# Autonomous Abnormal Behaviour Detection in Intelligence Surveillance and Reconnaissance Applications

R. Meo, R. Esposito, M. Botta, S. Viola, C.M. Choor, V. Mellano and F. Ciaramaglia

*Abstract*—**This paper describes a module that extracts rules or frequent patterns through data mining from a large database fed by targets detected by a Mission System installed on an unmanned airborne platform and the associated ground station to discover anomalies in local traffic. It has been demonstrated that the module is able to detect all tracks or targets present in the ground truth and also the paths followed by each tracks. Traffic anomalies can be detected by observing differences in extracted rules in reference missions compared to the current mission. The module will significantly reduce the operator workload as it can operate autonomously..**

*Index Terms* — **Aerospace Engineering, Command and Control Systems, Data Analysis, Data Mining, Knowledge Discovery**,

## I. INTRODUCTION

With evolving technologies and increase in the capability of sensors installed on an airborne platform, the amount of data being collected is huge and this creates a challenging task to the human operator in the ISR applications, i.e., for interpreting and processing the data into information in a quick and timely manner.

This paper describes a module that tries to help the operator in the above tasks attempting to discover, with autonomous capabilities, current relationships among data and events in the context of their environment.

The module uses data mining techniques to discover anomalies in the local traffic formed by targets. Targets observed data are collected and analyzed by a mission system installed on an unmanned airborne platform and the associated ground station. The details of the system are described in Section III.

Through the described tests, the paper will show the importance of a timely assessment of the situation with the scope to provide a timely response.

Rosa Meo, Marco Botta and Roberto Esposito are with Università degli Studi di Torino, Italy (meo@di.unito.it, botta@di.unito.it, roberto.esposito@unito.it)

Valter Mellano, Chee Ming Choor, Franco Ciaramaglia and Sergio Viola are With SELEX-ES, Italy

(walter.mellano@selex-es.com, C.M.Choor,cheeming.choor@selex-es.com, franco.ciaramaglia@selex-es.com, sergio.viola125@gmail.com)

## II. LITERATURE REVIEW

From the state of the art of data mining some previous approaches report results similar to this proposal but almost all of them, in contrast with our setting which is unsupervised, deal with supervised learning problems. The works in [1], [2] are an overview of the field of Data Mining and Machine Learning for knowledge discovery. The paper in [1] presents some results in the field of reinforcement learning for manufacturing and automated robot exploration and learning. The work in [2] deals with learning in a supervised setting in the context of a huge amount of data. The authors of [3] present algorithms and data structures for efficiently extracting some statistics on machine learning data-sets: this approach could be applied also to deriving the support of conjunctive rules. The work in [4] employs agents for information retrieval, rational communication and negotiation with applications in knowledge discovery. None of these articles employed a declarative query language for knowledge discovery with the exception of the study in [5] that uses a deductive logical language for the description of the domain and the extraction of frequent patterns. The works presented in [6] and [7] are extensions of deductive database systems with a mechanism for temporal, non-monotonic and non-deterministic reasoning. The early work in [6] is largely extended to extract spatio-temporal trajectories of vehicles from GPS in [10]. The works in [8] and [9] represent other examples of applying deductive and inductive logical languages to infer the class of observations in census spatial data and describe them by association rules.

In conclusion, at the best of our knowledge there are no previous works that attempt to extract knowledge in the form of association rules with query languages and try to generate new queries autonomously with the support of a domain ontology.

## III. DESCRIPTION OF THE MODULE

The module here described is part of the mission system of an Unmanned Aircraft System. A typical mission system for Unmanned Aircraft Systems (UAS) is split into two sections: one is installed onboard the aircraft and the other on the ground platform. The first one has two different kind of functions: the ones aimed to the conduct of the platform and the ones in charge of the peculiar mission execution, i.e, responsible of the sensing, mission command & control, situation awareness and interoperability. The ground platform facilitates the control and monitoring of the air platform and the exploitation of the information provided by the air section. Figure 1 shows the

architecture of the module which is part of both the air and ground sections. The component on the ground station works off-line for the rule extraction step (creation of reference) as it has a high computational workload. The second component is on board. It has the goal to trigger an alarm when different rules are found compared to the reference.
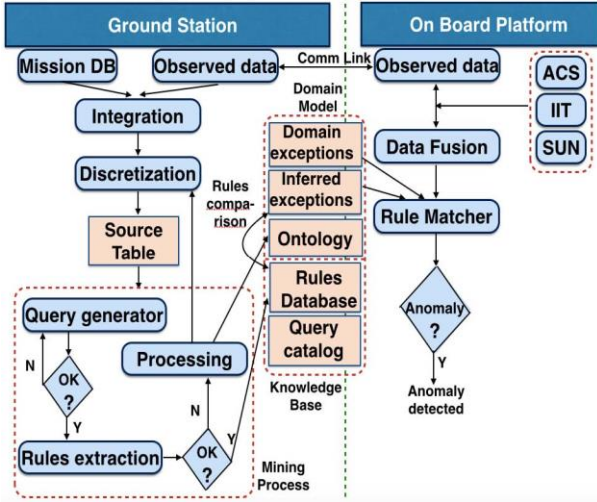


