

ProcK: Machine Learning for Knowledge-Intensive Processes

Tobias Jacobs¹[0000-0002-8130-4211], Jingyi Yu^{1,2}[0000-0002-1144-9765],
Julia Gastinger¹[0000-0003-1914-6723], and Timo Sztyler¹[0000-0001-8132-5920]

¹ NEC Laboratories Europe GmbH, Heidelberg, Germany
tobias.jacobs@neclab.eu, julia.gastinger@neclab.eu, timo.sztyler@neclab.eu
² Faculty of Electrical Engineering and Information Technology, RWTH Aachen
jingyi.yu@rwth-aachen.de

Abstract. We present a novel methodology to build powerful predictive process models. Our method, denoted *ProcK* (*Process & Knowledge*), relies not only on sequential input data in the form of event logs, but can learn to use a knowledge graph to incorporate information about the attribute values of the events and their mutual relationships. The idea is realized by mapping event attributes to nodes of a knowledge graph and training a sequence model alongside a graph neural network in an end-to-end fashion. This hybrid approach substantially enhances the flexibility and applicability of predictive process monitoring, as both the static and dynamic information residing in the databases of organizations can be directly taken as input data. We demonstrate the potential of ProcK by applying it to a number of predictive process monitoring tasks, including tasks with knowledge graphs available as well as an existing process monitoring benchmark where no such graph is given. The experiments provide evidence that our methodology achieves state-of-the-art performance and improves predictive power when a knowledge graph is available.

Keywords: predictive process management · neural networks

1 Introduction

We introduce *ProcK* (*Process & Knowledge*), a pipeline for predictive process monitoring. ProcK combines the usage of two complementary data representations in a novel way.

Predictive process monitoring deals with the task of forecasting properties of business processes that are currently under execution. This includes the type and occurrence time of future events as well as the process outcome. The primary input to predictive process models are logs recorded during process execution, and it is best practice in the process mining community to model them as sets of discrete *events*. Each event is characterized by its case identifier, activity type, timestamp, and potentially further data; see Van Der Aalst *et al.* [2011a].

Machine learning methods for predictive process monitoring described in literature can be separated into two main approaches. The more traditional approach is to hand-engineer a set of feature extraction functions that operate on top of event sequences. The set of features, after some further pre-processing, is then used as input to train a machine learning model like e.g. an SVM (Leontjeva *et al.* [2016]). The second approach relies on deep learning, feeding the raw event log directly into a deep neural network which builds meaningful features automatically during training. Because of the sequential nature of the input data, it is a natural choice to apply a recurrent neural network, like done by Tax *et al.* [2017]. More recently, feedforward networks have been demonstrated by Mauro *et al.* [2019] and other authors to achieve superior performance in many cases.

Our work follows the deep learning paradigm, and, to the best of our knowledge, we are the first to complement the event sequence data model with an additional representation of the available input data as a knowledge graph. This idea is rooted in the first fundamental step of typical practical process mining projects: to extract the event log from the data lake of an organization (see Reinkemeyer [2020]), which is often structured in the form of one or more relational databases. When selecting data with the goal of building a high-quality prediction model, the limitations of the event sequence view become apparent: only a subset of the relevant data can be naturally expressed in the form of events with case identifier and timestamp.

Example 1. For predicting the success of a loan repayment process at the time when the first rate has been paid, it might be relevant to take into account the bank account from which the rate was transferred. Specifically, the economic stability of the bank account’s country might be an indicator.

In Example 1, the relevant piece of information (economic stability) is not an event, and it is neither a primary attribute of the event. It is rather an indirect attribute of the transfer event, which has to be derived via a specific semantic path (transfer \rightarrow bank \rightarrow country \rightarrow economic stability). Domain experts could re-define such derived attributes by hand as primary event information, but this counteracts the benefit of deep learning to be applicable on top of the raw data.

To address that issue, ProcK takes a knowledge graph as additional input. It stacks a sequence model for events on top of a graph neural network in order to compute meaningful event representations. In Example 1, information about the *economic stability* can be propagated backwards across the path *economic stability* \leftarrow *country* \leftarrow *bank* \leftarrow *transfer*, where *economic stability*, *bank*, *country* are knowledge graph nodes and *transfer* is a time-stamped event of a particular process. Having been integrated into the representation of *transfer*, the information is then further processed by the event sequence model in order to predict the success probability of the repayment process.

In this work, we present the conceptual architecture of ProcK together with an implementation based on deep learning models for graph-structured and sequential data. We also present results of an experimental study based on four datasets, three of them including a knowledge graph, from different applica-

tion domains. Our experiments demonstrate that ProcK achieves state-of-the-art predictive performance, which is further improved by utilizing the additional knowledge graph input.

