

In this section we define the components of the metric to measure the validity of a sequence.

0.0.1 Semi-Structured Damerau-Levenshtein Distance

Before discussing some of the similarity measures it is important to briefly introduce the Damerau-Levenshtein distance. This distance function is a modified version of the Levenshtein distance[3], which is a widely used to compute the edit-distance of two discrete sequences[1, 4]. The most important applications are within the Natural Language Processing (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 1, we schematically show two sequences and their distance.

In Equation 1 depicts the recursive formulation of the distance. The distance computes the costs of transforming the sequence a to b , by computing the minimum of five separate terms.

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

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. For Process Mining, it becomes necessary to modify the distance further.

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 zero, as we can recognise one substitution at the second character in both strings. However, this

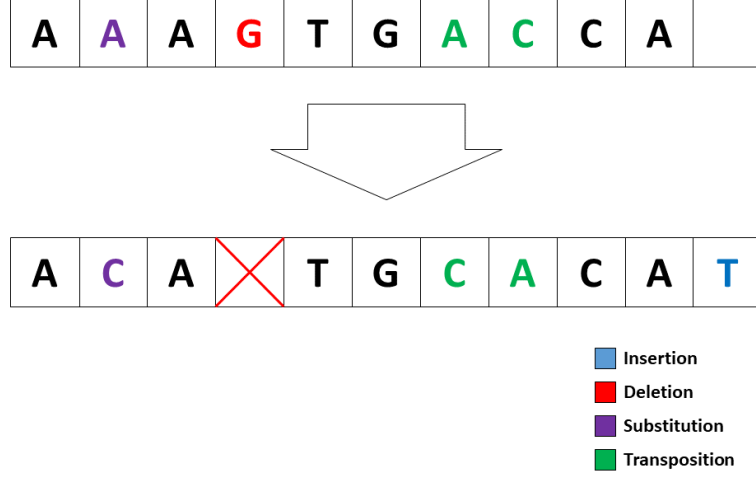


Figure 1: Shows two sequences. The edit distance is the sum of multiple operations. Blue shows an insert, red a deletion, purple a substitution and green a transposition. Therefore the edit distance is 4.

representation is insufficient for process instances as they may also contain attribute values. Therefore, we characterise the sequences as process events in Equation 2.

$$s^1 = \{a, a, b, a\} \quad (2)$$

$$s^2 = \{a, a^*, b, a\} \quad (3)$$

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

$$s^4 = \{a, b, a\} \quad a, b, c \in \mathbb{R}^3 \quad (5)$$

$$a = \begin{bmatrix} 2 \\ 1 \\ 4 \end{bmatrix} \quad a^* = \begin{bmatrix} 3 \\ 3 \\ 4 \end{bmatrix} \quad b = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad c = \begin{bmatrix} 3 \\ 1 \\ 4 \end{bmatrix} \quad (6)$$

If we do not consider attribute values, it becomes clear that s^2 , s^3 and s^4 have an edit distance of 0, 1 and 1 to s^1 . However, with attribute values

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. These examples show that we can neither disregard attribute values nor events, while computing the edit distance of two process instances. We show this in ???. In other words, we cannot simply assume a static cost of 1 for each edit operation. Instead, we have to define a cost function which takes attribute variables into account. In the following sections, we will establish distances which use a modified Damerau-Levenshtein distance approach. Here, the cost of each edit-operation will be weighted with a distance function that considers the difference between event attributes. In simplified terms, we can say that s^1 and s^2 are identical, if we only consider the activity. However, taking attribute values into account, s^1 and s^2 actually differ on two accounts. 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 7. For the remainder, we will denote this edit-distance as Semi-structured Damerau-Levenshtein distance (SSDLD).

