactual outcome[39]. If we compare the difference of the counterfactual prediction score with the factual prediction score, we can determine an increase or decrease. Ideally, we want to increase the likelihood of the desired outcome. We refer to this value as *delta*. However, the binary case introduces some noteworthy considerations.

Within this task setting, we have to consider multiple cases. First, prediction score which is typically limited to a domain within 0 and 1, which we can interpret as probability distribution. Hence, if the model score is 0.6, then the model has a confidence of 60% that the input can be categorised as belonging to class 1. For instance, within a medical process we could say, the model is 75% confident that the patient can be cured. Conversely, there's a 25% percent confidence that the process instance belongs to class 0. We can

make decisions by using a threshold. Typically this threshold is lies at 50%. Hence, we determine that a patient can be cured if the model's confidence exceeds 0.5². If we want to determine a soft version of the delta measure, we have to take this decision threshold into account.

We identify 2 cases:

Case 1: A counterfactual generator *flips* the prediction score to the opposite side of the decision threshold. Then, we archieve our general aim and the difference between the scores is a direct indicator of the counterfactual's success.

Case 2: A counterfactual does not change the factual decision. For instance, when the counterfactual and the factual prediction score for a patients recovery chance are below 0.5 or both are above 0.5. Then we have to consider the whether the is moving towards the desired outcome or away from it.

[2.1] If the prediction for the factual decides an outcome of 0 but the predictions score for the counterfactual is even lower, then we did not change the prediction at all. In fact, we increase the chance of the factual outcome. That situation is worse than what we desire. For instance, a patient would not want to pursue a counterfactual situation in which his odds of recovery are worse than his current.

[2.2] In contrast, if a prediction model's score leads to an outcome of 0 but the counterfactual returns a higher prediction score than for the factual predictions score, a patient might still be interested in the counterfactual. In some situations even a small improvement is desirable.

The subcases of 2 go in both ways. Hence, we have to incorporate each case differently in the delta score.

3.3.5 Discussion

Given the current viability function we can already determine the optimal counterfactual:

The optimal counterfactual flips the strongly expected factual outcome of a model to a desired outcome, maintaining the same trajectory as the factual in terms of events, with minimal changes

²Obviously, the domain of the application decides where this threshold lies. One can always aregue that a confidence of 51% is close to randomly guessing.

its event attributes, while remaining feasible according to the data.

The elements that fulfill these criteria make up the pareto surface of this multi-valued viability function. If each of the values are scaled a range between 0 and 1, the theoretical ceiling is 4. This value is only possible if we can flip the outcome of a factual sequence without changing it. As this is naturally impossible for deterministic model predictions, the viability has to be lower than 4.

Furthermore, we can already postulate, that a viability of 2 is an important threshold. If we score the viability of a factual against itself, a normalised sparsity and similarity value have to at its maximal value of 1. In contrast, the improvement has to be 0. The feasibility is 0 depending on whether the factual was used to estimate the data distribution or not. With these observations in mind, we determine that any counterfactual with a viability of at least 2 is already better than the factual.

3.3.6 Differences to DiCE4EL

Hsieh, Moreira, and Ouyang follow a very similar pattern of assessing the quality of their counterfactuals. The authors also focus on the aspects similarity, sparsity, feasibility and likelihood improvement. However, they incorporate and operationalize them differently. Their approach is mostly apparent in their loss function.

Similarity: Similar to our approach, the authors use a distance function and optimize it using gradient descent. They evaluate the quality of their counterfactuals using the same function³. However, we use a modified Damerau-Levenshtein distance algorithm to also incorporate structural differences such as the sequence lengths or transposed events.

Sparsity: The authors do not optimize towards sparsity, but assess it during their evaluation.

Feasibility: This quality criterion is embodied by two loss functions: The category loss and scenario Loss. The category loss ensures that categorical variables remain categorical after generation. The scenario loss adds emphasis on only generating counterfactuals that are in the event log. Unlike our probabilistic interpretation, they treat the existence of feasible counterfactuals as a binary criterion⁴.

³They call it proximity during evaluation

⁴They call it plausibility during evaluation

Likelihood: Similar to the authors’ scenario loss, they treat the improvement of a class as a binary state. Either the counterfactual changes the model’s prediction to the desired class or it does not.

The details of each criterion’s operationalisation, are explained in [26]. By assessing their interpretation of quality criterions, we see the clear distinction between our approach and the approach of Hsieh, Moreira, and Ouyang.

First, their viability measure decisively discourages the generation of counterfactuals that are not present in the dataset. In constast, our approach treats this aspect as a soft constraint.

Second, while our approach also acknowledges general improvements in likelihoods, DiCE4EL treats all counterfactuals that do not lead to better desired as detrimental solutions. However, one can argue that improving the likelihood of a desired outcome just slightly is already beneficial.

Third, [26] do not optimize sparsity, while we include it within our framework. One can argue that similarity automatically incorporates aspects of sparsity, but we disagree with this notion. We can see this by employing a simple example: Let factual A have features signifying the biological sex

(binary), the income (normalized) and the age (normalized) $\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$ as event

attributes. Let counterfactual B have the same event attributes with $\begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$.

Let’s assume the distance measure uses the L1-norm. Then, a counterfactual

C with event attributes $\begin{pmatrix} 1 \\ 0.5 \\ 0.5 \end{pmatrix}$, would have the same distance to factual A

than B has. However, C requires the change of two event attributes, while B only requires 1 change. In a scenario, in which we seek to reduce the number of edits, B is more preferable than C, regardless of the distance to A.

The last difference stems from the fact that Hsieh, Moreira, and Ouyang do not include structural sequence characteristics in their similarity measure. A sequence **XXZXX** might be more similar to **XXXZX**, than **XXXXZ**. The former requires only a transposition, while the latter requires two changes. Both have two positions that are not correct.

3.4 Prediction Model: LSTM

The architecture of the prediction model is shown in Figure 3.4. The model architecture is inspired by Hsieh, Moreira, and Ouyang. However, we do not

separate the input into dynamic features and static features.

One input consists of an 2-dimensional event tensor containing integers. The second input is a 3-dimensional tensor containing the remaining feature attributes. The first dimension in each layer represents the variable batch size and *None* acts as a placeholder.

The next layer is primarily concerned with preparing the full vector representation. We encode each activity in the sequence into a vector-space. We chose a dense-vector representation instead of a one-hot representation. We also create positional embeddings. Then we concat the activity embedding, positional embedding and the event attribute representation to a final vector representation for the event that occurred.

Afterwards, we pass the tensor through a LSTM module. We use the output of the last step to predict the outcome of a sequence using a fully connected neural network layer with a sigmoid activation as this is a binary classification task.

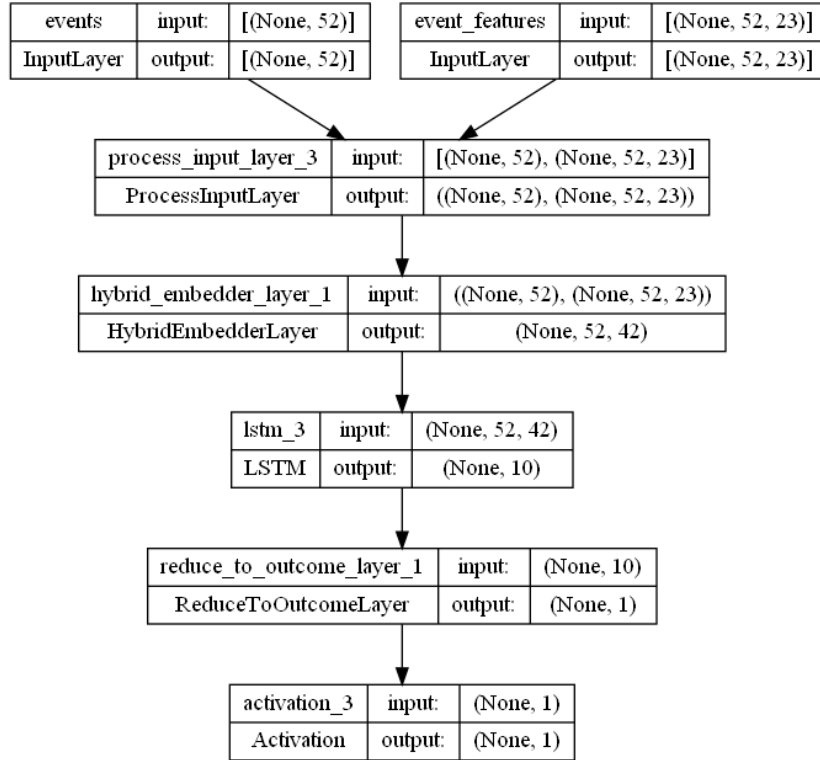


Figure 3.4: The different components of the LSTM architecture. Each elements contains information about the input and output of a layer. None is a placeholder for the batch size.

3.5 Counterfactual Generators

3.5.1 Baseline Model: Random Generator

This model acts as one of the baseline methods. Here, we generate a random sequence of events. Afterwards we generate event attributes, randomly. This approach is reasonably fast, but expected to perform poorly.

As explained earlier, a randomly sampling a possible sequence of events becomes more and more unlikely the longer the sequence is and the more events are possible. One generally has a chance of $\frac{1}{A^T}$ to randomly find an event, that is generally possible given the process model. The chances decrease even more if one also generates event attributes randomly. Therefore, we expect most models to perform better on average.

3.5.2 Baseline Model: Sample-Based Generator

The last baseline resembles the random baseline. However, we use the feasibility model to guide the random search for the generation of counterfactuals. We refer to the model specified in Equation 3.4. The sampling procedure utilises the model structure for the sampling process. We first generate a random seed of possible starting events ($p(e_0)$). Afterwards, we randomly sample subsequent events by iteratively sampling new activities according to the transition probabilities we gathered from the data ($\prod_1^T p(e_t | e_{t-1})$). Given the sequence, we simply sample the features per event from $p(f_t | e_t)$.

3.5.3 Baseline Model: Case-Based Generator

Case-based techniques leverage the data by using example instances. The idea is to find suitable candidates that fulfill the counterfactual criterions the best. We treat this model as a baseline. Therefore we keep this approach simple. We find candidates by searching by randomly sampling cases from the Log and then, evaluating them using the viability measure.

Inherently, this approach is restricted by the *representativeness* of the data. It is not possible to generate counterfactuals that have not been seen before. This method works for cases, in which the data holds enough information about the process. If this condition is not met, it is impossible to produce suitable candidates.

Note, that this approach will automatically fulfill the criterion of being feasible, as the counterfactuals are drawn from the Log directly. Hence, we expect their feasibility to often be higher than for other methods.

3.5.4 Generative Model: Evolutionary Algorithm

We introduced most of the operators in section 2.8. In this section, describe the operators in detail and select a subset that we want to explore further.

Operators

We implemented a number of different evolutionary operators. Each one belongs to one of five categories. The categories are initiation, selection, crossing, mutation and recombination.

Initiation

DI: The Default-Initiator generates an initial population entirely randomly.

SBI: The *Sampling-Based-Initiation* generates an initial population using a distribution estimated from the data.

CBI: *Case-Based-Initiation* uses examples of the data as initial population.

Selection

RWI: *Roulette-Wheel-Selection* Selects individuals randomly, but proportionate to their fitness score.

TS: *Tournament-Selection* Compares two or more individuals and selects a winner among them.

ES: *Elitism-Selection* selects each individual solely on their fitness

Crossing

OPC: *One-Point-Crossing* Chooses one point in the sequence and creates offspring by taking everything from or after that point from another individual.

TPC: *Two-Point-Crossing* Chooses two points in the sequence and creates offspring by taking everything between or outside these points from another individual.

UCx: *Uniform-Crossing* Uniformly selects positions in the sequence to take from another individual. The amount of selected positions is determined by a crossing-rate between 0 and 1.

Mutation

RM: *Random-Mutation* creates entirely random features for inserts and substitution.

SBM: *Sampling-Based-Mutation* creates sampled features based on data distribution for inserts and substitution.

Recombination

FSR: *Fittest-Survivor-Recombination* Determines the survivor among the mutated offsprings and the population.

BBR: *Best-of-Breed-Recombination* Determines better than average survivors among the mutated offsprings and adds them to the population.

RR: *Ranked-Recombination* Determines survivors based on sorting each individual by a determined order. The order we chose is to first focus on feasibility, then delta, then sparsity and at last similarity.

We use abbreviations to refer to them in figure, tables and so one. For instance, *CBI-RWI-OPC-RM-PR* refers to an evolutionary operator configuration, that samples its initial population from the data, probabilistically samples parents based on their fitness, crosses them on one point and so on. For the *Uniform-Crossing* operator we additionally indicate its crossing rate using a number. For instance, *CBI-RWI-UC3-RM-PR* is a model using the *Uniform-Crossing* with a child receiving roughly 30% of the genom of one parent and 70% of another parent.