Fig. 1. Overview of the system

Both components receive the data collected by the airborne mission system and related to the targets located in the ground scenario: targets position, speed, direction of displacement, radiated frequency spectrum, target classification, etc... More in detail, this data is captured by the sensing capability that includes the sensors and their management and it is specific to the mission. For example in a high range surveillance mission, the suitable sensors will be the radar (with the SAR-Synthetic Aperture Radar and GMTI-Ground Moving Target Indicator sub modes) and the ESM-Electronic Support Measure, while in a close range mission, the suitable sensors will be the EO-Electro Optical or the IR-Infrared depending from the day or night conditions. Another important capability of the airborne mission system is the situation awareness that includes the sensor data exploitation and the data fusion. The data fusion will be required if the needed mission performances cannot be reached directly by the sensor; in this case data from different sensors have to be fused to reach the goal. For example if the target geolocation accuracy supplied by the ESM does not reach the mission goal of 50 m, a fusion between ESM and radar data is necessary.

The rules composed by frequent patterns are extracted by mining queries that fulfil the user's requests. The queries are expressed in a specialized, declarative language allowing the specification of constraints and the selection of rule attributes that vary according to the user's analysis goals and specific scenarios (an example is reported in a following subsection). The main components of the module are the followings.

### Mining Process

A mining query is generated automatically by the mining process component. The autonomous generation of new queries works by refinement of an initial set of queries provided by the user as examples of requests for acquiring knowledge from the observed events. The refinement of a query by the mining process module occurs by addition or modification of the constraints or of the rule attributes by selection of attributes describing the domain concepts at lower levels of the concept taxonomy.

### Knowledge base

Once a mining query is executed, the history of the system is updated in the query catalogue and the extracted knowledge is integrated into the knowledge base in the rule database.

The concept taxonomy (see the related subsection) is part of the domain ontology which is compiled by the domain expert and is part of the knowledge base. It constitutes a high-level description of the observed domain.

### Rule Comparator

The analysis of the acquired knowledge occurs by comparison with a referential time frame. The rule comparator works off-line and compares two sets of rules returned by the same mining query but in different time frames. The rules that were absent in the original time frame but were generated subsequently constitute examples of exceptions: when such a rule is matched on new data it triggers the presence of an anomaly. The time frame definition is highly dependent on the specific domain and on the user/analyst experience. In our experiments we set the mission observation time frame equal to 4 consecutive hours (either in the morning or in the afternoon). This is assumed to be the minimum observation time necessary to obtain an overview of the usual traffic. However, in dependence on the dynamics of the traffic a history composed by more consecutive missions (on more days) could be adopted to define the usual traffic.

In the experiments we decided that the reference model was defined by a previous mission (4 hours long) and that it was sufficient in a following mission to extract a set of patterns which had a statistically significant difference with the previous one to determine the decision about the existence of a discrepancy between the reference model and an anomalous situation.

For the definition of the statistically significant difference we applied the Kullback-Leibler divergence [11] that detects a difference between two probability distribution functions P (x) and Q(x) as follows:

$$D_{KL}(P,Q) = D_{KL}(P(x)\,||\,Q(x)) \qquad (1)$$

$$= -\sum_x P(x)\log Q(x) + \sum_x P(x)\log P(x) \qquad (2)$$

In the context of our work, the two probability distribution functions $P(x)$ and $Q(x)$ are the distributions of the patterns extracted respectively by the two missions (the reference one and the current one) with value $x$ of the attributes observed in the patterns. For instance, the attribute whose values would be easily observed in the patterns could be the location of the targets (equal to the spatial grid cell identifier) or the pair of

location and time stamp of the targets. Please refer to the paragraph *Mining query example* for examples of patterns.

In order to detect if the difference between the two pattern distributions is statistically significant we applied the Heoffding's inequality [12] for the determination of the threshold value of the difference, $\delta$ :

$$D_{KL}(P,Q) > \delta = \sqrt{-\frac{\log(\alpha)}{2n}} \cdot C \qquad (3)$$

where $\alpha$ is the level of confidence assumed for the confi¬dence interval (usually set to 0.01 or 0.005), $n$ is the number of values $x$ and $c = \log(\frac{b}{a}) - \log(\frac{a}{b})$ with a and b respectively the minimum and maximum values of the distributions.

