

# CDESf Tutorial

August 19, 2020

## Abstract

Concept Drift in Event Stream Framework (CDESf) is a completely available tool for the detection of drifts in event streams. Moreover, CDESf supports several online Process Mining tasks, such as process model evolution and anomaly detection. In this tutorial, we present a practical guide to the main aspects of CDESf and how to explore the results.

## 1 Introduction

CDESf is built for Python 3.6 or higher versions. The library dependencies required to run CDESf are the following: networkx  $\sim$  2.4, numpy  $\sim$  1.18.4, setuptools  $\sim$  41.2.0, pandas  $\sim$  1.0.4, scipy  $\sim$  1.5.1, matplotlib  $\sim$  3.2.2, seaborn  $\sim$  0.10.1, pygraphviz  $\sim$  1.6, graphviz  $\sim$  2.38.0. The first step is to clone the repository<sup>1</sup> into your local machine. Following, an example on how to trigger the processing of event streams in CDESf.

```
1  from cdesf2.utils import read_csv
2  from cdesf2.core import CDESf
3
4  path = "demo"
5  filename = "hospital_billing.csv"
6
7  event_stream_test = read_csv(f"{path}/{filename}")
8  cdesf = CDESf(name='hospital_billing',
9               time_horizon=259200,
10              lambda_=0.05,
11              beta=0.2,
12              epsilon=0.2,
13              mu=4,
14              stream_speed=100,
15              n_features=2,
16              gen_metrics=True,
17              gen_plot=True)
18  cdesf.run(event_stream_test)
```

The hyperparameters that control event processing are crucial to execution of CDESf. We list them in detail next:

---

<sup>1</sup><https://github.com/gbrltv/cdesf2>

- *path*: path of the folder containing the event log files.
- *filename*: event log file to be processed.
- *name*: the name of the process you want to analyze (also used to save the output).
- *time\_horizon*: window that controls both model updating and forgetting mechanism. Smaller time horizons trigger more frequent updates, hence, a higher rate of concept drifts. Larger time horizons obtain fewer updates, hence, the process model is more stable and less drift warnings will be triggered. A good value is the mean time of a case execution. Time horizon uses seconds as a time unit.
- *lambda*: the decay factor impacts on how much historical data is relevant. Higher values implies in less importance to historical data. The optimal range lies between 0.05 and 0.4.
- *beta*: controls the promotion of micro-clusters, consequently impacting on drift sensibility. Higher values are more robust to drifts. Instead, lower values trigger more drift warnings. Usually the range of values is between 0.1 and 0.6.
- *epsilon*: represents the maximum actuation range of a micro-cluster. Higher values enable bigger micro-clusters, and consequently more cover area. Optimal values range from 0.05 to 0.5.
- *mu*: used in conjunction with *beta*, controls the minimum density of a micro-cluster to be considered core behavior.
- *stream\_speed*: controls the application of the decay factor, i.e., influences how fast the micro-clusters weight decay.
- *n\_features*: the number of dimensions DenStream must consider, here it is set to 2 since we have two attributes (graph distance and time distance). However, CDESf can consider multiple dimensions too.
- *gen\_metrics*: controls metrics generation and saving. Useful to investigate case labels and micro-cluster behavior.
- *gen\_plot*: control the generation of feature space plots after each event. Useful for detailed analysis of the micro-clusters evolution.

## 2 Reading files

CDESf framework is capable to read and analyze event streams to perform real-time tasks, such as process discovery, anomaly detection, clustering and concept drift detection. For that, it has to read the log file and extract the data. This operation is made by the following reading function:

```

1 def read_csv(path: str) -> np.ndarray:
2     event_stream = pd.read_csv(path,
3                               usecols=['case', 'activity', '
4                                     timestamp'],
5                                     parse_dates=['timestamp'],
6                                     infer_datetime_format=True,
7                                     memory_map=True)
8     event_stream['activity'].replace(' ', '_', regex=True, inplace=
    True)
    return event_stream.values

```

This function uses a pandas component to read the file at the given path. Moreover, it also converts timestamps to the datetime object and transform names of activities to make them suitable for CDESf execution. After that, the function returns the stream of events as an array. In detail for pandas function *read\_csv* we use the parameters

- *usecols*: it is required that the log file contains the columns *case*, *activity* and *timestamp*. Otherwise the file will not be read.
- *parse\_dates*: selects the column containing the date and time when events were executed.
- *infer\_datetime\_format*: infers the datetime string format and transforms to a python object.
- *memory\_map*: improves performance by mapping the file directly to memory.