2 Related work

2.1 Predictive process monitoring

Our work presents a new approach for *predictive process monitoring*, the task to predict future properties of processes from their execution logs. A considerable range of machine learning techniques have been studied in the context of process predictions. A review of seven methods that fall into the category of traditional machine learning (decision trees, random forests, support-vector machines, boosted regression, all with heavy feature engineering) has been published by Teinmaa *et al.* [2019].

Considering deep learning methods, due to the sequential nature of process logs, it is a straightforward approach to apply models designed for sequences. Tax *et al.* [2017] study the usage of LSTM neural networks for various prediction tasks, including the next activities and the remaining process time. Further approaches based on RNN and LSTM networks have been presented by Evermann *et al.* [2017]; Tello-Leal *et al.* [2018]; Camargo *et al.* [2019].

More recently, it has been demonstrated that feedforward networks often outperform recurrent neural networks for predictive process monitoring tasks. Al-Jebrni *et al.* [2018] employ 1D-convolutional networks, Mauro *et al.* [2019] study stacked inception CNNs, and Pasquadibisceglie *et al.* [2019] present a method where traces and their prefixes are first mapped onto a 2D image-like structure and then 2-dimensional CNNs are applied. Finally, Taymouri *et al.* [2020] present a prediction approach leveraging generative adversarial networks.

Although the main idea presented and evaluated in this work is independent from the particular choice of the sequence processing model, we share the experience of Al-Jebrni *et al.* [2018] and other authors that feedforward networks are more reliable to achieve good results in the process monitoring domain.

A direction that has been followed by several researchers is to make use of explicit process models. More than two decades of research on process mining has yielded sophisticated algorithms to create graph-shaped models of processes from event log data, often in the form of Petri nets (Van Der Aalst *et al.* [2011a]). From the viewpoint of such a process model, events trigger state changes of process instances, and predictive process monitoring models can take the current process state as input. Van der Aalst *et al.* [2011b] propose a solution to the problem of predicting the completion time, simply by calculating the mean remaining time for each state. Prediction from a combination of partial process models and event annotations has been performed by Ceci *et al.* [2014]. Folino *et al.* [2014] combine some of the aforementioned ideas, first clustering events and processes to achieve more abstract process models and then applying cluster-specific prediction models. Recently, Theis and Darabi [2019] have presented a

methodology to annotate states with time and other information, and use this information as input to a deep learning model.

What makes the idea of using a graphical process model powerful is that it builds upon a mature topic in the process mining domain, where the graph can be constructed on top of the existing event log. Our work, in contrast, is built on the hypothesis that, in concrete applications, *additional* non-sequential data is available, and that data is naturally modeled as a graph. In other words, process modeling is a specific type of feature pre-processing, while our approach taps a previously unused resource of data.

2.2 Combined sequence and graph models

The machine learning model we employ in our work uses a combination of graph-structured and sequential data as input. To our knowledge, we are the first to apply this approach to the domain of predictive process monitoring. Nevertheless, the benefits of combining such complementary views have been demonstrated in other domains. Fang *et al.* [2018] construct a graph based on spatial distances between mobile cells. Their model for cellular demand prediction uses an LSTM to compute a feature vector for each cell based on its demand history, then a graph convolutional network is used to model influences between nearby cells. Hu *et al.* [2019] address the problem of predicting the freezing-of-gaits symptom of Parkinson disease patients from video segments, where the first layer of their model is used to compute a representation of anatomic joints and their interactions as a graph. The authors introduce specialized LSTM cells to model both time-based interactions between subsequent video segments and interactions between joints. Wang *et al.* [2020] combine the graph-based representation of molecules with the SMILES string representation. Two individual models are trained for those two input representations, and the models are combined using ensemble techniques. An application in the retail domain has been proposed by Chang *et al.* [2021]. Here the first step is to convert the sequence of past user interactions with items into a graph with nodes representing items, and graph-based models are then applied to predict the interest of users.

While a variety of approaches to combine graphical and sequential input have been proposed in the body of work mentioned above, our approach has some unique and novel features. Firstly, each of the above methods performs some sort of node or graph classification which is only assisted by the sequential input. In our work, event sequences constitute the primary input, and the knowledge graph is used to assist the model to interpret the event data. Secondly and most importantly, the knowledge graph is used in our work to capture overall knowledge about the domain and application context, while specific information about each instance of the prediction problem is represented by its event sequence.

Another direction of related work is machine learning for dynamic knowledge graphs, called temporal knowledge graphs, where nodes, attributes, or relations change over time. The area of temporal knowledge graph reasoning can be divided into the *interpolation* and the *extrapolation* setting. As described by Jin *et*

al. [2020], in the interpolation setting, new facts are predicted for time steps up until the current time step, taking into account time information from past and current time steps; see García-Durán *et al.* [2018]. The methods in the extrapolation setting predict facts for future time steps. Recent work in the extrapolation setting includes RE-NET by Jin *et al.* [2020] and CluSTeR by Li *et al.* [2021].