$$d_{a,b}(i, j) = \min \begin{cases} d_{a,b}(i-1, j) + cost(\mathbf{0}, b_j) & \text{if } i > 0 \\ d_{a,b}(i, j-1) + cost(a_i, \mathbf{0}) & \text{if } j > 0 \\ d_{a,b}(i-1, j-1) + cost(a_i, b_j) & \text{if } i, j > 0 \\ & \& \bar{a}_i = \bar{b}_j \\ d_{a,b}(i-1, j-1) + cost(a_i, \mathbf{0}) + cost(\mathbf{0}, b_j) & \text{if } i, j > 0 \\ & \& \bar{a}_i \neq \bar{b}_j \\ d_{a,b}(i-2, j-2) + cost(a_i, b_{j-1}) + cost(a_{i-1}, b_j) & \text{if } i, j > 1 \\ & \& \bar{a}_i = \bar{b}_{j-1} \\ & \& \bar{a}_{i-1} = \bar{b}_j \\ 0 & \text{if } i = j = 0 \end{cases} \quad (7)$$

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. $cost(a, b)$ is a cost function which takes the attribute values of a and b into account.

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*.

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.

0.0.2 Discussing the Semi-Structured DL Distance

There are two important characteristics, which require a discussion.

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

Second, to explain the intuition behind most terms, it is necessary to establish a common understanding between the relationship of an event with its event attributes. 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 will 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 cost functions, 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 $(\begin{smallmatrix} 2 \\ 3 \end{smallmatrix})$ and $(\begin{smallmatrix} 2 \\ 3 \end{smallmatrix})$, these events may be seen as more similar than two \bar{a} and \bar{a} with attribute values $(\begin{smallmatrix} 2 \\ 3 \end{smallmatrix})$ and $(\begin{smallmatrix} 5 \\ 0 \end{smallmatrix})$.

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 an event at daytime. 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. Second, the unmodified Damerau-Levenshtein distance applies to discrete sequences, such as textual data with atomic words or characters. By treating each event as a discrete sequence element, we remain faithful to the original function.

[Change to a more consistent order of introducing the measures.]

0.0.3 Delta-Measure

For this measure, we can evaluate the precision of a counterfactual trace by determining whether a counterfactual would lead to the desired outcome. For this purpose, we use the predictive model, which will output a prediction based on the counterfactual sequence. However, it is often difficult to force a deterministic model to produce a different result. We can relax the condition by maximising the likelihood of the counterfactual outcome. If we compare the likelihood of the desired outcome under the factual sequence with the counterfactual sequence, we can determine an increase or decrease. Ideally, we want to increase the likelihood. We can compute the odds or the difference between the two likelihoods. We choose to use the delta. **[Mention that Delta is a replacement name to avoid confusion.]**

$$improvement = p(o^*|a^*) - p(o^*|a) \quad (8)$$

Here, $p(o|a)$ describes the probability of an outcome, given a sequence of events. **[Describe the more complicated view on this measure]**

0.0.4 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 events themselves. 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, we can assume that the data is representative of the process dynamics. Hence, we can simply 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. Hence, we compute the joint probability of a case by using the forward algorithm [CITE forward algorithm](#). In Equation 9 shows the formulation.

$$\begin{aligned}
feasibility_e &= p(a_i | a_{i-1} \dots a_1, \theta) \approx p(a_i | a_{i-1}, \theta) \\
feasibility_f &= p(f_i | a_i, \theta) \\
feasibility &= \prod_{i=1}^n [p(f_i | a_i, \theta) * p(a_i | a_{i-1}, \theta)]
\end{aligned} \tag{9}$$

Here, a and f are the activity and features of a particular event. Likewise, θ is the data sample which is used to determine the parameters of the density function. The first equation shows the approximation based on the markov assumption.

0.0.5 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 (subsection 0.0.1), we want to take the structural distance and the feature distance into account. Henceforth, we will use the previously established SSDLD. The similarity distance uses a cost function as specified in Equation 10.

$$\begin{aligned}
cost(a_i, b_j) &= L2(a_i, b_j) \\
a_i, b_j &\in \mathbb{R}^d
\end{aligned} \tag{10}$$

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

0.0.6 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 11.

$$\begin{aligned} cost(a_i, b_j) = \sum_d \mathbb{I}(a_{id} = b_{jd}) \\ a_i, b_j \in \mathbb{R}^d \end{aligned} \tag{11}$$

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.

0.0.7 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 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 a considerable counterfactual.]**