### 3 Case and Activity Structures

This section covers the basic structures of CDESf. As presented previously, the reading function prepares the log file and identifies cases, activities and timestamps. CDESf has two data structures to represent basic Process Mining concepts: Case and Activity. The activities and timestamps taken are respectively activity name and timestamp attributes of Activity objects. Moreover, the case column provides the identifiers of Case objects. In details, the Case object has the following attributes:

- *id*: the case identifier, taken and assigned from the event log.
- *activities*: a list of activities of type Activity executed within the case.
- *graph\_distance*: value produced by the comparison between the case flow and the process model.
- *time\_distance*: value produced by the comparison of time between case activities and the process model.

Case has also two properties. The first one calculate its point, which is a two-element array containing *graph\_distance* and *time\_distance*. The point is used later in the clusterization step. The property's structure is the following:

```

1 @property
2 def point(self) -> np.array:
3     return np.array([self.graph_distance, self.time_distance])

```

The second property retrieves the last event timestamp:

```

1 @property
2 def last_time(self) -> datetime:
3     return self.activities[-1].timestamp

```

Furthermore, Case has also other useful methods:

- `set_activity(activity_name, activity_timestamp)`

Allows the attachment of a new Activity object to the Case.

- `get_trace()`

Retrieves the list of activities within that case.

- `get_timestamp()`

Retrieves the list of activities timestamps.

## 4 Running CDESf

With the basic structures presented, we can proceed to analyze how CDESf runs. As explained previously, the first step is reading the log file and storing its return in a variable. Then, a *run* function can be triggered with the stream and the hyperparameters. The function simulates the event stream by iterating through the stream returned from the reading. Each new event is processed (*process\_event*) using the case, the activity name and the timestamp as input from the stream. If it is the first event in the stream, the function sets the activity timestamp as the first check point. Once the stream is over, an output with the total number of drifts and the drift points is shown on the console. Moreover, a *visualization* function creates a plot with all drifts to provide further insights. The detailed function is:

```

1 for index, event in enumerate(stream):
2     self.event_index = index
3     case_id, activity_name, activity_timestamp
4     = (event[0], event[1], event[2])
5     if index == 0:
6         self.check_point = activity_timestamp
7     self.process_event(case_id, activity_name,
8                       activity_timestamp)

```

### 4.1 Processing events

The *process\_event* function is the core of CDESf. For new events that belong to never seen cases, it instantiates a new case, otherwise, it adds the new activity and timestamp to an existing case through the *set\_case* function. In details *set\_case* as follows:

```

1     index = self.get_case(case_id)
2     if index is None: # new case
3         self.cases.append(Case(case_id))
4         self.check_point_cases += 1
5         index = self.get_case(case_id)
6     self.cases[index].set_activity(act_name, act_timestamp) #
7     # appends new activity
8     self.cases.insert(0, self.cases.pop(index))
9     return index

```

It checks through *get\_case* if a case exists and it creates one if it does not exist. After creating the case, CDESf adds the activity with the timestamp to the case list in the first position, maintaining a first in first out logic. Then *process\_event* checks if the difference between the last timestamp of the modified case and the check point is greater than the *time\_horizon* hyperparameter. If the condition is met for the first time, the initialization function (*initialize\_cdesf*) is triggered. The goal is to create the first process model, extract the graph distances of the existing cases and initialize the online clustering feature space. First, *initialize\_cdesf* sets system variables and creates the first Process Model Graph (PMG) and, with the extracted case distances to the graph, CDESf initializes the micro-clusters used by the DenStream algorithm in the online clustering step. If the hyperparameters *gen\_metrics* and *gen\_plot* are set to *True*, it also initializes metrics and graphic functions, respectively. The part of the code regarding the initialization of the system variables and DenStream is the following:

```

1 self.nyquist = self.check_point_cases * 2
2 self.process_model_graph = initialize_graph(nx.DiGraph(), self.
3     cases)
4 self.initialize_case_metrics()
5 self.denstream.dbscan(self.cases)
6 groups = self.denstream.generate_clusters()
7 for group in groups[0]:
8     for cluster in group:
9         self.active_core_clusters.add(cluster.id)