The predictive process management application in our work requires to treat events as independent input. Modeling each event as a knowledge graph element would be technically possible, but only at the price of scalability, as typical applications include tens or hundreds of millions of events. Thus, in our work events do not have a direct interpretation as nodes or edges, but they instead contain attributes in the form of references to graph nodes. Conversely, triples in our knowledge graphs are in general not necessarily interpretable as events.

To summarize the discussion of our work in light of state-of-the-art, we are the first to utilize a knowledge graph as additional input to predictive process monitoring models to help interpreting the event data. The idea is realized by a new type of neural network architecture which takes events as primary input and learns to utilize an additional knowledge graph to interpret the event data.

3 Preliminaries

Following the process mining terminology, an *event log* $\mathcal{L} = (L, C, T, A)$ consists of an event set L , a set C of cases, a set T of possible event types, and a set A of additional event attributes. Each event $\ell = (c_\ell, t_\ell, \tau_\ell, \alpha_\ell) \in L$ is a 4-tuple characterized by its case identifier $c_\ell \in C$, its event type $t_\ell \in T$, a timestamp $\tau_\ell \in \mathbb{N}$, and a partial assignment function $\alpha_\ell : A \rightarrow V$ which specifies the values of a subset of attributes. For each case $c \in C$ we define $L_c := \{\ell \in L \mid c_\ell = c\}$ as the subset of events belonging to case c .

The domain V of possible attribute values is arbitrary in general; in this work it will be assumed that V is the node set of a knowledge graph. This assumption represents only a mild limitation, because categorical attribute values that do not appear in the given knowledge graph can simply be interpreted as isolated nodes. In fact, our experimental study includes one dataset where no knowledge graph is given at all. Extending ProcK with the ability to incorporate numerical attributes *directly* (i.e. without discretizing them to categorical attributes) remains for future work.

A knowledge graph (also called *knowledge base*) $\mathcal{G} = (V, R, E)$ is a directed graph defined by the node set V , relation types R , and edges E , where each edge $(v, r, v') \in E$ is a triple containing the head node $v \in V$, relation type $r \in R$, and tail node $v' \in V$.

As mentioned above, in typical practical applications, the input data will originate from the databases of the organization that performs predictive process monitoring. Relational databases are likely to contain time-stamped records that can become events. At the same time, mutual references between different tables are a fundamental element of relational data models, which makes it straightforward to interpret a part of the database as a knowledge graph. For

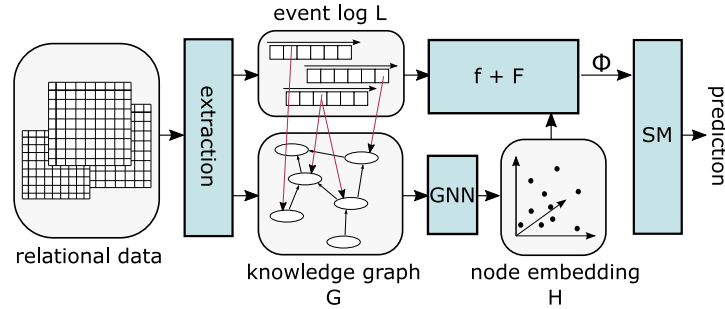


Fig. 1. ProcK conceptual architecture. From relational or tabular data, an event log and a knowledge graph is extracted. Then, the four components GNN, F , f , and SM of the neural network model subsequently compute node embedding H , event embeddings Φ , and the final prediction.

the purpose of our experimental study we have developed tools to extract both event data and a knowledge graph from a database dump.

4 ProcK architecture

We first specify, in Section 4.1, the conceptual architecture which is composed of four functional components GNN, f , F , and SM. Then, in Section 4.2, our implementation of each of the components is described.

4.1 Conceptual architecture

The architecture of ProcK, depicted in Figure 1, combines a graph neural network with a model for sequential data. From the bottom to the top of the network, input elements like nodes, edges, timestamps, and events will be encoded by embedding vectors. We employ a fixed embedding dimensionality $d \in \mathbb{N}$ across all layers.

The first component of the ProcK model is a graph neural network GNN which computes embedding vectors $H = (h_v)_{v \in V}$ containing an embedding $h_v \in \mathbb{R}^d$ for every node v of the knowledge graph. Formally,

$$h_v := \text{GNN}(\mathcal{G}, v), v \in V, \quad (1)$$

where $\mathcal{G} = (V, R, E)$ is the given knowledge graph.