Hyper Parameters

As with all models, the evolutionary approach comes with a number of hyper parameters. Generally, it is hard to determine the best set of hyperparameters as they interact and depend on the task setting. Nevertheless, it is important to discuss them.

We first discuss the model configuration. As show in this section, there are a number of ways to combine different operators. Depending on each individual operator, we might see very specific behaviours. For instance, it is obvious, that initiating the population with a random set of values can hardly converge at the same speed as a model which leverages case examples. Similarly, the selection of only the fittest individuals is heavily prone to local optima issues. The decision of the appropriate set of operators is by far the most important in terms of convergence speed and result quality.

The next hyperparameter is the *termination point*. Eventually, most correctly implemented evolutionary algorithms will converge to a local optimum. Especially if only the best individuals are allowed to cross over. If the termination point was chosen too early, then the generated individual will most likely underperform. In contrast, choosing a termination point too far in the future might yield optimal result at the cost of time performance. Furthermore, the existence of local optima may result in very similar solutions in the end. Optimally, we find a termination point, which finds a reasonable middle point.

The *mutation rate* is another important hyperparameter. It signifies how much a child can differ from its parent. Again, choosing a rate that is too low does not explore the space as much as it could. In turn, a mutation rate that is too high significantly reduces the chance to converge. The optimal mutation rate allows for exploring novel solutions without immediately pursuing suboptimal solution spaces. Our case is special, as we have four different mutation rates to consider. The change rate, the insertion rate, the deletion rate and the transposition rate. Naturally, these strongly interact. For instance, if the deletion rate is higher than the insertion rate there's a high chance that the sequence will be shorter, if not 0, at the end of its iterative cycles. Mainly, because we remove more events, than we introduce. However, we cannot assume this behaviour across the board as other hyperparameters interplay. Most prominently, the fitness function. Say, we have a high insertion rate but the fitness function rewards shorter sequences. Subsequently, both factors cancel each other out. Hence, the only way to determine the best set of mutation-rates requires an extensive search.

Chapter 4

Evaluation

In this chapter, we are going to establish most of the methods, the results section covers. In detail, we discuss the datasets, we use, the preprocessing pipeline and the final representation for each of the algorithms.

4.1 Datasets

In this thesis, we use a multitude of datasets for generating the counterfactuals. All of the data sets were taken from Teinemaa, Dumas, La Rosa, and Maggi. Each dataset consists of log data and contains labels which signify the outcome of a process. We focus on binary outcome predictions. Hence, each dataset will provide information about one of two possible outcomes associated with the case. For instance, a medical process might be deemed a success if the patient is cured or a failure if the patient remains ill. A loan application process might deem granting the loan a success or the rejection as failure. The determination of the outcome depends on the use-case and the stakeholders involved. An insurance provider might deem a successful claim as a failure, while the client deems it as a success.

BPIC12 The first dataset is the popular BPIC12 dataset. This dataset was originally published for the Business Process Intelligence Conference and contains events for a loan application process. Each individual case relates to one loan application process and can be accepted (regular) or cancelled (deviant).

Sepsis The next dataset is the Sepsis-Dataset. It is a medical dataset, which records of patients with life-threatening sepsis conditions. The outcome describes whether the patient returns to the emergency room within 28 days from initial discharge.

TrafficFines Third, we apply our approach to the Traffic-Fines-Dataset. This dataset contains events related to notifications sent related to a fine. The dataset originates in a log from an Italian local police force.

Dice4EL Lastly, we include a variation of the BPIC dataset. It is the dataset which was used by Hsieh, Moreira, and Ouyang. The difference between this dataset and the original dataset is two-fold. First, Hsieh, Moreira, and Ouyang omit most variables except two. Second it is primarily designed for next-activity prediction and not outcome prediction. We modified the dataset, to fit the outcome prediction model.

Dataset	#Cases	Min Len	Max Len	% Unique Traces	#Unique Ev.	#Data Columns	#Event Attr	#Regular	#Deviant
Dice4EL	3 051	12	25	0.000328	23	9	7	1 853	1 198
BPIC12-25	866	15	25	0.001155	32	23	21	682	184
BPIC12-50	3 728	15	50	0.000268	36	25	23	2 111	1 617
BPIC12-75	4 461	15	75	0.000224	36	25	23	2 379	2 082
BPIC12-100	4 628	15	100	0.000216	36	25	23	2 420	2 208
Sepsis25	707	5	25	0.001414	15	75	73	610	97
Sepsis50	770	5	47	0.001299	15	76	74	662	108
Sepsis75	777	5	66	0.001287	15	76	74	667	110
Sepsis100	779	5	88	0.001284	15	76	74	669	110
TrafficFines	129 615	2	20	0.000008	10	40	38	70 602	59 013

Table 4.1: All datasets used within the evaluation. Dice4EL is used for the qualitative evaluation and the remaining are used for quantitative evaluation purposes.

For more information about these datasets we refer to Teinemaa, Dumas, La Rosa, and Maggi’s comparative study[49]. We list all the important descriptive statistics in Table 4.1.

4.2 Preprocessing

To prepare the data for our experiments, we employed basic tactics for pre-processing. First, we split the log into a training and a test set. The test set will act as our primary source for evaluating factials, that are completely unknown to the model. We further split the training set into a training set and validation set. This procedure is a common tactic to employ model selection techniques. In other words, Each dataset is split into 25% Test and 75 remaining and from the remaining we take 25 val and 75 train.

First, we filter out every case, whose’ sequence length exceeds 25. We keep this maximum threshold for most of the experiments that forucs on the evolutionary algorithm. The reason is . Furthermore, two components of the proposed viability measure have a run time complexity of at least 2. Hence, limiting the sequence length saves a substantial amount of ressources.

Next, we extract time variables if they are provided in the log. Then, we normalise the values. For a proper time-format, we encode all information

from seconds to a year. If the full log occurs during one time-unit only, e.g. every event happened within a year, drop the column that was extracted. Afterwards, we standard scale all remaining time features.

Each categorical variable is converted using binary encoding. Binary encoding is very similar to onehot encoding. However, it is still distinct. Binary encoding uses a binary representation for each class encoded. This representation saves a lot of space as binary encoded variables are less sparse, than one-hot encoded variables.

We also add an offset of 1 to binary and categorical columns to introduce a symbol which represents padding in the sequence. All numerical columns are standardized to have a zero mean and a standard deviation of 1.

We omit the case id, the case activity and label column from this preprocessing procedure, for reasons explained in section 2.5. The case activity is label-encoded. Hence, every category is assigned to a unique integer. The label column is binary encoded, as we focus on outcome prediction.

4.3 Experimental Setup

As mentioned in section 2.3, counterfactual generation is notorious for their lack of a standardised evaluation procedure. Nonetheless, we attempt to address our research questions with the following experiments.

Experiment 1: Model Selection

Before comparing models, it is important to reduce the number of possible models that *can* be compared. Especially, the evolutionary generator has a number of free parameters. These range from structural configurations to general hyperparameters. In terms of operators, we introduced 3 initiators, 3 selectors, 3 crossers, 2 mutators and 3 recombiners. Hence, comparing all possible evolutionary operator combinations requires to examine a total of 54 different models. Furthermore, each model has hyperparameters, we have to define, too. Therefore, the first set of experiments are dedicated to choose among a subset of operator combinations and subsequently select appropriate hyperparameters.

First, we compute all possible configurations, without changing any hyperparameter. To avoid confusion, we refer to each unique operator combination as a model-configuration. For instance, one model-configuration would consist of a *SamplingBasedInitiator*, an *ElitismSelector*, a *OnePointCrosser*, *SamplingBasedMutator* and a *FittestSurvivorRecombiner*. For the sake of brevity, we refer to a specific model-configuration in terms of its abbreviated

operators. For instance, the earlier example is denoted as *SBI-ES-OPC-SBM-FSR*.

Afterwards, we explore the hyperparameters of the model. We start with the termination point. Hence, we want to explore the effects of the iterative cycles that each evolutionary algorithm will run for. The goal is to find a stopping criterion which yields reasonably good counterfactuals, while reducing the computation time. We will only consider the number of iterative cycles as a stopping criterion. We refer to each different criterion as termination point. Hence, a termination point at 5 means the algorithm, will not proceed to optimize its results, further after reaching the fifth iteration. We can choose the termination point by inspecting how the average population viability evolves across each cycle. We keep every other experimental setting as established beforehand.

For determining the mutation rate for every mutation type, we choose the best evolutionary algorithm and run the configuration with 6 rates from 0 to 0.5 in steps of 0.1. We omit everything beyond 0.5 to preserve information about the parent. For instance, if we use a change rate of 0.9, we mutate 90% of the genes the child inherited. This would defeat the purpose of evolving better counterfactuals through breeding. We use the termination point established in the prior experiment. We keep every other experimental setting as established beforehand.

After, executing all preliminary experiments we choose the evolutionary generators and compare them with all baseline models in all subsequent experiments.

Experiment 2: Model Comparison

First, we assess the viability of all the chosen evolutionary generators and the baseline generators. For this purpose, we sample 10 factials and use the models to generate 50 counterfactuals. We determine the median viability across the counterfactuals. With this experiment, we show that a model which optimizes quality criteria of counterfactuals produces better results than models, which do not. Hence, we expect the evolutionary algorithm to perform best, as it can directly optimize multiple viability criterions. In the following we list all models, we are going to compare:

In accordance with *RQ1-H1* and *RQ1-H2* we expect the evolutionary algorithms to outperform the baselines, when it comes to viability.

Experiment 3: Comparing with alternative Literature

The model comparison is not enough to establish the validity of our solution, as defined proposed the viability measure ourselves. Therefore, we also assess each model based on the evaluation criterions of an alternative work. More precisely, we quantify the viability of our models using the metrics employed by Hsieh, Moreira, and Ouyang. Hence, we measure the sparsity by computing the average Levenshtein difference and proximity using the L2-Norm. Furthermore, we compute the average intra-list-diversity and plausibility as well as the models capability of changing the prediction to a desired one.

Similar to Hsieh, Moreira, and Ouyang, we focus on the *activities* that are generated by each model and its accompanying *resource* event-attribute. For diversity and plausibility we remain close to the original evaluation protocol by Hsieh, Moreira, and Ouyang as we also treat each counterfactual trace sequence as a symbol. Hence, a sequence ABC is treated as a completely different symbol than $ABCD$.

The goal is to show that models, which optimise viability criterions, perform better, even if viability is assessed differently as stated in *RQ2-H1* of our research question (section 1.4).

Experiment 4: Qualitative Assessment

For the last assessment, we follow Hsieh, Moreira, and Ouyang’s procedure of assessing the models qualitatively. We use the dataset as the authors do. However, as we focus on outcome prediction, we attempt to answer one of two questions:

1. *what would I have had to change to prevent the cancellation/rejection of the loan application process*
2. *what would I have had to change to get cancelled/rejected of the loan application process*

The goal is to show, that the results are viable despite not having a standardized protocol to measure their viability.

Chapter 5

Results

This chapter presents the results of each evaluation step. Furthermore, we analyse the results.

5.1 Experiment 1: Model Selection

5.1.1 Model Configuration

Results

As there are many ways to combine each configuration, we select a few configurations by examining them through simulations.

The set of model-configuration contains 54 elements. We choose to run each model-configuration for 50 evolution cycles. For all model-configurations, we use the same 4 factual process instances, which are randomly sampled from the test set. We ensure, that the outcomes of these factuals are evenly divided. We decide to limit the population size to a maximum of 1000 counterfactuals. Within each evolutionary cycle, we generate 100 new offsprings. We keep the mutation rate at 0.01 for each mutation type. Hence, across all cases that are mutated, the algorithm deletes, inserts, and changes 1% of events per cycle. We collect the mean viability and its components across the iterative cycles of the model.