```

## 4.2 Graph operations

As you can see from the code above, the PMG graph structure relies on the networkx DiGraph object. CDESf uses it to build PMG as directed graphs since it successfully represents directly-follows relations. There are three main graph operations:

- *normalize\_graph*: makes a normalization of time and weight for each edge in the graph. From a trace perspective, the normalization is the occurrence of a specific transition divided by the value of the most occurred transition. The result of this operation is referred to as *weight\_normalized*. For the time perspective, edge normalization computes the mean time of a transition between activities. The result of this operation is referred to as *time\_normalized*.

```

1 max_weight = max([attributes['weight']
2                   for n1, n2, attributes
3                   in graph.edges(data=True)])
4 for node1, node2, data in graph.edges(data=True):
5     data['weight_normalized'] = data['weight'] / max_weight
6     data['time_normalized'] = data['time'] / data['weight']
7 return graph
8

```

- *initialize\_graph*: this function initializes a graph based on the weights and time differences from a list of cases. The resulting graph has activities as nodes and *time* and *weight* as edges' attributes. It creates the edge if it does not already exist, otherwise it increments its weight by 1 and time by time difference for each directly-follows relation. The normalized graph is returned.

```

1     trace_list, time_list = extract_cases_time_and_trace(
2         case_list)
3     time_list = time_difference(time_list)
4
5     for trace, time in zip(trace_list, time_list):
6         for i in range(len(trace)-1):
7             edges = (trace[i], trace[i+1])
8             if edges not in graph.edges:
9                 graph.add_edge(*edges, weight=1, time=time[i])
10            else:
11                graph[edges[0]][edges[1]]['weight'] += 1
12                graph[edges[0]][edges[1]]['time'] += time[i]
13
14 return normalize_graph(graph)

```

- *merge\_graph*: this function receives two graphs and merge them. Typically the two graphs are the PMG and the graph built using only new cases, i.e., cases that appeared in the last time horizon, named as Check Point Graph (CPG). First of all this function decays the 5% of the PMG's edges weight. For each edge in CPG, it checks if there is a corresponding edge in the PMG. If it already exists, it updates the weight and the time of the PMG edge with the corresponding weight and time of the CPG edge. If it does not exists in the PMG, it creates the new edge and assigns the weight and the time value taken from CPG. Finally the function returns the normalized (PMG).

```

1 for node1, node2, data in process_model_graph.edges(data=True):
2     :
3     data['weight'] *= 0.95
4
5 for node1, node2, data in check_point_graph.edges(data=True):
6     path = (node1, node2)
7     if path in process_model_graph.edges:
8         process_model_graph[node1][node2]['weight'] += data['weight']
9         process_model_graph[node1][node2]['time'] += data['time']

```

```

9         else:
10             process_model_graph.add_edge(*path, weight=data['
               weight'], time=data['time'])
11
12     return normalize_graph(process_model_graph)
13

```

### 4.3 Case metrics

After process model creation, case metrics are initialized in the function *initialize\_case\_metrics*, which extracts case distances,  $GD_{trace}$  (*graph\_distance*) and  $GD_{time}$  (*time\_distance*), for cases in the first time horizon cycle. Here the detailed code:

```

1 def initialize_case_metrics(self) -> None:
2     for case in self.cases:
3         graph_distance, time_distance = extract_case_distances(self
               .process_model_graph, case)
4         case.graph_distance = graph_distance
5         case.time_distance = time_distance

```

The *extract\_case\_distances* function receives a graph and a case and returns graph and time distances based on Equations 1 and 2:

$$GD_{trace}(tr) = \frac{\sum_{i=1}^{T_{tr}} 1 - PMG_{weight}(tr[i])}{T_{tr}} \quad (1)$$

$$GD_{time}(tr) = \log_{10} \left( \frac{\sum_{i=1}^{T_{tr}} |PMG_{time}(tr[i]) - tr[i]_{time}|}{\sum_{j=1}^{T_{tr}} PMG_{time}(tr[j])} \right) \quad (2)$$

Given a trace  $tr$ ,  $T_{tr}$  is the total amount of edges in  $tr$ ,  $tr[i]$  corresponds to the  $i$ th edge in  $tr$ ,  $tr[i]_{time}$  is the time delta in  $tr[i]$  and  $tr[i]_{weight}$  is the weight of  $tr[i]$ .