Next, consider a single event $\ell = (c_\ell, t_\ell, \tau_\ell, \alpha_\ell)$ from the event log \mathcal{L} . Let F be an aggregation function for sets of d -dimensional vectors. The first step of constructing the event embedding is to compute

$$\beta_\ell := F(\{h_v \mid v \in \alpha_\ell(A)\}) \in \mathbb{R}^d. \quad (2)$$

Recall from the previous section that α_ℓ is a partial assignment function which specifies values for a subset of the attributes from A . In the above equation, $a_\ell(A)$ is the set of those attribute values.

Further, we employ a timestamp embedding function $f : \mathbb{N} \rightarrow \mathbb{R}^d$ and compute the final event embedding as

$$\phi_\ell := \beta_\ell + f(\tau_\ell) . \quad (3)$$

Having computed the sequence representation as a series of d -dimensional vectors $\Phi = (\phi_\ell)_{\ell \in L_c}$ for the given case $c \in C$, we feed the event representations into a sequence model SM to compute the prediction for the given case:

$$P_c := \text{SM}(\{\phi_\ell \mid \ell \in L_c\}) . \quad (4)$$

The structure of P_c depends on the prediction target; e.g. it can be a single real number for regression tasks or a vector of probabilities for a classification task. All four functions GNN, F , f , SM are potentially parameterized by trainable vectors $\Theta_{\text{GNN}}, \Theta_F, \Theta_f, \Theta_{\text{SM}}$, respectively.

4.2 Implementation

The bottom layer of our GNN implementation is a trainable embedding vector $h_v^0 \in \mathbb{R}^d$ for each node $v \in V$, as well as an embedding vector $h_r^0 \in \mathbb{R}^d$ for every relation $r \in R$. We compute higher-level node embeddings using graph convolution layers, where we adopt a simplified version of the compGCN architecture proposed by Vashishth *et al.* [2019] as described below. Given the layer i embeddings $(h_v^i)_{v \in V}, (h_r^i)_{r \in R}$, two transformations are applied for each $v \in V$:

$$h_{v,\text{self}}^{i+1} := W_{\text{self}}^{i+1} \cdot h_v^i \quad (5)$$

$$h_{v,\text{adj}}^{i+1} := W_{\text{adj}}^{i+1} \cdot \sum_{(v,r,v') \in E} \text{cmp}(h_r^i, h_{v'}^i) , \quad (6)$$

where $\text{cmp} : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the *composition operator*. Our implementation supports the composition operators of addition and element-wise multiplication; we employ the latter throughout our experiments. The node and relation embeddings on layer $i + 1$ are then computed via

$$h_v^{i+1} = \text{relu}(h_{v,\text{self}}^{i+1} + h_{v,\text{adj}}^{i+1}) , v \in V \quad (7)$$

$$h_r^{i+1} = W_{\text{rel}}^{i+1} \cdot h_r^i , r \in R . \quad (8)$$

The computations on each layer are parameterized by $W_{\text{self}}^{i+1}, W_{\text{adj}}^{i+1}, W_{\text{rel}}^{i+1} \in \mathbb{R}^{d \times d}$. Equation 6 specifies backward flow across the edges of the directed graph. In Vashishth *et al.* [2019], forward flow with independent parameterization is additionally specified. This is also supported by our implementation, but we only consider backward flow in our experiments. The reason is that the knowledge

graphs in our datasets contain nodes with a huge number of incoming links, and we experienced that summing up over them during the calculation of forward flow de-stabilizes the model and does not scale well.

Having computed k layers of graph convolution, the final node embeddings are the GNN output:

$$h_v = \text{GNN}(\mathcal{G}, v) := h_v^k, v \in V. \quad (9)$$

For the aggregation function F we employ mean pooling across the nodes referenced in each event:

$$\beta_\ell = F(\{h_v \mid v \in \alpha_\ell(A)\}) := \frac{1}{|\alpha_\ell(A)|} \sum_{v \in \alpha_\ell(A)} h_v. \quad (10)$$

For the timestamp embedding function f our implementation supports parameterized embedding and non-parameterized embedding based on sinusoids as described by Vaswani *et al.* [2017]. We have however found that, for the prediction tasks included by our experimental study, the timestamp input is not essential; thus we applied the constant zero function there for most datasets.

We now describe our implementation of the sequence model SM. First, a linear transformation is applied to the embeddings $(\phi_\ell)_{\ell \in L_c}$:

$$\phi'_\ell = W_1 \cdot \phi_\ell, \ell \in L_c. \quad (11)$$

After this initial transformation, we aggregate over the events of the sequence using mean pooling:

$$\phi'' = \frac{1}{|L_c|} \sum_{\ell \in L_c} \phi'_\ell. \quad (12)$$

A final fully connected hidden layer connects the aggregated events with the output:

$$\phi''' = \text{relu}(W_2 \cdot \phi''), \quad (13)$$

$$P_c = \text{SM}(\{\phi_\ell \mid \ell \in L_c\}) := g(W_3 \cdot \phi'''). \quad (14)$$

The dimensionality of the matrices is $W_1, W_2 \in \mathbb{R}^{d \times d}$ and $W_3 \in \mathbb{R}^{o \times d}$. For binary classification problems, $o = 1$, and g is the sigmoid activation function. For multi-class classification, o corresponds to the number of classes and g is the softmax function. Finally, for regression problems, $o = 1$ and g is the identity function.