Figure 5.1 shows the bottom and top-5 model-configurations based on the viability after the final iterative cycle. We also show how the viability evolves for each iteration. The results reveal a couple of patterns. First, all of the top-5 algorithms use either *Case-Based-Initiator* as initiation operation. In contrast, the bottom-5 all use *Random-Initiator* as initialisation. Hence, the initiation appears to be majorly important for the algorithm. Second, we see that most of the top-5 algorithms use the *Elitism-Selector*. The complete

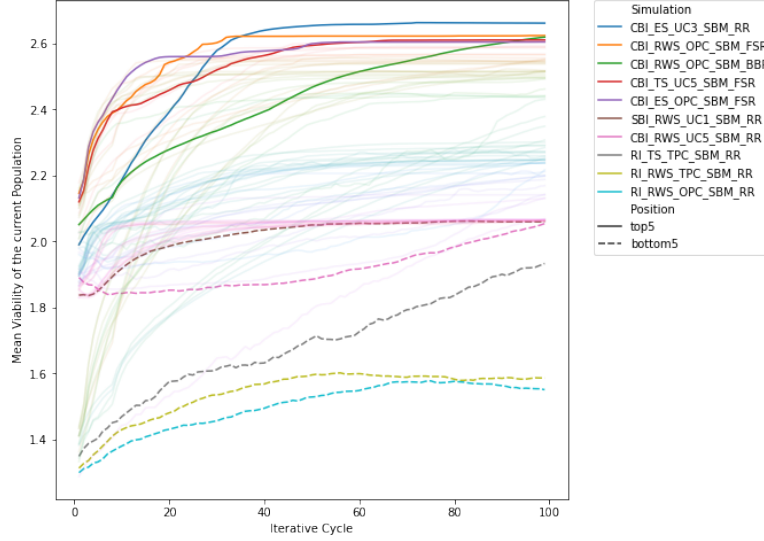


Figure 5.1: This figure shows the average viability of the 5 best and worst model-configurations. The x-axis shows how the viability evolves for each evolutionary cycle.

table of results is in ??.

In Figure 5.2, we see the effects of each operator type.

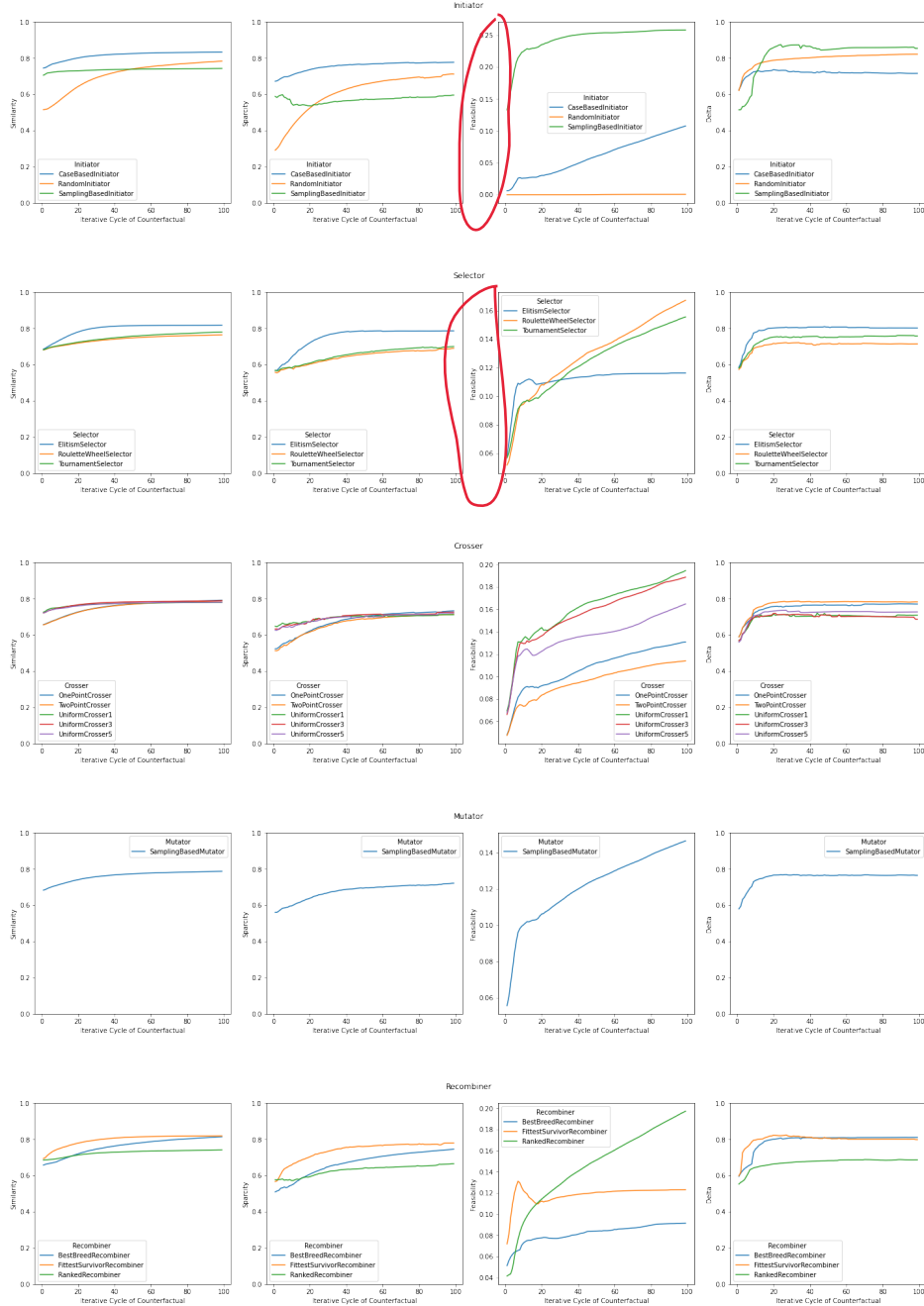
Starting with some commonalities across operator-type and measure, the figure shows that the initiator heavily determines the starting point for each measure. For instance, the *Random-Initiator* starts well below the other initiator operations when it comes to sparsity and similarity. Another noteworthy general observation is the delta measure. Here, for each operator type we see a movement towards the highest possible delta value. Hence, most configurations are capable of changing the source class to the desired class.

In terms of feasibility Figure ??, shows an increase for most operators. This is especially true if the operator has a random component or if it optimizes for feasibility. Similar holds for recombination with *Fittest-Survivor-Recombiner*.

The results for the selection operator type are undeniably in favor of *Elitism-Selector* for all viability measures. The same holds for the recombination operation. Here, the *FittestSurvivorRecombiner* yields better results.

When it comes to the crossing operation, the results indicate, the differences between *One-Point-Crosser* and *Two-Point-Crosser* are inconclusive for all viability measures except feasibility. One can explain that by noting, that both operations are very similar in nature. However, cutting the sequence only once produces less impossible sequences for the child sequences.

Interesting, feasibility seems to continuously improve while the other three measures seem to stabilize. Any explanation?



NICE results!

Figure 5.2: The evolution of each viability measure over the entire span of iterative cycles. Each figure adjust the respective operator type by taking the average over all other operator types.

Discussion

Moving forward, we have to choose a set of operators. We consider the following operators: We choose the *Case-Based-Initiator* as it might increase our chances to generate feasible variables.

For selection, we use the *Elitism-Selector*, as it generally appears to return better results.

Furthermore, we choose to move forward with the *One-Point-Crosser*. This crossing operation is slightly better in yielding feasible results.

For selection and recombination, we use the *Elitism-Selector* and the *FittestSurvivorRecombiner*, respectively. They all outcompete their alternatives for all similarity, sparsity and feasibility.

In the next experiment we vary mutation rates, using **SBI**-ES-OPC-SBM-FSR Generator .

CBI?

5.1.2 Model Termination Point

Results

For the experiment we chose a termination point of *200* which is twice the length of the previous simulation. We keep the mutation rate at *0.01* for each mutation type. The remaining procedure follows the process described in section 5.1.1.

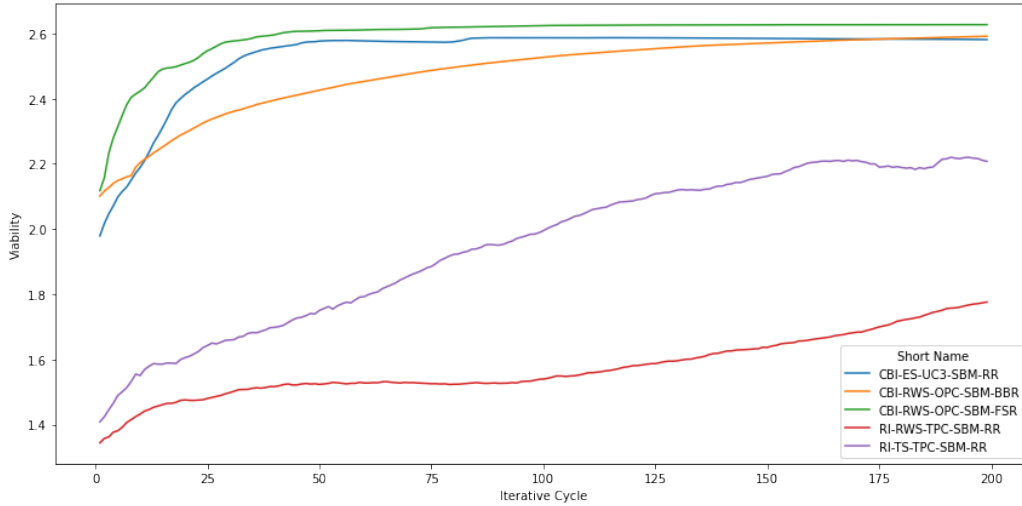


Figure 5.3: This figure shows the viability of across the iteration cycles.

In Figure 5.3, we see a general increase in viability for each termination point. It shows, that increasing the termination point also yields better

results at the end of the generation process. We see that *CBI-ES-UC3-SBM-RR* returns the best results in the shortest time span. The model converges after roughly 50 iterative cycles. *CBI-RWS-OPC-SBM-BBR* appears to have not reached convergence. The randomly initiated models have not reached convergence as well. However, they remain far below models that use a more sophisticated method to initialize their population.

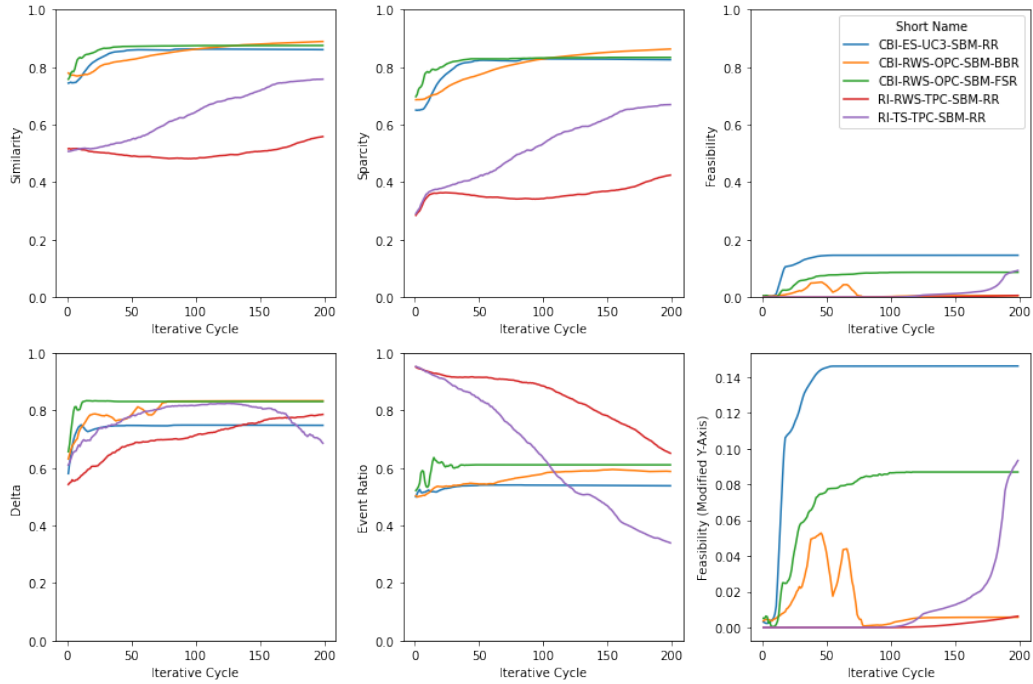


Figure 5.4: This figure shows the **remaining** measure components. Additionally, we show the ratio of events within the population. We also show a **magnified** version of the feasibility measure.

Figure 5.4 shows a decomposed view on how the viability measure evolves. Furthermore, we show the average amount of events within a generated counterfactual. In terms of similarity and sparsity all models behave similar. This is no surprise as both measures are inherently interlinked. We see that the randomly initiated models (RI-x) decrease the amount of events they generate. Case-based initiated models appear to slightly gain more events. Although, *CBI-RWS-OPC-SBM-BBR* appears that reaches its saturation point significantly later (100th cycle). Interestingly, the *CBI-RWS-OPC-SBM-BBR* model struggles to maintain feasibility and collapses to near 0 after the 100th iterative cycle. Another surprise is the steep ascent of only model that uses tournament selection (*RI-TS-TPC-SBM-RR*) towards the end of the generation process. The model even overtakes the model that leads the model-configurations in terms of *viability*. Furthermore, we see that *CBI-*

ES-UC-SBM-RR has the highest feasibility among all models. However, it also quickly converges after 50 iterative cycles.