### 4.4 Online clustering

After the case metrics are initialized, *initialize\_cdesf* triggers DenStream, which starts its micro-clusters using the DBSCAN algorithm. First of all a case is taken and all other cases in the buffer passed in input, The input for this function is the list of cases currently active. One by one, cases are analyzed and their euclidean distances are measured. A micro-cluster is a group of cases with a distance lower than the *epsilon* hyperparameter. This process is repeated recursively until all cases are used. If the formed group of cases has enough weight, i.e., the weight is higher than  $\beta * \mu$  (hyperparameters), then a micro-cluster is created with them. If weight is lower than  $\beta * \mu$ , the cases are anomalous, since they belong to a low density region. The complete code of DBSCAN is the following:

```

1 def dbscan(self, buffer: List) -> None:
2     used_cases = set()
3     for case in (case for case in buffer if case.id not in
               used_cases):

```

```

4         used_cases.add(case.id)
5         group = [case]
6         for other_case in (case for case in buffer if case.id not
7         in used_cases):
8             dist = self.euclidean_distance(case.point, other_case.
9             point)
10            if dist <= self.epsilon:
11                group.append(other_case)
12
13            weight = len(group)
14            if weight >= self.beta * self.mu:
15                new_p_mc = MicroCluster(id_=self.mc_id,
16                n_features=self.n_features,
17                creation_time=self.time,
18                lambda_=self.lambda_)
19
20            for case_ in group:
21                used_cases.add(case_.id)
22                new_p_mc.update(case_)
23                self.all_cases[case_.id] = new_p_mc.id
24            self.mc_id += 1
25            self.p_micro_clusters.append(new_p_mc)
26        else:
27            used_cases.remove(case.id)

```

Once CDESf is initialized, for each new event, `process_event()` calculates the distances of the case and assigns them to its attributes. Then it triggers DenStream using the function `train` to accomodate the case in the online clustering space.

```

1 def train(self, case: Case):
2     micro_cluster_updated = self.add_point(case)
3     self.all_cases[case.id] = micro_cluster_updated
4
5     self.no_processed_points += 1
6     if self.no_processed_points % self.stream_speed == 0:
7         self.time += 1
8         for mc in self.p_micro_clusters:
9             mc.decay()
10
11     if self.time % self.tp == 0:
12         for i, mc in enumerate(self.p_micro_clusters):
13             if mc.weight < self.beta * self.mu:
14                 self.p_micro_clusters.pop(i)
15
16         for i, mc in enumerate(self.o_micro_clusters):
17             to = mc.creation_time
18             e = ((math.pow(2, - self.lambda_ * (self.time - to +
19             self.tp)) - 1) /
20                 (math.pow(2, - self.lambda_ * self.tp) - 1))
21             if mc.weight < e:
22                 self.o_micro_clusters.pop(i)

```

This way, `train` adds a new point to a micro-cluster through the `add_point` function based using distances in a sequence of steps:

- first, Denstream tries merge the case into the closest potential micro-cluster:



```

1 i, closest_p_mc, _ = self.find_closest_mc(case.point, self.
    p_micro_clusters)
2 if closest_p_mc.radius_with_new_point(case.point) <= self.
    epsilon:
3     closest_p_mc.update(case)
4     return closest_p_mc.id
5

```

- if it fails, it tries to merge the case outlier micro-cluster and eventually try to promote it to a potential micro-cluster. To have a promotion the outlier micro-cluster must have its weight higher than  $\beta\mu$ :

```

1 i, closest_o_mc, _ = self.find_closest_mc(case.point, self.
    o_micro_clusters)
2 if closest_o_mc.radius_with_new_point(case.point) <= self.
    epsilon:
3     closest_o_mc.update(case)
4     if closest_o_mc.weight > self.beta * self.mu:
5         self.o_micro_clusters.pop(i)
6         self.p_micro_clusters.append(closest_o_mc)
7         return closest_o_mc.id
8

```

- if the merge fails, a new outlier micro-cluster is created with this case:

```

1 new_o_mc = MicroCluster(id_=self.mc_id,
2                           n_features=self.n_features,
3                           creation_time=self.time,
4                           lambda_=self.lambda_)
5 self.mc_id += 1
6 new_o_mc.update(case)
7 self.o_micro_clusters.append(new_o_mc)
8 return new_o_mc.id
9