We remark that the design choice of the sequence model is a result of exploratory experiments with various architectures. During those experiments we observed that, throughout the datasets, more sophisticated architectures (recurrent networks, transformer) did not lead to better results and introduced stability problems. This is in line with the finding, reported by Al-Jebrni *et al.* [2018] and Mauro *et al.* [2019], that feedforward networks outperform the LSTM architecture for predictive process monitoring tasks.

dataset	knowledge graph	event log	reference
OULAD	240K nodes, 1.1M edges	33K cases, 11M events	Kuzilek <i>et al.</i> [2017]
PKDD99	200K nodes, 1.1M edges	4.5K cases, 1M events	Berka [1999]
BPI12	–	13K cases, 180K events	van Dongen [2012]
DBLP	370K nodes, 1.2M edges	29K cases, 1.7M events	Tang <i>et al.</i> [2008]

Table 1. Summary of the event logs and knowledge graphs extracted from the datasets.

5 Experiments

5.1 Data

Our experiments encompass six prediction tasks using four datasets; see Table 1 for a summary.

Three tasks are based on the Open University Learning Analytics (OULAD) dataset provided by Kuzilek *et al.* [2017]. The dataset has the structure of a relational database dump, consisting of seven tables that represent information about students registering for courses, interacting with the study material, taking assessments and exams. From this data we extracted a knowledge graph as described in Figure 2. We further extracted an event log where each case $c \in C$

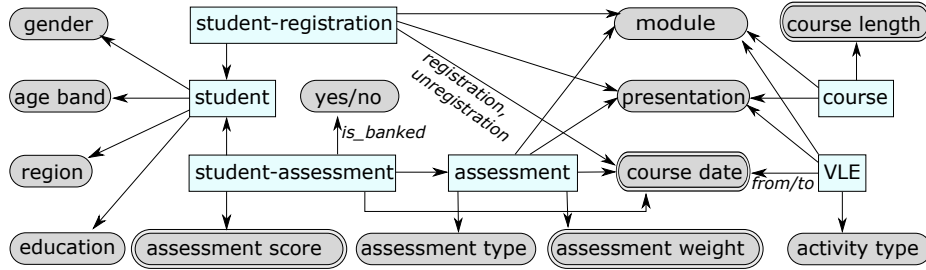


Fig. 2. Schema of the knowledge graph generated from the OULAD dataset. Box-shaped meta-nodes represent nodes generated from table rows, single-lined oval meta-nodes originate from categorical values, and double-lined meta-nodes represent discretized numerical values. Each arrow represents a distinct relation type (sometimes two); annotations have been added only where the type is not self-explanatory.

represents one student taking part in one course. There are five event types: *case info*, containing links to student, module, and semester, *assessment*, containing the assessment submission date as the timestamp, and a link to the corresponding student assessment node in the knowledge graph, *student registration* and *deregistration*, containing the date of (de)registration as the timestamp, and a link to the student registration node in the knowledge graph, and finally *VLE interaction*, containing as timestamp the time of interaction with material of the

Virtual Learning Environment (VLE), as well as a link to the knowledge graph node representing the material.

For the OULAD dataset we distinguish between three prediction targets: *dropout* (predict whether the student will drop out from the course), *success* (predict whether the student will finish the course successfully), and *exam score* (regression task to predict the final exam score, a number between 0 and 100). We further consider three different time horizons: *late prediction* is a variant of the prediction task where all events that happened during the course (except for the final exam and the deregistration event) are available as input. *Early prediction* refers to predictions taking into account only the events that happened before the 60th day of the course (all course modules take between 234 and 269 days), while in *very early prediction* no events have been recorded and only the *case info* is available. Across all tasks and their variants, we uniformly selected 20% of the data for validation and 5% for testing, like done by Jha *et al.* [2019].

A dataset from the financial domain was provided by Berka [1999] for the *PKDD99* challenge. This dataset also comes in the form of multiple interconnected tables, representing bank accounts, financial transactions, clients, geographical districts, and loans. The task is to predict the status of a loan (noncritical or critical) given the history of transactions and information about the loan and the client. The triples of the knowledge graph relate loans to accounts, accounts to districts, and bank orders as well as transactions to banks. Categorical and numerical attributes (after discretization) of those entities are additionally represented as neighbor nodes of them. The two event types extracted from the dataset are *case info*, including links to the loan and the account node, and *transaction*, containing a link to the node representing the transaction in the knowledge graph. We chose 20% of the data for validation and another 20% for testing. Because the number of samples in this dataset is comparably small and the dataset is rather imbalanced with only 10% of the loans having a critical status, we used stratified sampling to enforce the same balance of positive and negative examples in the training, validation, and test set.