Discussion

The results are not surprising. The longer the algorithm runs the closer it gets to a local minimum. We expect every evolutionary algorithm to converge at some point, as only the best within the population are chosen for the next iteration. If the model does not include enough non-deterministic components the results collapse to one optimal case in terms of structure. Hence, the counterfactual activities remain unchanged for the rest of the generation process. The events ratio should optimally approach a number around 0.5 if the factuals are evenly distributed in length. All model-configurations seem to follow this trajectory. However, models (*RI-TS-TPC-SBM-RR*) fall below this level. This coincides with its sharp rise in feasibility. We assume this behavior relates to a bias of the feasibility measure towards shorter sequences. The rise and decline of *CBI-RWS-OPC-SBM-BBR* shortly before overtaking all other models in terms of similarity and sparsity indicate a trade-off between how close the counterfactual is to the factual and how feasible it is.

For the next experiments we are going to use **50** as a termination point. It appears to be a reasonable point in which most models reach their highest viability yield and have not converged yet. We do not seek convergence, as it we want to maintain the diversity of our counterfactuals.

5.1.3 Model Parameters

Results

As explained earlier, for this simulation, we run the same configuration as beforehand established. However, this time we vary the rate with which we apply a mutation type.

As we can see in Figure 5.5, that a mutation rate of 0 yields better results on average. Suggesting that mutating the children might impede the model. For model-configurations that use the Fittest-Survivor-Recombination we observe a sharp **apattern** of convergence before the 10th iterative cycle.

Figure 5.6 reveals the reason for this behavior. In all plots of a **shapr** change right before the 10th iterative cycle. However, the feasibility measure also displays a sudden stop of improvement for all mutation rates except 0.0 and 0.4. These exceptions also change their rate of growth, but improve shortly after the 30th iterative cycle. The figure also shows that a mutation

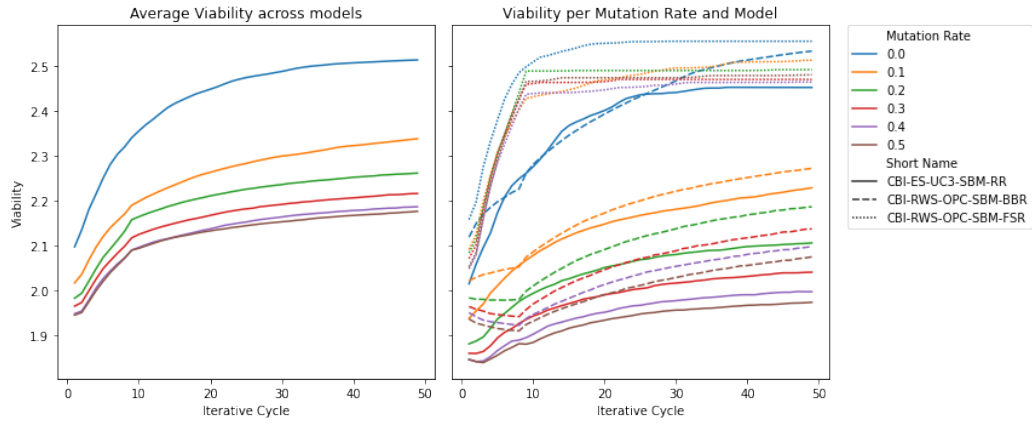


Figure 5.5: This figure shows the viability for each model and mutation rate per iterative cycle. The first plot shows the average across models. The second figure shows the same information per model. The x-axis shows how the viability evolves for each evolutionary cycle. The color indicates the mutation rate. The line-type marks each model tested.

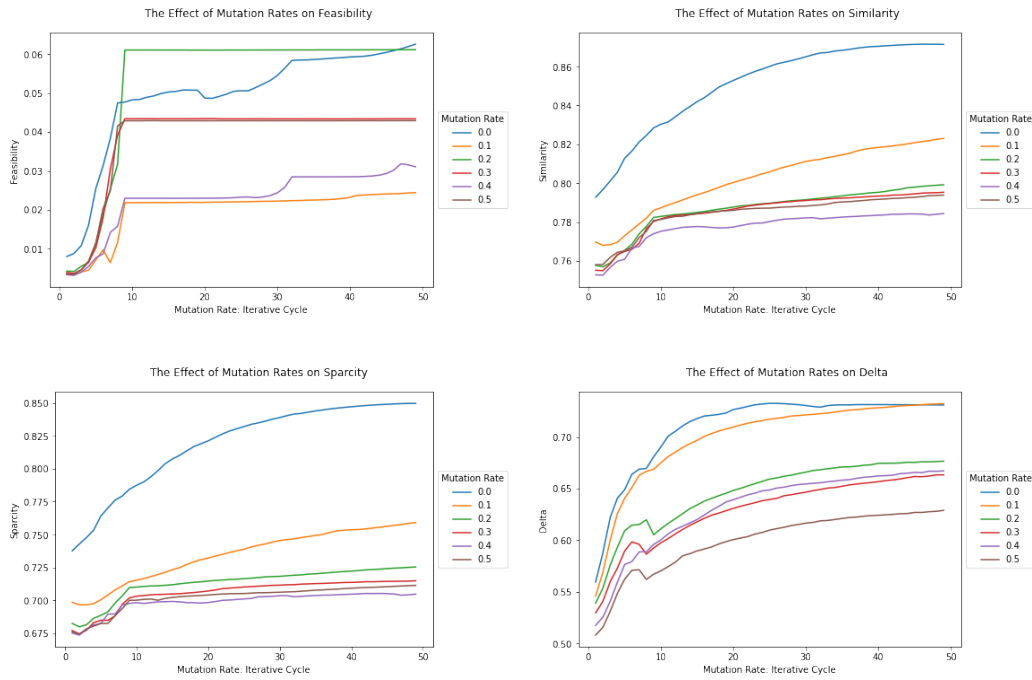


Figure 5.6: Shows all components of the viability measure.

rate of 0.2 reaches the highest feasibility among the other edit rates. However, after 48 cycles the mutation rate 0.0 overtakes 0.2.

Discussion

While it is expected, that every rate-configuration eventually converges towards an optimal value, it remains surprising that most rate-configurations suddenly converge around the 10 iteration. There are a multitude of possible reasons for this phenomenon. As the viability measure incorporates structural information and event-related information, we assume that the algorithm focuses on finding a structural optimum first.

Hence, the algorithm first prioritizes finding the best sequence in terms of activities. After finding a activity sequence, the model mostly focuses on improving the event attributes. Another explanation could be the ratio between the number of generated children and the population threshold. In this experiment, we generated 200 new children while limiting the population size to 1000.

With these observations in mind, we choose to set the mutation rate to 0.01. This decision implicates that mutations occur very rarely. Therefore, the main driving force for finding the best counterfactual is now the crossing operation. With this setting, we maintain the models ability to improve beyond 50th iterative cycle.

5.1.4 Model Candidates

To concude this section, we summarize the model selection, by choosing the models and their respective hyperparameters. Furthermore, we provide a quick overview of their characteristics. All models use the same mutator. Namely, the *Sample-Based-Mutator*.

- CBI-ES-UC3-SBM-RR This model initializes the first population actual using process instances. For each iterative cycle the individuals with the highest viability will go on to cross over their genom. Every child will receive 30% and 70% of its parents respectively. After the mutation phase, the generator reranks the full population and discards all individuals that have not meet reached the threshold. We choose this model as it promises to return the most feasible counterfactuals. However, the model most likely does not return the highest viability compared to other generators.
- BI-RWS-OPC-SBM-FSR This model initializes the first population using actual process instances. For each iterative cycle the individuals that pass on their genes are

probablistically selected based in proportion to their viability. For every child a crossover point decides how much of a parents genom are inherited. After the mutation phase, the generator selects the most viable individuals as the next population. We choose this model as it promises to return the highest value in terms of viability. However, the model is prone to reaching convergence very quickly.

5.2 Experiment 2: Model Comparison

In this section we examine the results of each model’s average viability across all datasets.

5.2.1 Results

In this comparison, we employ the other models mentioned in **??**. Namely, the *Case-Based Generator*, the *Sample-Based Generator* and the *Random Generator*.

We randomly sample 20 factuals from the test set and use the same factuals for every generator. We ensure, that the outcomes are evenly divided. The remaining procedure follows the established procedure of previous experiments.

the results shown in Figure 5.7 show that the evolutionary algorithm *CBI-ES-UC3-SBM-RR* slightly returns better results when it comes to the median viability. The worst model is the random generated model. The Case-Based model appears to be evenly and normally distributed at a viability of 2.25. The *CBI-RWS-OPC-SBM-FSR* has outliers that far exceed or underperform against other evolutionary algorithms.

Figure 5.7 also displays the vast difference in computation time for the evolutionary algorithms. Only the model using the *Ranking-Recombination* reranking version seems to be slightly faster than the ones using Best-Breed and Fittest-Survivor as recombination methods.

Table 5.1 shows the detailed results.

Table 5.1: Table shows the result of Experiment 4. The colors indicate the model configurations that were examined. The results are based on the average viability each counterfactual a model produces across all factuals that were tested.

Model (Abbr. Name)	Prediction Score	Viability	Sparcity	Similarity	Feasibility	Delta	Num. Paddings	Processing Time (sec.)	Max. Seq. Length
CBG-CBGW	0.514867	2.230507	0.764022	0.818115	0.014585	0.633786	14.584000	9.414627	27.000000
CBI-ES-UC3-SBM-RR	0.497746	2.678977	0.870874	0.896964	0.087737	0.823403	15.448000	588.550365	27.000000
CBI-RWS-OPC-SBM-BBR	0.445966	2.612767	0.851280	0.882271	0.095409	0.783807	15.560000	631.307437	27.000000
CBI-RWS-OPC-SBM-FSR	0.463966	2.728961	0.870071	0.899039	0.160373	0.799478	15.432000	625.714404	27.000000
RG-RGW	0.569685	1.554904	0.338077	0.578003	0.000000	0.638824	1.034000	8.175288	27.000000
SBG-SBGW	0.487669	2.151321	0.717582	0.755577	0.171964	0.506198	25.016000	9.927904	27.000000