```

A dedicated function (*find\_closest\_mc*) to find the closest micro-clusters is needed for these operations. It takes the point and the micro-cluster list and it uses the euclidean distance between the point and the cluster's centroid to find the closest cluster.

```

1 if len(micro_cluster_list) == 0:
2     raise NoMicroClusterException
3
4 distances = [(i, self.euclidean_distance(point, np.array(cluster.
    centroid)))
5               for i, cluster in enumerate(micro_cluster_list)]
6 i, dist = min(distances, key=lambda i_dist: i_dist[1])
7 return i, micro_cluster_list[i], dist

```

After adding the point, the *train* function decays the potential micro-clusters weights. This decay frequency is controlled by the *lambda* hyperparameter. Outlier micro-clusters with weights lower than *epsilon* are deleted.

## 4.5 Check point

In *process\_event*, if the difference between the event timestamp and the check point is greater than the *time\_horizon*, a check point is triggered. The main goal of a check point is to update the PMG and control memory consumption by releasing older cases. Moreover, the PMG is saved at each check point. The following code controls check point execution:

```
1 if (current_time - self.check_point).total_seconds() > self.  
    time_horizon:  
2     self.check_point = current_time  
3     self.check_point_update()  
4  
5     self.metrics.save_pmg_on_check_point(self.process_model_graph,  
        self.cp_count)  
6     if self.gen_metrics:  
7         self.metrics.save_case_metrics_on_check_point()  
8         self.metrics.save_cluster_metrics_on_check_point()
```

The *check\_point\_update* releases older cases according to the Nyquist frequency theorem and updates the PMG using the *merge\_graphs* function:

```
1 def check_point_update(self):  
2     self.cp_count += 1  
3     if len(self.cases) > self.nyquist:  
4         self.release_cases_from_memory()  
5         if self.check_point_cases > 5:  
6             self.nyquist = self.check_point_cases * 2  
7             check_point_graph = initialize_graph(nx.DiGraph(), self.  
                cases)  
8             self.process_model_graph = merge_graphs(self.  
                process_model_graph, check_point_graph)  
9             self.check_point_cases = 0
```

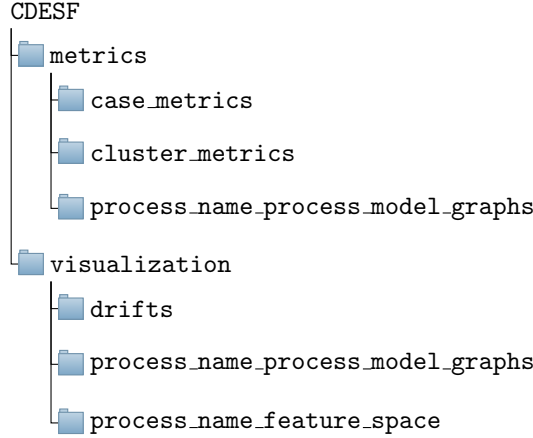
If the length of the list of cases is greater than the Nyquist value, older cases are released. This operation is guaranteed by *release\_case\_from\_memory* function:

```
1 def release_cases_from_memory(self) -> None:  
2     self.cases = self.cases[:self.nyquist]
```

Also *check\_point\_update* updates the Nyquist value according to the theorem using the frequency of new cases arrives in the last time horizon.

## 5 Metrics and visualization

Metrics and plotting are triggered in CDESF by two flags: *gen\_metrics* and *gen\_plot*. Two root folders contain CDESF produced outputs:



Note that “process\_name” will be filled with the *name* hyperparameter. The case metrics (stored inside “metrics/case\_metrics”) contains the stream index position, the timestamp, the check point number, the case identifier, the graph and time distances and the label of a case. The label is either normal or anomalous. If a case belongs to a low density region, DenStream classifies it as anomalous, otherwise, it is a normal execution. The cluster metrics (stored inside “metrics/cluster\_metrics”) contains the stream index position, the timestamp, the check point, the cluster identification, position, radius, weight and type (core, potential or outlier). The metrics files are useful to understand the evolution of objects in the stream processing. For instance, a cluster can be traced since its discovery in the feature space to analyze how it behaves in the events processing. Metrics functions are started in *initialize\_cdesf*. Further, it is controlled by *save\_case\_metrics\_on\_check\_point* and *save\_cluster\_metrics\_on\_check\_point*. Moreover, the PMG is saved at each check point (*save\_pmg\_on\_check\_point*), making possible to observe its evolution.