Our second dataset from the financial domain represents an established process mining benchmark, but it comes only with an event log. The events contained in the *BPI12* dataset (van Dongen [2012]) were recorded during the application procedure at a financial institution. Due to the lack of a knowledge graph, we treat all event attributes as isolated nodes of a graph without edges. To make our experiments comparable to Tax *et al.* [2017], we only consider events that are marking the completion of manually executed subprocesses, and we used the latest 30% among the sequences as the test set. Another 20% of the sequences were chosen as the validation set uniformly at random. As an additional pre-processing step, we treat every prefix of every sequence as one sample where the task is to predict the type of the next event. This step was done separately for the training, validation, and test set.

Our final prediction problem is the number of future citations of papers published in the year 2000 as reported in the DBLP dataset introduced by Tang *et al.* [2008]. We extracted a knowledge graph consisting of relations between

task	dropout rate	l2-weight	GC layers	time embedding
OULAD (dropout)	0.1	0.01	1	none
OULAD (success)	0.1	0.01	1	none
OULAD (score)	0.7	0	1	none
PKDD99	0.1	0.03	3	none
BPI12	0.1	0.01	1	parameterized
DBLP	0.5	0	0	none
all tasks, LSTM model	0.25	0.01	-	implicit

Table 2. Hyperparameters used in the experiments. Note that the LSTM model implicitly takes into account the event position by design.

papers, authors, and venues. To prevent information leakage, events related to a publication after 2000 are not considered for the knowledge graph construction. For each published paper, the history of previous publications of all authors is used as the sequence of events. We uniformly selected 20% of the data for validation and 5% for testing.

5.2 Setup

The machine learning models were implemented in Tensorflow 2.5.0, and the computational experiments were performed on GPUs (Nvidia GeForce GTX 1080 Ti). For all classification tasks we used the cross-entropy loss function for training and selected the model having the highest accuracy on the validation set among all 200 training epochs. The learning rate was set to 0.01, and we chose an embedding width of $d = 100$. We also applied dropout and l2-regularization, which required different strategies for different tasks. Our implementation supports dropout with uniform rate after each aggregation layer in the graph convolutional network and after the final fully-connected dense layer. Table 2 lists the chosen strategy used for ProcK for every task, as well as for the LSTM baseline, where we found a different configuration to work best.

5.3 Results

The results of our experiments are displayed in Table 3. The top part of the table contains results for classification problems, where the accuracy and the Area Under the Curve (AUC, only for binary classification problems) metric are reported. On most datasets a knowledge graph is available, and we compare models trained with it to models trained without. As a baseline we also evaluated an LSTM model which was trained using the event sequence as input. Whenever available, the table also contains the best results reported in literature.

For the AUC metric, it turns out that the availability of additional data in form of a graph improves the ability of the model to correctly separate positive and negative test samples, both in comparison with ProcK without knowledge graph and the LSTM model. The improvement is observable across most of the

prediction task	model	accuracy	AUC
OULAD (dropout, late)	ProcK	0.86	0.93
	ProcK (no KG)	0.86	0.93
	LSTM	0.86	0.92
	GBM (Jha <i>et al.</i> [2019])	-	0.91
OULAD (dropout, early)	ProcK	0.83	0.84
	ProcK (no KG)	0.81	0.82
	LSTM	0.81	0.82
OULAD (dropout, very early)	ProcK	0.68	0.58
	ProcK (no KG)	0.69	0.60
	LSTM	0.70	0.57
OULAD (success, late)	ProcK	0.87	0.91
	ProcK (no KG)	0.88	0.88
	LSTM	0.86	0.86
	GBM(Jha <i>et al.</i> [2019])	-	0.93
OULAD (success, early)	ProcK	0.73	0.74
	ProcK (no KG)	0.73	0.73
	LSTM	0.75	0.73
OULAD (success, very early)	ProcK	0.69	0.58
	ProcK (no KG)	0.69	0.56
	LSTM	0.68	0.57
PKDD99	ProcK	0.89	0.71
	ProcK (no KG)	0.89	0.71
	LSTM	0.89	0.50
BPI12 (KG not available)	ProcK	0.83	-
	LSTM	0.71	-
	LSTM(Tax <i>et al.</i> [2017])	0.76	-
prediction target	model	RMSE	
OULAD (score, late)	ProcK	18.93	
	ProcK (no KG)	18.95	
	LSTM	20.35	
OULAD (score, early)	ProcK	19.88	
	ProcK (no KG)	19.88	
	LSTM	21.08	
OULAD (score, very early)	ProcK	20.10	
	ProcK (no KG)	20.44	
	LSTM	20.13	
DBLP	ProcK	3.98	
	LSTM	4.01	

Table 3. Experimental results

problems. For dropout prediction on the OULAD dataset, all three deep learning models outperform the results of Gradient Boosting Machine (GMB), reported by Jha *et al.* [2019], while GBM performs better for the success prediction task.