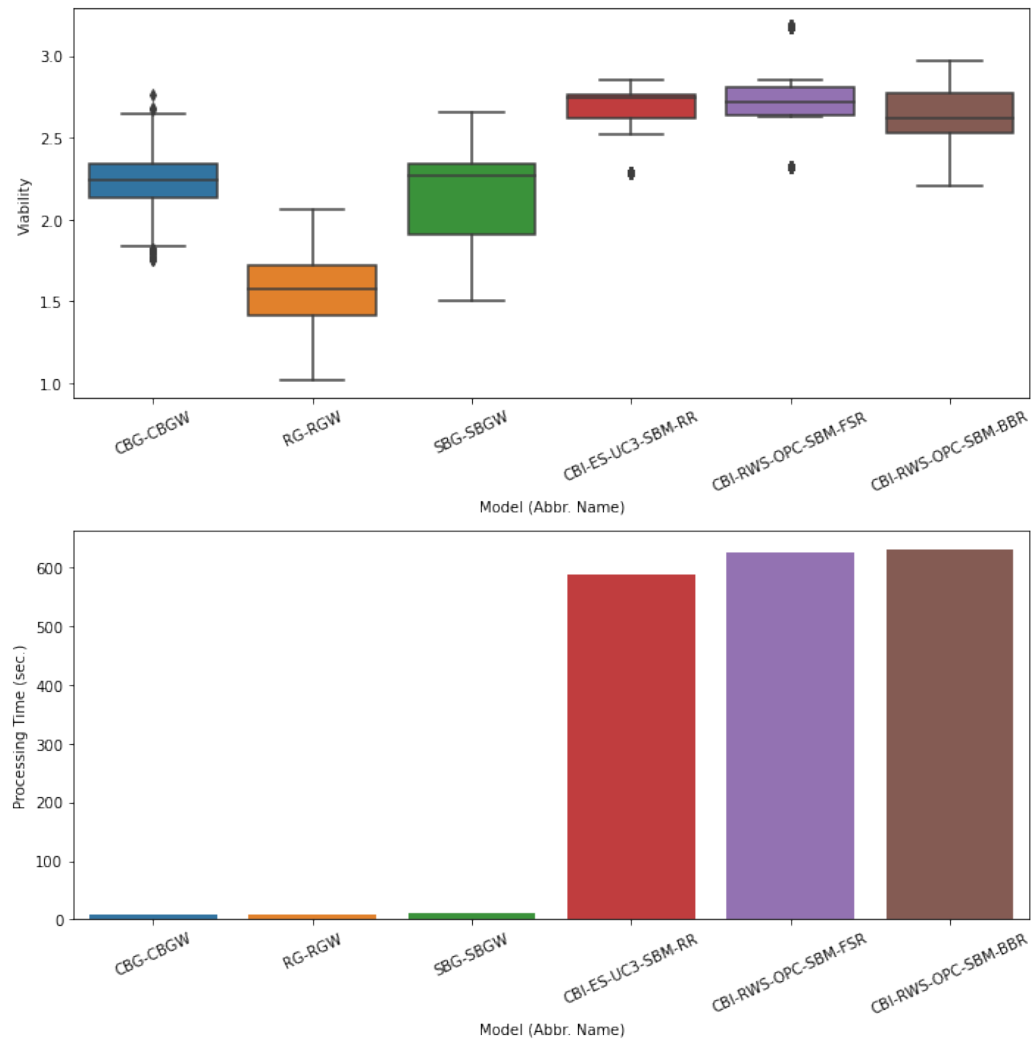
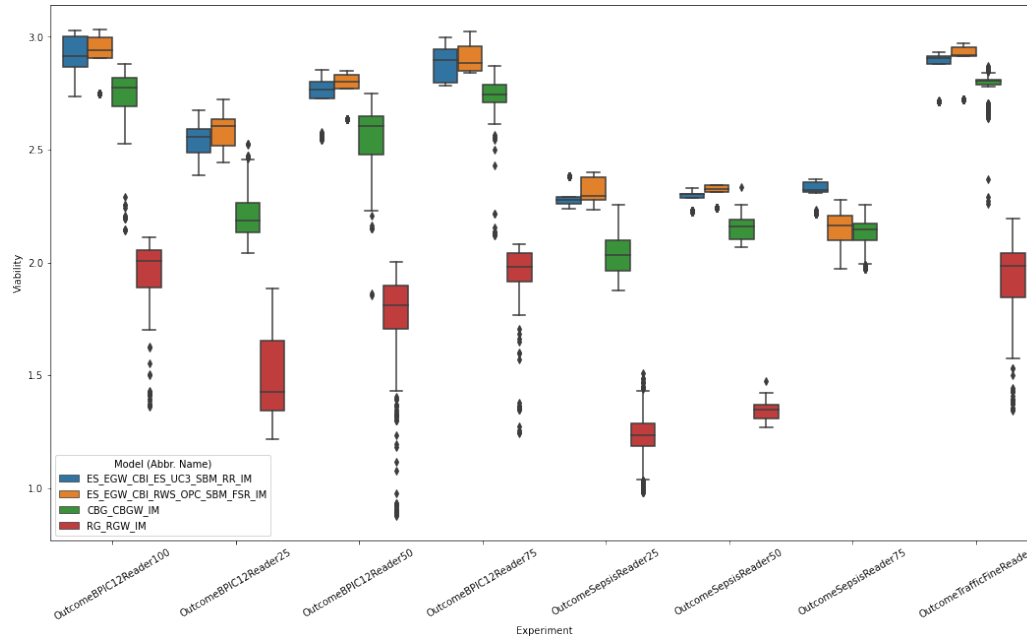


Figure 5.7: This figure shows boxplots of the viability of each models' generated counterfactual.

Figure 5.8 displays the results of running each algorithm on a set of factuals. The figure shows clear dominance of the evolutionary model all the models across all datasets.



Really NICE!

Any chance you can produce and add the same box plots for the four measures separately?

Figure 5.8: This figure shows boxplots of the viability of each models' generated counterfactuals across heterogeneous set of datasets.

Here, *CBI-ES-UC3-SBM-RR* and *CBI-RWS-OPC-SBM-FSR* display a higher median of viability across all datasets. This is surprising as the datasets are vastly different in sequence length and structure. The high median is reached for *CBI-ES-UC3-SBM-RR*.

5.2.2 Analysis

These results show that the model *CBI-RWS-OPC-SBM-FSR* is clearly superior to the other models. This result is unsurprising, as the baselines not actively search for an optimal solution. Furthermore, we see that the evolutionary models surpass their baselines by a wide margin.

The difference in computation time is most likely due to the similar and sparsity measures. The computation of the Damerau-Levenshtein distance is quadratic. As we also apply an additional custom cost function the computation times increase heavily, the longer the sequence. The evolutionary algorithm as described in section 2.8 is a sequential operation that also increases with the sequence length. However, we can deduce that the

I don't think it is very surprising that the EVO4EL approach outperform the baseline approaches, since the former uses the baselines as initiators if I understand correctly.

What I do find surprising is there seem to be significant improvements, the viability increased by 0.2. In 6 out of 8 data sets, orange seems to not even overlap with the green method. I would suggest highlight this. :)

Yes, agree! But it is surprising that the improvements are significant (two distributions do not overlap). Where do these significant improvements come from?

** It would be really nice to show and discuss some examples of counterfactuals from different approaches here. Just like the way you did in Section 5.4.

time difference between the *CBI-ES-UC3-SBM-RR* stems from either the *Ranking-Recombination* or the *Uniform-Crossing* operation. As those two are the only discernible operators.

In contrast, the baselines have been implemented in ways that vectorize most operations using numpy. Meaning, they can vastly decrease their computation time. The evolutionary algorithms, on the other hand, are subject to python's notorious ^{CITE s} low-looping operations. However, this is not a vital issue for two reasons. First, it is possible to run evolutionary algorithms in parallel manner ^{CITE}. Second, we have not explored more optimized implementations, of either the SSDLD or the evolutionary algorithm. However, we are certain, there are better and fast implementations available.

Knowing these results, a couple of questions remain. Namely, whether the results remain consistent for longer sequences and for other datasets? Furthermore, how does this procedure compare to other methods in the literature? The remaining experiments will address these questions.

The results for Figure 5.8 show that both evolutionary algorithms outperform the competition across all datasets and against all baselines. This is a remarkable result as it shows that the algorithm can outperform baselines regardless of the process log, and its length.

5.3 Experiment 3: Evaluation under a different Viability Measure

5.3.1 Results

Table 5.2 shows how each model scores under different operationalisations of viability aspects. They were derived from Hsieh, Moreira, and Ouyang's custom evaluation protocol and aim to provide a better comparison. Each value reflects the mean across all counterfactual results per model.

The results show that **diversity** is the highest for the evolutionary algorithm in terms of activity traces and resource traces. The Random-Search Generator displays **low diversity** for activities generated and a higher diversity for the resource.

Only the Casebased-Search Generator reaches a maximum score of 1 for plausibility. All the other models are far below or 0.

In terms of proximity, the Casebased-Search Generator has the lowest activity proximity. The average distance is 12.55. The SBI-ES-OPC-SBM-FSR Generator takes the second place. Interestingly, the gap between the proximity for activities is larger than the gap between proximities in terms of resources.

This is relatively difficult to follow.

I would suggest to shortly recap the meaning of each measure in the DiCE4EL approach for your readers.

Generator	Dimension	Model Property Iteration	Our Model Plausibility	Proximity	Sparsity	D4EL Plausibility	Proximity	Sparsity
CBG-CBGW-IM	Activity	0	0.320000	4.114943	9.000000	0.160000	4.178792	11.000000
		1	0.240000	3.862351	7.840000	0.120000	3.802004	6.420000
		2	0.160000	3.791798	7.680000	0.080000	3.766728	6.340000
		3	0.180000	4.179553	9.640000	0.090000	4.211097	9.320000
		4	0.280000	4.560320	12.260000	0.140000	4.625368	12.630000
	Resource	5	0.400000	4.258332	10.720000	0.200000	4.308616	10.360000
		0	0.000000	4.251903	17.920000	0.000000	4.724028	21.460000
		1	0.000000	3.818159	14.720000	0.000000	4.358569	19.360000
		2	0.000000	3.918192	15.500000	0.000000	4.557172	20.250000
		3	0.000000	4.283800	18.340000	0.000000	4.739976	20.670000
ES-EGW-CBI-ES-UC3-SBM-RR-IM	Activity	4	0.000000	4.681512	21.820000	0.000000	4.938832	22.410000
		5	0.000000	4.168343	17.400000	0.000000	4.682248	20.700000
		0	0.000000	4.123106	6.000000	0.000000	4.182873	9.500000
		1	0.000000	3.872983	2.000000	0.000000	3.807320	3.500000
		2	0.000000	3.741657	4.000000	0.000000	3.741657	4.500000
	Resource	3	0.000000	3.875524	3.020000	0.000000	4.059082	6.010000
		4	0.000000	4.582576	8.000000	0.000000	4.636496	10.500000
		5	0.000000	4.242641	7.000000	0.000000	4.300770	8.500000
		0	0.000000	4.123106	17.000000	0.000000	4.659629	21.000000
		1	0.000000	3.492392	11.200000	0.000000	4.195686	17.600000
ES-EGW-CBI-RWS-OPC-SBM-FSR-IM	Activity	2	0.000000	2.965685	8.800000	0.000000	4.080919	16.900000
		3	0.000000	4.242641	17.640000	0.000000	4.719397	20.320000
		4	0.000000	4.795832	23.000000	0.000000	4.995992	23.000000
		5	0.000000	4.000000	16.000000	0.000000	4.598076	20.000000
		0	0.000000	4.123106	6.000000	0.000000	4.182873	9.500000
	Resource	1	0.000000	3.741657	3.000000	0.000000	3.741657	4.000000
		2	0.000000	3.741657	4.000000	0.000000	3.741657	4.500000
		3	0.000000	4.000000	4.000000	0.000000	4.121320	6.500000
		4	0.000000	4.472136	5.000000	0.000000	4.581276	9.000000
		5	0.000000	4.242641	6.000000	0.000000	4.300770	8.000000
RG-RGW-IM	Activity	0	0.000000	4.000000	14.000000	0.000000	4.598076	19.500000
		1	0.000000	4.358899	16.000000	0.000000	4.628939	20.000000
		2	0.000000	2.236068	5.000000	0.000000	3.716110	15.000000
		3	0.000000	4.123106	15.000000	0.000000	4.659629	19.000000
		4	0.000000	4.582576	21.000000	0.000000	4.889364	22.000000
	Resource	5	0.000000	3.872983	15.000000	0.000000	4.534568	19.500000
		0	0.000000	4.832195	20.960000	0.000000	4.537418	16.980000
		1	0.000000	4.822574	20.380000	0.000000	4.282115	12.690000
		2	0.000000	4.856500	20.460000	0.000000	4.299079	12.730000
		3	0.000000	4.805536	20.800000	0.000000	4.524089	14.900000
RG-RGW-IM	Activity	4	0.000000	4.794502	21.740000	0.000000	4.742459	17.370000
		5	0.000000	4.731463	20.580000	0.000000	4.545181	15.290000
		0	0.000000	4.242641	18.000000	0.000000	4.719397	21.500000
		1	0.000000	3.741657	14.000000	0.000000	4.320318	19.000000
		2	0.000000	3.741657	14.000000	0.000000	4.468905	19.500000
	Resource	3	0.000000	4.242641	18.000000	0.000000	4.719397	20.500000
		4	0.000000	4.690416	22.000000	0.000000	4.943284	22.500000
		5	0.000000	4.358899	19.000000	0.000000	4.777526	21.500000

Table 5.2: A comparison between our model and D4EL

Again, the Casebased-Search Generator has the lowest sparsity with 9.34 in terms of activity but only remains slightly better than SBI-ES-OPC-SBM-FSR Generator in terms of resources.

The results suggest that the SBI-ES-OPC-SBM-FSR Generator is capable of producing very diverse counterfactual solutions, but cannot compete with the Casebased-Search Generator in terms of plausibility, proximity and sparsity. Hence, the Casebased-Search Generator is completely plausible given the data, is closer to the factual on average and displays less changes.

However, this only holds for the activities that are generated. In terms of resources that were generated, the Casebased-Search Generator is just slightly better.

5.3.2 Analysis

Based on these results, we can see that our model does seem to optimize properly for our viability function, but it does not compete under different operationalisations of counterfactual viability.

It is unsurprising, that the Casebased-Search Generator achieves the highest plausibility as all of the counterfactuals were drawn from the data itself.

In terms of the very similar resource proximity and sparsity scores, we assume that the model is able to identify the correct resource given the activity. For instance, if within a loan approval process only one person or machine executes the activity of checking the identity of an applicant, then SBI-ES-OPC-SBM-FSR Generator seems capable of learning this relationship.