- *compute\_case\_metrics*: this function takes in input from CDESF the event index, the timestamp, the current check point value, the case to be saved, and a boolean flag label to control if case is normal or anomalous. *compute\_case\_metrics* appends the case metrics to the *case\_metrics* attribute.
- *save\_case\_metrics\_on\_check\_point*: works in conjunction with *compute\_case\_metrics*. It saves the case metrics into the output file on every check point. Moreover, the function frees the *case\_metrics* attribute to save memory.
- *compute\_cluster\_metrics*: generates the cluster metrics and save them into the *cluster\_metrics* attribute. Cluster metrics come from the online processing, provided by DenStream. The saved values are the event index, the check point timestamp, the current check point number and cluster attributes.
- *save\_cluster\_metrics\_on\_check\_point*: works in conjunction with *compute\_cluster\_metrics*.

It takes the data written in the *cluster\_metrics* attribute and saves into a file, then, releases the *cluster\_metrics* attribute.

- *save\_pmg\_on\_check\_point*: this function takes the PMG as input, then, it saves the PMG in a *JSON* file and also plots it using the function for visualization. This function plots the graph using the *pygraphviz* library.

Another crucial visualization is provided by the *feature\_space* function, which plots the clusters and cases in the feature space after every new event in the stream. Core micro-clusters (black ellipses) represent common process behavior, potential micro-clusters (blue ellipses) represent potential core behavior and outlier micro-clusters (red ellipses) represent outlier behavior. Normal and anomalous cases are also shown in black dots and yellow markers, respectively. The x-axis is marked as  $GD_{trace}$  and the y-axis is marked as  $GD_{time}$ . The plots are saved under the “process\_name.feature\_space” folder.

## 6 Experiments

To test the effectiveness of the CDESf toolkit, we have analyzed a synthetic event log, which contains one drift and 20% of anomalous traces<sup>2</sup>.

The hyperparameters for the framework were set to: *time\_horizon*=259200, *lambda*=0.05, *beta*=0.2, *epsilon*=0.6, *mu*=4, *stream\_speed*=100, *n\_features*=2, *gen\_metrics*=True, *gen\_plot*=True.

The first output is the console drift warnings. This fast identification of drifts is interesting in scenarios where actions must be taken in a short timespan. Figure 1 shows a drift warning using the synthetic event log. Furthermore, Figure 1 points to a drift in event 2532 and that the new detected core micro-cluster is the number 4. To complement this information, the cluster metrics file can be analyzed to further investigate cluster 4 evolution. Figure 2 presents the metrics output for this event log. The first appearance of cluster number 4 is on event 2525 (row 1546), as an outlier micro-cluster. Note that on previous events there were only cluster number 0, 1, 2 and 3. From event 2525 on, all subsequent events maintain the cluster 4, as Figure 2 shows.

Once the drift is found, the feature space at that event can be analyzed for further insights. Figure 3 presents the feature space before and after the drift has occurred. Initially, the feature space is depicted before cluster 4 appears (Fig. 3a). Then, cluster 4 is detected as an outlier micro-cluster since its weight is still low at this point (Fig. 3b). After another event arrival (Fig. 3c), cluster 4 increases its weight, making the region more dense, and steps to a potential micro-cluster. Finally, on event 2532 (Fig. 3d), cluster 4 is promoted to a core micro-cluster, i.e., a new common behavior. This way, a drift has been detected. With that, the drift detection can be seen from various complementing perspectives (console output, metrics evaluation and feature space evolution).

