# EFFICIENT DISTANCE PRUNING FOR PROCESS SUFFIX COMPARISON IN PRESCRIPTIVE PROCESS MONITORING

Sarra Madad[1,2]

[1]Université de Technologie de Troyes, LIST3N Research Unit, Troyes, FRANCE
[2]QAD Process Intelligence, Courbevoie, FRANCE

## ABSTRACT

Prescriptive process monitoring seeks to recommend actions that improve process outcomes by analyzing possible continuations of ongoing cases. A key obstacle is the heavy computational cost of large-scale suffix comparisons, which grows rapidly with log size. We propose an efficient retrieval method exploiting the triangle inequality: distances to a set of optimized pivots define bounds that prune redundant comparisons. This substantially reduces runtime and is fully parallelizable. Crucially, pruning is exact: the retrieved suffixes are identical to those from exhaustive comparison, thereby preserving accuracy. These results show that metric-based pruning can accelerate suffix comparison and support scalable prescriptive systems.

## 1 INTRODUCTION : PRESCRIPTIVE PROCESS MONITORING FOR SUFFIX PREDICTION

Suffix prediction plays a central role in process analytics, as it enables estimating how an ongoing execution may unfold under different continuations. In the context of process mining, a suffix simply denotes the sequence of future activities of a case until its completion. Building on this capability, prescriptive approaches aim not only to anticipate future behavior but also to recommend interventions that improve process performance (Weinzierl et al. 2020). In this setting, the system must evaluate how different possible continuations of the current execution would affect key performance indicators (KPIs), enabling the recommendation of a next-best action. A central mechanism to achieve this is the search for contrasting suffixes: process continuations from past cases that diverge in their outcomes. By comparing these suffixes, the system can infer which actions tend to lead to favorable trajectories. This strategy entails numerous pairwise distance computations, which become prohibitive as event logs grow (Berti 2019). To mitigate this, we apply a triangle inequality–based pruning method (Jeromin and Körner 1989), where distances to a small set of pivots define bounds that allow discarding many redundant comparisons. This significantly reduces computation while preserving exactness.

## 2 TRIANGULAR INEQUALITY ACCELERATION

Let $\mathscr{S}$ be the countable set of process suffixes, $d : \mathscr{S} \times \mathscr{S} \to [0,\infty)$ a distance function (e.g., Euclidean or cosine distance), $\tau \in [0,\infty)$ a threshold, and $P = \{z_1,\dots,z_K\} \subset \mathscr{S}$ a finite pivot set chosen to cover $\mathscr{S}$ (e.g., minimizing $R(P) = \max_x \min_{z \in P} d(x,z)$). For any suffixes $x,y \in \mathscr{S}$ and a pivot $z \in \mathscr{S}$, the triangle inequality gives:

$$\big|d(x,z) - d(y,z)\big| \ \leq \ d(x,y) \ \leq \ d(x,z) + d(y,z).$$

By introducing a set of $K$ pivots $P = \{z_1,\dots,z_K\}$, $P = \{z_1,\dots,z_K\} \subset \mathscr{S}$, we can refine these bounds as follows:

$$\max_{1 \leq k \leq K} \big|d(x,z_k) - d(y,z_k)\big| \ \leq \ d(x,y) \ \leq \ \min_{1 \leq k \leq K} \big(d(x,z_k) + d(y,z_k)\big).$$

This allows pruning: if the lower bound already exceeds a search threshold $\tau$ (e.g., for $k$-nearest neighbors (Cover and Hart 1967)), then computing $d(x,y)$ explicitly is unnecessary. Similarly, if the upper bound is below $\tau$, the pair $(x,y)$ can be accepted directly.

The effectiveness of this approach depends on the choice of pivots. A common strategy is to select pivots that "cover" the space of suffixes. This can be formalized as a *k*-center problem:

$$\min_{P \subset \mathscr{S}, |P|=K} \max_{x \in \mathscr{S}} \min_{z \in P} d(x,z),$$

which seeks pivots that minimize the maximum distance of any suffix to its closest pivot. Although NP-hard, this objective can be approximated efficiently using a greedy farthest-point heuristic, which iteratively selects the suffix farthest from the already chosen pivots. This ensures that the selected pivots are well spread across the dataset, tightening the bounds and increasing pruning efficiency.

**Example.** Suppose we have two suffixes $x, y$, a threshold $\tau = 4$, and two pivots $z_1, z_2$ with precomputed distances:

$$d(x,z_1) = 2, \quad d(x,z_2) = 6, \quad d(y,z_1) = 7, \quad d(y,z_2) = 9.$$

Then the lower bound is

$$\max\{|2-7|, |6-9|\} = \max\{5,3\} = 5.$$

Since $5 > \tau$, the comparison between $x$ and $y$ can be discarded without computing $d(x,y)$.

At scale, this method relies on precomputing a distance matrix of size $|\mathscr{S}| \times K$ (suffixes $\times$ pivots), which can be reused across queries. In practice, a small number of well-chosen pivots often suffices to prune a large fraction of candidate comparisons, yielding substantial computational savings while preserving exactness.

## 3 EVALUATION

Processing the full dataset of about 150,000 suffixes originally required nearly 89 hours. With the new approach, batches of 500 suffixes take about 2.5 hours and the method is fully parallelizable. By construction, pruning is lossless: retrieved suffixes always match the baseline of exhaustive pairwise comparison, a result confirmed empirically with 100% accuracy.

## 4 CONCLUSION

This approach addressed the scalability challenge of suffix retrieval in prescriptive process monitoring. By leveraging the triangle inequality with optimized pivot selection and batching strategies, the proposed approach drastically reduces the number of distance computations while preserving accuracy. The method is fully parallelizable and therefore well-suited for large-scale event logs.

## ACKNOWLEDGMENTS

## REFERENCES

Alessandro Berti 2019. "Increasing Scalability of Process Mining using Event Dataframes: How Data Structure Matters" https://doi.org/https://doi.org/10.48550/arXiv.1907.12817.

Cover, T., and P. Hart. 1967. "Nearest neighbor pattern classification". *IEEE Transactions on Information Theory* 13(1):21–27 https://doi.org/10.1109/TIT.1967.1053964.

Jeromin, B., and F. Körner. 1989. "Triangle inequality and symmetry in connection with the assignment and the traveling salesman problem". *European Journal of Operational Research* 38(1):70–75 https://doi.org/https://doi.org/10.1016/0377-2217(89)90470-0.

Weinzierl, S., S. Dunzer, S. Zilker, and M. Matzner. 2020. *Prescriptive Business Process Monitoring for Recommending Next Best Actions*, 193–209. Springer International Publishing https://doi.org/10.1007/978-3-030-58638-6_12.