5.4 Experiment 4: Qualitative Assessment

5.4.1 Results

In the result tables you can see some of the factials that were generated by our model and the model of [26].

Factual Seq. Amount	Activity	Outcome	Resource	Our CF Seq. Amount	Activity	Outcome	Resource	DICE4EL CF Seq. Activity	Resource	Amount
	A-SUBMITTED	0	112	155	A-SUBMITTED	1	112			
	A-PARTLYSUBMITTED	0	112	14214	A-PARTLYSUBMITTED	1	112			
	A-PREACCEPTED	0	101	14715	A-PREACCEPTED	1	112			
	W-Afhandelen leads	0	101	15372	A-ACCEPTED	1	9	A-SUBMITTED	112	
	A-ACCEPTED	0	111	138	O-SELECTED	1	912	A-PARTLYSUBMITTED	112	
	O-SELECTED	0	111	14962	A-FINALIZED	1	912	A-PREACCEPTED	112	
	A-FINALIZED	0	111	14887	O-CREATED	1	111	A-ACCEPTED	1	
	O-CREATED	0	111	14597	O-SENT	1	103	O-SELECTED	1	
	O-SENT	0	111	15235	W-Completeren aanvraag	1	111	A-FINALIZED	1	
	W-Completeren aanvraag	0	111	15473	W-Nabellen offertes	1	111	O-CREATED	1	
	W-Nabellen offertes	0	111					O-SENT	1	
	O-CANCELLED	0	111					W-Completeren aanvraag	1	
	A-CANCELLED	0	111					O-SENT-BACK	11259	
	W-Nabellen offertes	0	111	14474	W-Nabellen offertes	1	111	W-Nabellen offertes	11259	
				14715	A-REGISTERED	1	111	O-ACCEPTED	9	

Table 5.3: A comparison between the CBI-ES-UC3-SBM-RR and D4EL

Factual Seq. Amount	Activity	Outcome	Resource	Our CF Seq. Amount	Activity	Outcome	Resource	DICE4EL CF Seq. Activity	Resource	Amount
	A-SUBMITTED	0	112		A-SUBMITTED	1	112			
	A-PARTLYSUBMITTED	0	112		A-PARTLYSUBMITTED	1	112			
	A-PREACCEPTED	0	101		A-PREACCEPTED	1	112			
	W-Afhandelen leads	0	101					A-SUBMITTED	112	
	A-ACCEPTED	0	111		A-ACCEPTED	1	111	A-PARTLYSUBMITTED	112	
	O-SELECTED	0	111		O-SELECTED	1	111	A-PREACCEPTED	112	
	A-FINALIZED	0	111		O-FINALIZED	1	111	A-ACCEPTED	1	
	O-CREATED	0	111		O-CREATED	1	111	O-SELECTED	1	
	O-SENT	0	111		O-SENT	1	111	A-FINALIZED	1	
	W-Completeren aanvraag	0	111		W-Completeren aanvraag	1	111	O-CREATED	1	
	W-Nabellen offeres	0	111		W-Nabellen offeres	1	111	O-SENT	1	
	O-CANCELLED	0	111		O-CANCELLED	1	111	W-Completeren aanvraag	1	
	A-CANCELLED	0	111		W-Nabellen offeres	1	111	O-SENT-BACK	11259	
	W-Nabellen offeres	0	111		O-ACCEPTED	1	629	W-Nabellen offeres	11259	
								O-ACCEPTED	9	

Table 5.4: A comparison between the CBI-RWS-OPC-SBM-FSR and D4EL

In this section we show how both models (*CBI-ES-UC3-SBM-RR* and *CBI-RWS-OPC-SBM-FSR*), that the models are capable of changing the outcome of the factual. Both models also return reasonable counterfactuals. However, *CBI-ES-UC3-SBM-RR* appears to be more consistent with the counterpart of [26]. Especially in terms of the activity sequence. For instance, both, our counterfactual and the D4EL counterfactual recognize that after O-SENT, there has to be an O-SENT-BACK that eventually leads to an acceptance of the counterfactual. Both evolutionary algorithms also manage to start the process with the correct sequence of A-SUBMITTED, A-PARTLYSUBMITTED and A-PREACCEPTED. Furthermore, our model appears to be much closer in terms of sequences than the model by Hsieh, Moreira, and Ouyang. *CBI-RWS-OPC-SBM-FSR* (the model that only chooses the fittest survivors) has gaps. These gaps are an indication that the model also attempts to align towards the correct structure of the factual model. We do not see that in *CBI-ES-UC3-SBM-RR*, as it ranks feasibility above similarity and sparsity. Introducing gaps automatically reduces the feasibility of the model.

We also see, that the value for *Amount* fluctuates for the evolutionary generators. Similar, holds for the resource field. The model tends to focus on event structure first and event attributes second. This might be seen as a limiting factor when it comes to event attributes. However, one could argue that the most revealing information the counterfactuals provide for sequences are within the sequence structure and less the event attributes.

5.4.2 Analysis

Most of the results are reasonable. Surprisingly, the the models did not necessarily create counterfactuals that are much shorter than their factual counterparts. In fact, most of Hsieh, Moreira, and Ouyang's counterfactuals are shorter in length. This characteristic can be an advantage for use cases, such as medicine. The fluctuations in the loan amount was expected, as well. We did not implement any safeguard option to keep certain attributes fixed.

But the "O-SENT-BACK" steps seem to be missing in your CFs in tables 5.3 and 5.4

Nice to read this!! A similar advantage as an approach in our community known as alignment.

Did you have to manually indicate the "milestones" for the DICE4EL approach?

Before this sentence, you were discussing that both your evo4el and the dice4el obtained similar good results.

Starting here you are discuss the advantages of your approach. Therefore, I would suggest to start a new paragraph and emphasize that.

The values that were produced are more or less an indication of what the model deems as useful to change the outcome at a specific step in the process.

We are also not surprised, all models manage to capture the first few activities. These are mostly the same across all cases. If our models had not recognized these its usefulness could be questioned.

All models successfully manage to flip the outcome of the prediction model and surprisingly close to the factual compared to the model by Hsieh, Moreira, and Ouyang. However, we have to keep in mind that these observations tell us more about the model rather than the true process. More specifically, our model is capable of showing, which events and attributes have to be present at a specific point within the process.

All in all, we claim that the generator model can teach us more about the model primarily. Further improvement might show even more nuance in the models behaviour. We discuss some of them in the discussion chapter.

Much better than the
first version!!!

Chapter 6

Discussion

In this chapter, we are going to reexamine many of the past decisions we made. We critically assess the results of experiments and how we interpret them. We also propose possible improvements and opportunities for future research.

6.1 Interpretation of Results

Our first two experiments showed, that we can optimize towards viability successfully. We defined **certain criteria** and showed that a model which optimizes towards those criteria can return superior results. Furthermore, we created models that are capable of optimizing complicated operationalisations of these criteria without the limitation of a function, that has a clearly defined gradient. ~~Therefore our results run contrary to the recent optimization paradigm which focuses on gradient based optimizations.~~

~~Furthermore,~~ we highlight how it is possible to modify the counterfactual generation based on **the decision criterion someone uses to optimize them**. Specifically, the model that selected iteration survivors based on a specifically sorted ranking created more feasible results. Those results reflected patterns within our log far more than the model that exclusively focused on improving the viability measure. In contrast, this model showed that structure can play a key role in understanding why a counterfactual might change the outcome of a process.

Based on the results, we have seen towards the latter experiments, we can confidently say that the ~~model~~ is capable of generating viable counterfactuals. Infact, compared to other methods in the literature we show that our counterfactuals attempt to be closer to the factual we desire to understand. We have to note that these counterfactuals are mostly a reflection of

Try to give a short map to help your readers structure the discussion. For example,

In the following, we discuss the results in four aspects: (1) the quality in terms of viability of the counterfactual sequences generated by our models, (2) their quality compared to two baseline approaches and the state-of-the-art DICE4EL approach, and (3)....

I don't see the reasoning behind this. Where in the results did you build the arguments for this point? Could you refer back to the sections or figures/tables in the results?

Be specific now since the reader should know your approach and results.

Suggestion:
We defined the viability comprising four measures (similarity, sparsity, feasibility, delta-in-likelihood) and showed that our evolutionary approach (ECF4EL or Evo4EL) which optimizes towards).

the underlying prediction model. One might argue that this does not translate to a real world scenario. However, a model never truly does. If our framework attempts to explain, how a prediction model behaves, then its applicability to real world scenarios is depends on the viability of the model itself. But regardless of the prediction model’s performance, we can clearly gain an understanding about its internal reasoning pattern.

The viability measure we proposed shows that structural difference can help us to better understand when and where we have to apply counterfactual changes. Other approaches often seem to overlook the importance of the sequence structure. However, the *CBI-RWS-OPC-SBM-FSR* model shows, that it may be reasonable to incorporate structural differences in our viability measures. Especially, if we talk about sequences and processes. The gaps within the counterfactuals that were produced are a clear indication of that. If a model attempts to align sequences, it becomes much easier to compare them side-by-side.

In contrast to the closest alternative approach by Hsieh, Moreira, and Ouyang, we show that we can create these counterfactuals without incorporating domain specific knowledge such as an understanding of milestone patterns. Obviously, domain knowledge can always help us create better solutions. However, we do not always have access to them. We believe, that showing it is possible to create viable counterfactuals without domain specific knowledge is our greatest contribution. Furthermore, our models are capable of generating solutions that are not currently present within the data. This is often an oversight for case-based solutions, as they obviously are heavily biased towards the data input. Second they can also fail to deliver the necessary structural nuance when it comes to understanding sequences.

6.2 Limitations

There were also a number of limitations to our approach. We begin with the most obvious flaw. The generation of counterfactuals is always hard to gauge, when it comes to their usefulness. There’s not standardized way to evaluate the viability of a counterfactual. In fact, this is still an open research question. Therefore, we often have to evaluate the counterfactuals in some subjective and qualitative way. In this thesis, we decided to compare the counterfactuals with another approach in the literature and the factual itself. Because our counterfactuals did not produce nonsensical results, we deemed them viable. A domain expert might strongly disagree. Therefore, we advice to also incorporate experts in the evaluation of such an approach. This is a clear limitation of our approach and we have to acknowledge it.

Next, we introduced a novel way to measure the viability of a multivariate sequence. However, we did not compare its result to other approaches in the literature. Mostly, because very few researchers have touched upon this topic. This lack of good multivariate sequence distances is something that needs to be explored further. However, our viability measure, does introduce new ideas to this sphere of research. Mainly, the idea of incorporating structure. We believe that this might benefit disciplines such as *Business Process Mining* the most.

The viability components we chose, showed, they were capable of leading to an optimized solution, but there are most likely better ways to operationalize viability criterions. However, what makes a good counterfactual and how can we quantify that is still a subject of debate. Many researchers fall back to defining their own evaluation methods. However, we believe that a good approach is a direct and qualitative comparison between two different approaches.

Furthermore, we did not take diversity into account. Our models strictly optimize towards the optimization goal. However, as we discussed, diversity can also help us understand factuals better.

When it comes to the evolutionary algorithm, we have to admit, that there are most likely more advanced and more efficient algorithms that utilize the notion of evolution. Our approach mainly followed the basic structure of an evolutionary algorithm. However, there are methods such as CMA-ES, that are capable of improving the efficiency of the evolutionary generation.

6.3 Improvements

There are a number of improvements we propose. First, the feasibility metric compared to the other metrics often appeared far lower. The small probabilities, we saw are emblematic of the probabilistic sphere. However, it would certainly help to find ways to operationalize feasibility and make it comparable to other viability components. Our ranking-based method showed, it is possible to overcome this issue, but a less opinionated solution would be more beneficial.

Furthermore, we would like to stress that our approach is only as good as the prediction model it attempts to explain. To gain further insights into *true* process models one would have to make sure that the prediction model is accurately reflecting the real world. Again, domain expert might help to deduce, which model is the best reflection of realistic phenomena.

6.4 Future Work

With regards to future directions, it is worth pointing out, whether it is beneficial to employ other components of the viability structure. The measure described here clearly operationalized a set of criterions. However, there may be more aspects to also consider and generate even better counterfactuals. A good example would be diversity. In terms of other evolutionary approaches, it would be interesting to apply modern state-of-the-art methods, with the same viability measure.

Chapter 7

Conclusion

As discussed, give a quick summary, and discuss your answers towards each research question:)

As a conclusion, we can say that the generation of counterfactual multivariate sequences is possible. We have shown that in our experiments. The viability of the results could be contested by domain experts. However, we believe that they primarily explain the model we attempt to understand. Therefore, they are a valid and transparent reflection of a particular model. Furthermore, we show it is worth pursuing more research and insights into the generation of processes. Many examples within this thesis showed that processes are a ubiquitous part of our life. Many things can be understood as a process. Hence, shying away from complicated problems like multivariate sequence problems heavily limits our progress and understanding about cause and effect relations within our daily lives.