When looking at the accuracy metric, the advantage of the knowledge graph input is not that clearly visible; only on two out of seven problems the full ProcK model with knowledge graph exhibits the best performance, on one problem it is outperformed by ProcK without knowledge graph input, and two tasks the LSTM model performs best. The only classification problem with more than two classes is next event type prediction on the BPI12 dataset. Here no knowledge graph is given, and ProcK outperforms the LSTM model by a large margin. This finding is consistent to results found in other works e.g. by Mauro *et al.* [2019].

We are the first to study variants of OULAD with different points of prediction time (early and very early prediction). It turns out that the length of the event log makes a significant difference, with AUC values decreasing to less than 0.60 when only the initial case information is available. However, the benefit of using knowledge graph input does not seem to depend on the length of the event log, which can be explained with the fact that less event input also means less information from the knowledge graph.

The bottom part of Table 3 contains the results for regression problems, including exam score (OULAD dataset) and number of citations (DBLP dataset). We do not include ProcK without knowledge graph input for the latter, because here the number of graph convolution layers is set to zero (see Table 2), making models with and without knowledge graph input equivalent. The table reports the root mean square error, and here again the benefits of the knowledge graph input can be demonstrated across the variants of the score prediction task.

6 Summary and conclusion

In this work, we introduced ProcK, a novel machine learning pipeline for data from knowledge-intensive processes. Within the pipeline, two complementary views of the available information are first extracted from raw tabular data and then re-combined as input to the downstream prediction model. We implemented prototypes of each pipeline component, and we tested their interplay on six prediction tasks on four datasets. We could demonstrate on the majority of classification tasks that ProcK achieves improved AUC values when having a knowledge graph available as input, but further investigation of the accuracy metric remains a task for future work. Also for regression tasks ProcK exhibits a small but consistent advantage in terms of the RMSE metric.

There are several additional directions for future work. One interesting question from a practical viewpoint is how machine learning can be employed to extract the knowledge base and event log from the source databases in an automatic and configuration-free manner. Furthermore, while ProcK has been applied for prediction tasks so far, the ability to make process recommendations (e.g. events that should happen in the future to positively influence the outcome of cases) will be an important next step.

Ethics discussion All experiments reported in this work are based on anonymized datasets (OULAD, PKDD, BPI12) or data actively published by the data subjects (DBLP). Nevertheless, the presented technology is applicable to ethically sensitive tasks, including assessment of loan applications and performance prediction of humans. A careful assessment of potential ethical issues has to be carried out prior to bringing this work to application.

Bibliography

- Abdulrhman Al-Jebrni, Hongming Cai, and Lihong Jiang. Predicting the next process event using convolutional neural networks. In *2018 IEEE International Conference on Progress in Informatics and Computing (PIC)*, pages 332–338. IEEE, 2018.
- Petr Berka. Workshop notes on Discovery Challenge PKDD’99. 1999.
- Manuel Camargo, Marlon Dumas, and Oscar González-Rojas. Learning accurate lstm models of business processes. In *International Conference on Business Process Management*, pages 286–302. Springer, 2019.
- Michelangelo Ceci, Pasqua Fabiana Lanotte, Fabio Fumarola, Dario Pietro Cavallo, and Donato Malerba. Completion time and next activity prediction of processes using sequential pattern mining. In *International Conference on Discovery Science*, pages 49–61. Springer, 2014.
- Jianxin Chang, Chen Gao, Yu Zheng, Yiqun Hui, Yanan Niu, Yang Song, Depeng Jin, and Yong Li. Sequential recommendation with graph neural networks. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 378–387, 2021.
- Joerg Evermann, Jana-Rebecca Rehse, and Peter Fettke. Predicting process behaviour using deep learning. *Decision Support Systems*, 100:129–140, 2017.
- Luoyang Fang, Xiang Cheng, Haonan Wang, and Liuqing Yang. Mobile demand forecasting via deep graph-sequence spatiotemporal modeling in cellular networks. *IEEE Internet of Things Journal*, 5(4):3091–3101, 2018.
- Francesco Folino, Massimo Guarascio, and Luigi Pontieri. Mining predictive process models out of low-level multidimensional logs. In *International conference on advanced information systems engineering*, pages 533–547. Springer, 2014.
- Alberto García-Durán, Sebastijan Dumancic, and Mathias Niepert. Learning sequence encoders for temporal knowledge graph completion. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4816–4821. Association for Computational Linguistics, 2018.
- Kun Hu, Zhiyong Wang, Wei Wang, Kaylena A Ehgoetz Martens, Liang Wang, Tieniu Tan, Simon JG Lewis, and David Dagan Feng. Graph sequence recurrent neural network for vision-based freezing of gait detection. *IEEE Transactions on Image Processing*, 29:1890–1901, 2019.
- Nikhil Indrashekhara Jha, Ioana Ghergulescu, and Arghir-Nicolae Moldovan. Oulad mooc dropout and result prediction using ensemble, deep learning and regression techniques. In *CSEdu (2)*, pages 154–164, 2019.
- Woojeong Jin, Meng Qu, Xisen Jin, and Xiang Ren. Recurrent event network: Autoregressive structure inference over temporal knowledge graphs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6669–6683, 2020.
- Jakub Kuzilek, Martin Hlostá, and Zdeněk Zdrahal. Open university learning analytics dataset. *Scientific data*, 4(1):1–8, 2017.