<sup>2</sup><https://ieee-dataport.org/open-access/synthetic-event-streams>



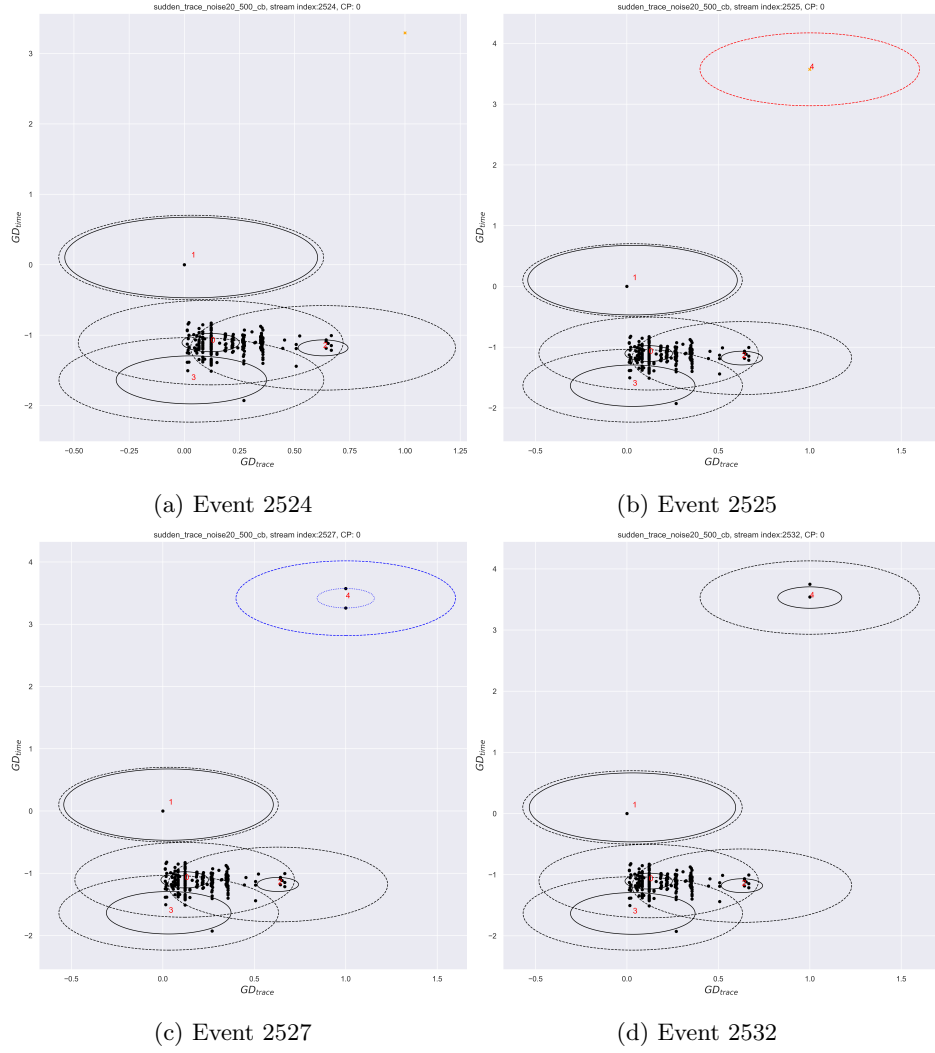


Figure 3: Evolution of the feature space detecting a concept drift.

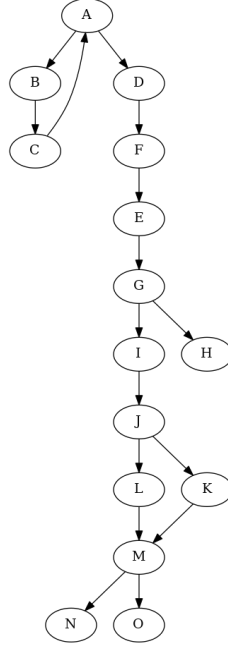


Figure 4: PMG of the synthetic event log at the first check point.

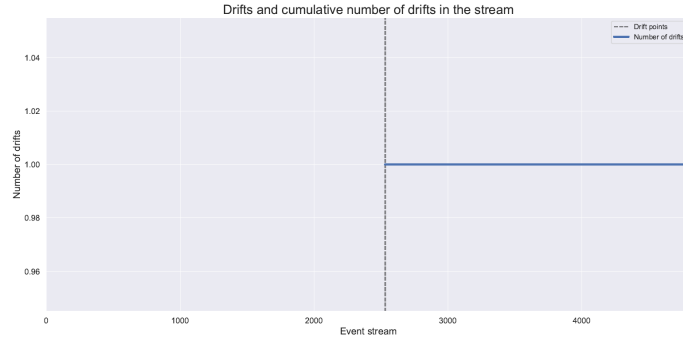


Figure 5: Relation of the event stream processing, number of detected drifts and their positions. Vertical lines represent the drift positions and the curve shows the cumulative number of drifts, which is 1 in this case.