And very very important, following the open science principle, Do not forget to add a link to your code and your raw results.

It would also be nice to have a tutorial/example on how to use your code to generate cfs, but this could be after the defense.

Appendices

Appendix A

Counterfactual Results

Factual Seq. Amount	Activity	Outcome	Resource	Our CF Seq. Amount	Activity	Outcome	Resource	DICE4EL CF Seq. Activity	Resource	Amount
150	A-SUBMITTED	1	112							
150	A-PARTLYSUBMITTED	1	112							
150	A-PREACCEPTED	1	112							
150	W-Completeren aanvraag	1	111							
150	W-Completeren aanvraag	1	111	15423	A-SUBMITTED	0	112			
150	A-ACCEPTED	1	111	15519	A-PARTLYSUBMITTED	0	112			
150	A-FINALIZED	1	111	109	A-PREACCEPTED	0	112			
150	O-SELECTED	1	111	154	A-ACCEPTED	0	972			
150	O-CREATED	1	111	161	A-FINALIZED	0	other			
150	O-SENT	1	111	15274	O-SELECTED	0	912			
150	W-Completeren aanvraag	1	111	15293	O-CREATED	0	111			
150	O-SENT-BACK	1	149	15973	O-SENT	0	101			
150	W-Nabellen offertes	1	149	14964	W-Completeren aanvraag	0	789	A-SUBMITTED	112	171
150	O-ACCEPTED	1	629	14487	O-SENT-BACK	0	149	A-PARTLYSUBMITTED	112	171
150	A-APPROVED	1	629					A-PREACCEPTED	881	171
150	A-REGISTERED	1	629	153	W-Nabellen offertes	0	899	W-Afhandelen leads	881	171
150	A-ACTIVATED	1	629					W-Completeren aanvraag	881	171
150	W-Valideren aanvraag	1	629	15832	W-Valideren aanvraag	0	899	W-Completeren aanvraag	881	171
								W-Completeren aanvraag	11119	171

Table A.1: A comparison between the CBI-ES-UC3-SBM-RR and D4EL

Factual Seq. Amount	Activity	Outcome	Resource	Our CF Seq. Amount	Activity	Outcome	Resource	DICE4EL CF Seq. Activity	Resource	Amount
150	A-SUBMITTED	1	112							
150	A-PARTLYSUBMITTED	1	112							
150	A-PREACCEPTED	1	112							
150	W-Completeren aanvraag	1	111							
150	W-Completeren aanvraag	1	111	1	A-SUBMITTED	0	112			
150	A-ACCEPTED	1	111	1	A-PARTLYSUBMITTED	0	112			
150	A-FINALIZED	1	111	1	A-PREACCEPTED	0	112			
150	O-SELECTED	1	111	1	W-Completeren aanvraag	0	929			
150	O-CREATED	1	111	1	W-Completeren aanvraag	0	932			
150	O-SENT	1	111	1	A-ACCEPTED	0	111			
150	W-Completeren aanvraag	1	111	1	A-FINALIZED	0	111			
150	O-SENT-BACK	1	149	1	O-SELECTED	0	111			
150	W-Nabellen offertes	1	149	1	O-CREATED	0	111	A-SUBMITTED	112	171
150	O-ACCEPTED	1	629	1	O-SENT	0	111	A-PARTLYSUBMITTED	112	171
150	A-APPROVED	1	629	1	W-Nabellen offertes	0	11259	A-PREACCEPTED	881	171
150	A-REGISTERED	1	629	1	A-DECLINED	0	138	W-Afhandelen leads	881	171
150	A-ACTIVATED	1	629					W-Completeren aanvraag	881	171
150	W-Valideren aanvraag	1	629	1	W-Valideren aanvraag	0	138	W-Completeren aanvraag	881	171
								W-Completeren aanvraag	11119	171

Table A.2: A comparison between the CBI-RWS-OPC-SBM-FSR and D4EL

Table A.3: A comparison between the CBI-ES-UC3-SBM-RR and D4EL

Table A.4: A comparison between the CBI-RWS-OPC-SBM-FSR and D4E

Table A.5: A comparison between the CBI-ES-UC3-SBM-RR and D4EL

Table A.6: A comparison between the CBI-RWS-OPC-SBM-FSR and D4E

Bibliography

- [1] K. Agrawal, “To study the phenomenon of the Moravec’s Paradox,” *ArXiv*, 2010.
- [2] A. Anastasiou, P. Hatzopoulos, A. Karagrigoriou, and G. Mavridoglou, “Causality Distance Measures for Multivariate Time Series with Applications,” *Mathematics*, vol. 9, no. 21, p. 2708, 21 Oct. 25, 2021, ISSN: 2227-7390. DOI: 10.3390/math9212708. [Online]. Available: <https://www.mdpi.com/2227-7390/9/21/2708> (visited on 02/03/2022).
- [3] E. Anderson, “The Species Problem in Iris,” *Annals of the Missouri Botanical Garden*, vol. 23, no. 3, pp. 457–509, 1936, ISSN: 0026-6493. DOI: 10.2307/2394164. JSTOR: 2394164.
- [4] A. Apostolico and R. Giancarlo, “Sequence Alignment in Molecular Biology,” *Journal of Computational Biology*, vol. 5, no. 2, pp. 173–196, Jan. 1998. DOI: 10.1089/cmb.1998.5.173. [Online]. Available: <https://www-liebertpub-com.proxy.library.uu.nl/doi/10.1089/cmb.1998.5.173> (visited on 04/22/2022).
- [5] E. Ates, B. Aksar, V. J. Leung, and A. K. Coskun, “Counterfactual Explanations for Multivariate Time Series,” in *2021 International Conference on Applied Artificial Intelligence (ICAPAI)*, Halden, Norway: IEEE, May 19, 2021, pp. 1–8, ISBN: 978-1-72815-934-8. DOI: 10.1109/ICAPAI49758.2021.9462056. [Online]. Available: <https://ieeexplore.ieee.org/document/9462056/> (visited on 03/01/2022).
- [6] J. Baker, J. Song, and D. R. Jones, “Closing the Loop: An Empirical Investigation of Causality in IT Business Value,” *undefined*, 2017. [Online]. Available: <https://www.semanticscholar.org/paper/Closing-the-Loop%3A-An-Empirical-Investigation-of-in-Baker-Song/df210060211bdc598f2d3382c68c615319287f71> (visited on 03/01/2022).

- [7] A. D. Bautista, L. Wangikar, and S. Akbar, “Process Mining-Driven Optimization of a Consumer Loan Approvals Process - The BPIC 2012 Challenge Case Study,” in *Business Process Management Workshops*, 2012. DOI: 10.1007/978-3-642-36285-9_24.
- [8] S. Bond-Taylor, A. Leach, Y. Long, and C. G. Willcocks, “Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models,” Apr. 14, 2021, [Online]. Available: <http://arxiv.org/abs/2103.04922> (visited on 10/01/2021).
- [9] D. V. Carvalho, E. M. Pereira, and J. S. Cardoso, “Machine Learning Interpretability: A Survey on Methods and Metrics,” *Electronics*, vol. 8, no. 8, p. 832, 8 Aug. 2019. DOI: 10.3390/electronics8080832. [Online]. Available: <https://www.mdpi.com/2079-9292/8/8/832> (visited on 11/09/2021).
- [10] J. Chung, K. Kastner, L. Dinh, K. Goel, A. Courville, and Y. Bengio, “A Recurrent Latent Variable Model for Sequential Data,” in *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, ser. NIPS’15, Cambridge, MA, USA: MIT Press, Apr. 6, 2016, pp. 2980–2988. [Online]. Available: <http://arxiv.org/abs/1506.02216> (visited on 02/03/2022).
- [11] (). “Counterfactual,” Oxford Reference, [Online]. Available: <https://www-oxfordreference-com.proxy.library.uu.nl/view/10.1093/oi/authority.20110803095642948> (visited on 02/10/2022).
- [12] F. J. Damerau, “A technique for computer detection and correction of spelling errors,” *Communications of the ACM*, vol. 7, no. 3, pp. 171–176, Mar. 1, 1964, ISSN: 0001-0782. DOI: 10.1145/363958.363994. [Online]. Available: <https://doi.org/10.1145/363958.363994> (visited on 04/15/2022).
- [13] S. Dandl, C. Molnar, M. Binder, and B. Bischl, “Multi-Objective Counterfactual Explanations,” in *Parallel Problem Solving from Nature – PPSN XVI*, T. Bäck, M. Preuss, A. Deutz, H. Wang, C. Doerr, M. Emmerich, and H. Trautmann, Eds., ser. Lecture Notes in Computer Science, Cham: Springer International Publishing, 2020, pp. 448–469, ISBN: 978-3-030-58112-1. DOI: 10.1007/978-3-030-58112-1_31.
- [14] (). “Definition of PROCESS,” [Online]. Available: <https://www.merriam-webster.com/dictionary/process> (visited on 02/17/2022).

- [15] E. Delaney, D. Greene, and M. T. Keane, “Instance-Based Counterfactual Explanations for Time Series Classification,” in *Case-Based Reasoning Research and Development*, A. A. Sánchez-Ruiz and M. W. Floyd, Eds., ser. Lecture Notes in Computer Science, Cham: Springer International Publishing, 2021, pp. 32–47, ISBN: 978-3-030-86957-1. DOI: 10.1007/978-3-030-86957-1_3.
- [16] L. Deng, “The MNIST Database of Handwritten Digit Images for Machine Learning Research [Best of the Web],” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, Nov. 2012, ISSN: 1558-0792. DOI: 10.1109/MSP.2012.2211477.
- [17] R. A. Fisher, “The Use of Multiple Measurements in Taxonomic Problems,” *Annals of Eugenics*, vol. 7, no. 2, pp. 179–188, 1936, ISSN: 2050-1439. DOI: 10.1111/j.1469-1809.1936.tb02137.x. [Online]. Available: <http://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-1809.1936.tb02137.x> (visited on 04/21/2022).
- [18] M. Fraccaro, S. K. Sønderby, U. Paquet, and O. Winther, “Sequential neural models with stochastic layers,” in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, ser. NIPS’16, Red Hook, NY, USA: Curran Associates Inc., Dec. 5, 2016, pp. 2207–2215, ISBN: 978-1-5108-3881-9.
- [19] W. N. Francis and H. Kucera, “Brown corpus manual,” Department of Linguistics, Brown University, Providence, Rhode Island, US, 1979. [Online]. Available: <http://icame.uib.no/brown/bcm.html>.
- [20] S. L. Frank and M. H. Christiansen, “Hierarchical and sequential processing of language,” *Language, Cognition and Neuroscience*, vol. 33, no. 9, pp. 1213–1218, Oct. 21, 2018, ISSN: 2327-3798. DOI: 10.1080/23273798.2018.1424347. [Online]. Available: <https://doi.org/10.1080/23273798.2018.1424347> (visited on 04/22/2022).
- [21] L. Girin, S. Leglaive, X. Bie, J. Diard, T. Hueber, and X. Alameda-Pineda, “Dynamical Variational Autoencoders: A Comprehensive Review,” *Foundations and Trends® in Machine Learning*, vol. 15, no. 1-2, pp. 1–175, Dec. 1, 2021, ISSN: 1935-8237, 1935-8245. DOI: 10.1561/22000000089. arXiv: 2008.12595. [Online]. Available: <http://arxiv.org/abs/2008.12595> (visited on 03/02/2022).
- [22] K. M. Hangos, G. Szederkényi, R. Lakner, and M. Gerzson, *Intelligent Control Systems: An Introduction with Examples*. Springer Science & Business Media, 2001, 332 pp., ISBN: 978-1-4020-0134-5.