- Anna Leontjeva, Raffaele Conforti, Chiara Di Francescomarino, Marlon Dumas, and Fabrizio Maria Maggi. Complex symbolic sequence encodings for predictive monitoring of business processes. In *International Conference on Business Process Management*, pages 297–313. Springer, 2016.
- Zixuan Li, Xiaolong Jin, Saiping Guan, Wei Li, Jiafeng Guo, Yuanzhuo Wang, and Xueqi Cheng. Search from history and reason for future: Two-stage reasoning on temporal knowledge graphs. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers)*, pages 4732–4743. Association for Computational Linguistics, 2021.
- Nicola Di Mauro, Annalisa Appice, and Teresa Basile. Activity prediction of business process instances with inception cnn models. In *International conference of the italian association for artificial intelligence*, pages 348–361. Springer, 2019.
- Vincenzo Pasquadibisceglie, Annalisa Appice, Giovanna Castellano, and Donato Malerba. Using convolutional neural networks for predictive process analytics. In *2019 international conference on process mining (ICPM)*, pages 129–136. IEEE, 2019.
- Lars Reinkemeyer. *Process Mining in Action*. Springer International Publishing, 2020.
- Jie Tang, Jing Zhang, Limin Yao, Juanzi Li, Li Zhang, and Zhong Su. Arnetminer: extraction and mining of academic social networks. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 990–998, 2008.
- Niek Tax, Ilya Verenich, Marcello La Rosa, and Marlon Dumas. Predictive business process monitoring with lstm neural networks. In *International Conference on Advanced Information Systems Engineering*, pages 477–492. Springer, 2017.
- Farbod Taymouri, Marcello La Rosa, Sarah Erfani, Zahra Dasht Bozorgi, and Ilya Verenich. Predictive business process monitoring via generative adversarial nets: the case of next event prediction. In *International Conference on Business Process Management*, pages 237–256. Springer, 2020.
- Irene Teinemaa, Marlon Dumas, Marcello La Rosa, and Fabrizio Maria Maggi. Outcome-oriented predictive process monitoring: Review and benchmark. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 13(2):1–57, 2019.
- Edgar Tello-Leal, Jorge Roa, Mariano Rubiolo, and Ulises M Ramirez-Alcocer. Predicting activities in business processes with lstm recurrent neural networks. In *2018 ITU Kaleidoscope: Machine Learning for a 5G Future (ITU K)*, pages 1–7. IEEE, 2018.
- Julian Theis and Houshang Darabi. Decay replay mining to predict next process events. *IEEE Access*, 7:119787–119803, 2019.
- Wil Van Der Aalst, Arya Adriansyah, Ana Karla Alves De Medeiros, Franco Arcieri, Thomas Baier, Tobias Blickle, Jagadeesh Chandra Bose, Peter Van Den Brand, Ronald Brandtjen, Joos Buijs, et al. Process mining manifesto.

- In *International Conference on Business Process Management*, pages 169–194. Springer, 2011.
- Wil MP Van der Aalst, M Helen Schonenberg, and Minseok Song. Time prediction based on process mining. *Information systems*, 36(2):450–475, 2011.
- Boudewijn van Dongen. Bpi challenge 2012, Apr 2012.
- Shikhar Vashishth, Soumya Sanyal, Vikram Nitin, and Partha Talukdar. Composition-based multi-relational graph convolutional networks. *arXiv preprint arXiv:1911.03082*, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- Zhengyang Wang, Meng Liu, Youzhi Luo, Zhao Xu, Yaochen Xie, Limei Wang, Lei Cai, Qi Qi, Zhuoning Yuan, Tianbao Yang, et al. Advanced graph and sequence neural networks for molecular property prediction and drug discovery. *arXiv preprint arXiv:2012.01981*, 2020.