- [23] C. Hitchcock, “Causal Models,” in *The Stanford Encyclopedia of Philosophy*, E. N. Zalta, Ed., Summer 2020, Metaphysics Research Lab, Stanford University, 2020. [Online]. Available: <https://plato.stanford.edu/archives/sum2020/entries/causal-models/> (visited on 02/10/2022).
- [24] S. Hochreiter and J. Schmidhuber, “Long Short-Term Memory,” *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, Nov. 1, 1997, issn: 0899-7667. DOI: 10.1162/neco.1997.9.8.1735. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735> (visited on 03/27/2022).
- [25] B. F. A. Hompes, A. Maaradji, M. La Rosa, M. Dumas, J. C. A. M. Buijs, and W. M. P. van der Aalst, “Discovering Causal Factors Explaining Business Process Performance Variation,” in *Advanced Information Systems Engineering*, E. Dubois and K. Pohl, Eds., ser. Lecture Notes in Computer Science, Cham: Springer International Publishing, 2017, pp. 177–192, ISBN: 978-3-319-59536-8. DOI: 10.1007/978-3-319-59536-8_12.
- [26] C. Hsieh, C. Moreira, and C. Ouyang, “DiCE4EL: Interpreting Process Predictions using a Milestone-Aware Counterfactual Approach,” in *2021 3rd International Conference on Process Mining (ICPM)*, Eindhoven, Netherlands: IEEE, Oct. 31, 2021, pp. 88–95, ISBN: 978-1-66543-514-7. DOI: 10.1109/ICPM53251.2021.9576881. [Online]. Available: <https://ieeexplore.ieee.org/document/9576881/> (visited on 11/04/2021).
- [27] R. E. Kalman, “A New Approach to Linear Filtering and Prediction Problems,” *Transactions of the ASME–Journal of Basic Engineering*, vol. 82, pp. 35–45, Series D 1960.
- [28] J. Klimek, J. Klimek, W. Kraskiewicz, and M. Topolewski, “Long-term series forecasting with Query Selector – efficient model of sparse attention,” Aug. 17, 2021, [Online]. Available: <http://arxiv.org/abs/2107.08687> (visited on 11/09/2021).
- [29] J. Krause, A. Perer, and K. Ng, “Interacting with Predictions: Visual Inspection of Black-box Machine Learning Models,” in *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, ser. CHI ’16, New York, NY, USA: Association for Computing Machinery, May 7, 2016, pp. 5686–5697, ISBN: 978-1-4503-3362-7. DOI: 10.1145/2858036.2858529. [Online]. Available: <https://doi.org/10.1145/2858036.2858529> (visited on 02/26/2022).

- [30] R. Krishnan, U. Shalit, and D. Sontag, "Structured Inference Networks for Nonlinear State Space Models," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, no. 1, 1 Feb. 13, 2017, ISSN: 2374-3468. [Online]. Available: <https://ojs.aaai.org/index.php/AAAI/article/view/10779> (visited on 02/22/2022).
- [31] A. Lambora, K. Gupta, and K. Chopra, "Genetic Algorithm- A Literature Review," in *2019 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing (COMITCon)*, Feb. 2019, pp. 380–384. DOI: 10.1109/COMITCon.2019.8862255.
- [32] S. Leglaive, X. Alameda-Pineda, L. Girin, and R. Horaud, "A Recurrent Variational Autoencoder for Speech Enhancement," Feb. 10, 2020, [Online]. Available: <http://arxiv.org/abs/1910.10942> (visited on 02/07/2022).
- [33] V. Levenshtein, "Binary codes capable of correcting deletions, insertions, and reversals," *undefined*, 1965. [Online]. Available: <https://www.semanticscholar.org/paper/Binary-codes-capable-of-correcting-deletions%2C-and-Levenshtein/b2f8876482c97e804bb50a5e2433881ae31> (visited on 04/15/2022).
- [34] F. Mannhardt and D. Blinde, "Analyzing the trajectories of patients with sepsis using process mining: RADAR + EMISA 2017," *RADAR+EMISA 2017, Essen, Germany, June 12-13, 2017*, CEUR Workshop Proceedings, pp. 72–80, 2017. [Online]. Available: <http://www.scopus.com/inward/record.url?scp=85022001209&partnerID=8YFLogxK> (visited on 04/22/2022).
- [35] M. P. Marcus, M. A. Marcinkiewicz, and B. Santorini, "Building a large annotated corpus of English: The penn treebank," *Computational Linguistics*, vol. 19, no. 2, pp. 313–330, Jun. 1, 1993, ISSN: 0891-2017.
- [36] D. Martens and F. Provost, "Explaining data-driven document classifications," *MIS Quarterly*, vol. 38, no. 1, pp. 73–100, Mar. 1, 2014, ISSN: 0276-7783. DOI: 10.25300/MISQ/2014/38.1.04. [Online]. Available: <https://doi.org/10.25300/MISQ/2014/38.1.04> (visited on 02/26/2022).
- [37] I. Melnyk, C. N. dos Santos, K. Wadhawan, I. Padhi, and A. Kumar, "Improved Neural Text Attribute Transfer with Non-parallel Data," Dec. 4, 2017, [Online]. Available: <http://arxiv.org/abs/1711.09395> (visited on 02/28/2022).

- [38] R. Mitton, “Fifty years of spellchecking,” *Writing Systems Research*, vol. 2, no. 1, pp. 1–7, 2010. DOI: 10.1093/wsr/wsq004. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-78649452350&doi=10.1093%2fwsr%2fwsq004&partnerID=40&md5=5b9a37202101a18bb4b78b7cdeb52c>.
- [39] C. Molnar, *Interpretable Machine Learning. A Guide for Making Black Box Models Explainable*. 2019. [Online]. Available: <https://christophm.github.io/interpretable-ml-book/>.
- [40] R. K. Mothilal, A. Sharma, and C. Tan, “Explaining Machine Learning Classifiers through Diverse Counterfactual Explanations,” in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, Jan. 27, 2020, pp. 607–617. DOI: 10.1145/3351095.3372850. arXiv: 1905.07697 [cs, stat]. [Online]. Available: <http://arxiv.org/abs/1905.07697> (visited on 07/25/2022).
- [41] T. Narendra, P. Agarwal, M. Gupta, and S. Dechu, “Counterfactual Reasoning for Process Optimization Using Structural Causal Models,” in *Business Process Management Forum*, T. Hildebrandt, B. F. van Dongen, M. Röglinger, and J. Mendling, Eds., ser. Lecture Notes in Business Information Processing, Cham: Springer International Publishing, 2019, pp. 91–106, ISBN: 978-3-030-26643-1. DOI: 10.1007/978-3-030-26643-1_6.
- [42] M. Oberst and D. Sontag, “Counterfactual Off-Policy Evaluation with Gumbel-Max Structural Causal Models,” Jun. 6, 2019, [Online]. Available: <http://arxiv.org/abs/1905.05824> (visited on 09/22/2021).
- [43] J. Pearl, M. Glymour, and N. P. Jewell, *Causal Inference in Statistics: A Primer*. Chichester, West Sussex: Wiley, 2016, 136 pp., ISBN: 978-1-119-18684-7.
- [44] M. S. Qafari and W. M. P. van der Aalst, “Case Level Counterfactual Reasoning in Process Mining,” in *Intelligent Information Systems*, S. Nurcan and A. Korthaus, Eds., ser. Lecture Notes in Business Information Processing, Cham: Springer International Publishing, 2021, pp. 55–63, ISBN: 978-3-030-79108-7. DOI: 10.1007/978-3-030-79108-7_7.
- [45] M. Robeer, F. Bex, and A. Feelders, “Generating Realistic Natural Language Counterfactuals,” in *Findings of the Association for Computational Linguistics: EMNLP 2021*, Punta Cana, Dominican Republic: Association for Computational Linguistics, Nov. 2021, pp. 3611–3625. DOI: 10.18653/v1/2021.findings-emnlp.306. [Online]. Available:

- <https://aclanthology.org/2021.findings-emnlp.306> (visited on 02/26/2022).
- [46] C. L. Shook, D. J. Ketchen Jr., G. T. M. Hult, and K. M. Kacmar, “An assessment of the use of structural equation modeling in strategic management research,” *Strategic Management Journal*, vol. 25, no. 4, pp. 397–404, 2004, ISSN: 1097-0266. DOI: 10.1002/smj.385. [Online]. Available: <http://onlinelibrary.wiley.com/doi/abs/10.1002/smj.385> (visited on 03/01/2022).
 - [47] W. Starr, “Counterfactuals,” in *The Stanford Encyclopedia of Philosophy*, E. N. Zalta, Ed., Summer 2021, Metaphysics Research Lab, Stanford University, 2021. [Online]. Available: <https://plato.stanford.edu/archives/sum2021/entries/counterfactuals/> (visited on 02/09/2022).
 - [48] N. Tax, I. Verenich, M. La Rosa, and M. Dumas, “Predictive Business Process Monitoring with LSTM Neural Networks,” in *Advanced Information Systems Engineering*, E. Dubois and K. Pohl, Eds., ser. Lecture Notes in Computer Science, Cham: Springer International Publishing, 2017, pp. 477–492, ISBN: 978-3-319-59536-8. DOI: 10.1007/978-3-319-59536-8_30.
 - [49] I. Teinemaa, M. Dumas, M. La Rosa, and F. M. Maggi, “Outcome-Oriented Predictive Process Monitoring: Review and Benchmark,” Oct. 23, 2018, [Online]. Available: <http://arxiv.org/abs/1707.06766> (visited on 05/07/2022).
 - [50] S. Tsirtsis, A. De, and M. Gomez-Rodriguez, “Counterfactual Explanations in Sequential Decision Making Under Uncertainty,” Jul. 6, 2021, [Online]. Available: <http://arxiv.org/abs/2107.02776> (visited on 09/09/2021).
 - [51] W. van der Aalst, A. Adriansyah, A. K. A. de Medeiros, F. Arcieri, T. Baier, T. Blickle, J. C. Bose, P. van den Brand, R. Brandtjen, J. Buijs, A. Burattin, J. Carmona, M. Castellanos, J. Claes, J. Cook, N. Costantini, F. Curbera, E. Damiani, M. de Leoni, P. Delias, B. F. van Dongen, M. Dumas, S. Dustdar, D. Fahland, D. R. Ferreira, W. Gaaloul, F. van Geffen, S. Goel, C. Günther, A. Guzzo, P. Harmon, A. ter Hofstede, J. Hoogland, J. E. Ingvaldsen, K. Kato, R. Kuhn, A. Kumar, M. La Rosa, F. Maggi, D. Malerba, R. S. Mans, A. Manuel, M. McCreesh, P. Mello, J. Mendling, M. Montali, H. R. Motahari-Nezhad, M. zur Muehlen, J. Munoz-Gama, L. Pontieri, J. Ribeiro, A. Rozinat, H. Seguel Pérez, R. Seguel Pérez, M. Sepúlveda, J. Sinur, P. Soffer, M. Song, A. Sperduti, G. Stilo, C. Stoel, K. Swenson, M. Talamo, W.

- Tan, C. Turner, J. Vanthienen, G. Varvaressos, E. Verbeek, M. Verdonk, R. Vigo, J. Wang, B. Weber, M. Weidlich, T. Weijters, L. Wen, M. Westergaard, and M. Wynn, “Process Mining Manifesto,” in *Business Process Management Workshops*, F. Daniel, K. Barkaoui, and S. Dustdar, Eds., ser. Lecture Notes in Business Information Processing, Berlin, Heidelberg: Springer, 2012, pp. 169–194, ISBN: 978-3-642-28108-2. DOI: 10.1007/978-3-642-28108-2_19.
- [52] P. A. Vikhar, “Evolutionary algorithms: A critical review and its future prospects,” in *2016 International Conference on Global Trends in Signal Processing, Information Computing and Communication (ICGT-SPICC)*, Dec. 2016, pp. 261–265. DOI: 10.1109/ICGTSPICC.2016.7955308.
- [53] S. Wachter, B. Mittelstadt, and C. Russell, “Counterfactual Explanations Without Opening the Black Box: Automated Decisions and the GDPR,” *ArXiv*, 2017. DOI: 10.2139/ssrn.3063289.
- [54] J. Wang, S. Song, X. Zhu, and X. Lin, “Efficient recovery of missing events,” *Proceedings of the VLDB Endowment*, vol. 6, no. 10, pp. 841–852, Aug. 1, 2013, ISSN: 2150-8097. DOI: 10.14778/2536206.2536212. [Online]. Available: <https://doi.org/10.14778/2536206.2536212> (visited on 04/21/2022).
- [55] K. Wang, H. Hua, and X. Wan, “Controllable Unsupervised Text Attribute Transfer via Editing Entangled Latent Representation,” Dec. 12, 2019, [Online]. Available: <http://arxiv.org/abs/1905.12926> (visited on 02/28/2022).
- [56] Z. Wang, J. Zhang, H. Xu, X. Chen, Y. Zhang, W. X. Zhao, and J.-R. Wen, “Counterfactual Data-Augmented Sequential Recommendation,” in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, ser. SIGIR ’21, New York, NY, USA: Association for Computing Machinery, Jul. 11, 2021, pp. 347–356, ISBN: 978-1-4503-8037-9. DOI: 10.1145/3404835.3462855. [Online]. Available: <https://doi.org/10.1145/3404835.3462855> (visited on 09/09/2021).