### Rule Matcher

The rule matcher instead works on-line. Its goal is to match particular rules provided by the expert that describe the domain exceptions: once an exception rule is matched on the data it triggers an alarm. For this type of rules, an alarm is raised on the targets movements whose data match the rule completely (they are 100% matches).

### Spatial and Temporal Discretization

It might occur that the source data has been prepared in such a way that the chosen spatial and temporal discretization does not allow to find a satisfactory set of rules. Therefore, the module of data processing requires a change of focus and a new discretization for the spatial or temporal dimension is recalculated that changes the granularity with which the observations are stored and the knowledge is later extracted. Similarly, the choice of the statistical thresholds for the extraction of valid and frequent rules can be changed and can be automatically inferred by the system starting from the statistics of the occurrences of the rule attributes in the analyzed targets.

### Mining query example

One of the main queries we tested to monitor the target movements on the territory is a mining query that asks the rules associating subareas visited by a consistent number of targets in their traffic paths. In order to evaluate the occurrence of the paths in time (and determine in addition their direction) the query asks the specification in the rules together with each subarea visited by a target, of the time of the visit.

An example of the extracted rules is:

*[a44,8:50],[a44,9:00],[a35,9:10] -> [a33,9:20]*

where *aij* are identifiers of subareas in the spatial dimension followed by the time stamp representing the start of the time interval in which the targets visited the subarea. In the case of the example a time interval is 10 minutes long. The executed mining query is:

MINE RULE FrequentPaths AS

SELECT DISTINCT 1..n subarea id, d time AS BODY, subarea id, d time AS HEAD,

SUPPORT, CONFIDENCE

WHERE BODY.d time<HEAD.d time

FROM sourceData

GROUP BY track id

EXTRACTING WITH SUPPORT: 0.01, CONFIDENCE: 0.4

The meaning of the query is the following. MINE RULE clause introduces the name of the table used by the system to store the extracted rules. SELECT clause introduces the data structure of the rules specifying the attributes whose values compose the rules. In the case of this query both the antecedent part of the rules (named BODY) and the consequent one (HEAD) are composed by the values of pairs of attributes: subarea_id that represents the identifier of the visited subarea by a target and d_time that is the discretized time of the visit. The option DISTINCT 1..n specifies that the rules are composed by a varying number of distinct pairs (subareas visited by the same targets in the time instants). SUPPORT and CONFIDENCE are the statistics that will be used to evaluate the rules. FROM clause specifies the input table that contains the data on the targets movements; GROUP BY clause specifies the attribute used to partition the input table into groups of rows. In this case it is track_id that identifies the targets. The rules will be generated from the content of the groups and they need to be present in at least a minimum number of groups (targets). EXTRACTING clause specifies the minimum thresholds of the evaluating statistics (support and confidence). Notice that in this case the rules need not to have 100% confidence (that would represent the logical consequence rules). These are rules with statistical validity, suitable to catch the regularities that occur with some uncertainty in a noisy environment.

The above query has a constraint (WHERE BODY.d_time < HEAD.d_time) that requires that the subareas in the antecedent part of the rules must have a value of the discretized time that precedes the time value of the consequent.

In these examples the rules serve as a natural way to monitor the visits of the targets in terms of the carried out trajectories. However, the adoption of a query language allows flexibility to the specification of the input data and of the constraints that need to be satisfied by the observed entities allowing a certain degree of flexibility to the descriptive capabilities of the patterns. For instance, the query language allows the selection of more and diverse attributes that could satisfy the needs of the user/analysts, such as the vehicle type, its speed, the terrain type, etc. At the same time the query language allows generality to the underlying module of rules extraction that will continue to perform its work disregarding the selected attributes.

### A. Taxonomy of the concepts

The autonomous query generation could specialize the above query with addition of a constraint that asks to extract the association rules that associate subareas in which the targets have been observed on fields (an opposite situation w.r.t. the more common situations of the streets). The new mining query would be:

MINE RULE FrequentPathsOnFields AS

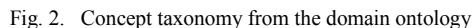SELECT DISTINCT 1..n subarea id, d time AS BODY, subarea id, d time AS

HEAD,

SUPPORT, CONFIDENCE

WHERE BODY.d time<HEAD.d time

AND BODY.target path='fields' AND HEAD.target path='fields'
FROM sourceData

GROUP BY track id

EXTRACTING WITH SUPPORT: 0.01, CONFIDENCE: 0.4

The above new query has an additional constraints (here highlighted in bold: BODY.target_path='fields' AND HEAD.target_path='fields') that impose that the new rules will regard traffic paths on fields (for the subareas of both the BODY and the HEAD).

The generation of the above query corresponds to browsing in the concept taxonomy following the highlighted path (shown in bold in Figure 2). The query asks to monitoring the events occurring in the space and time dimensions where the values of the time dimension already correspond to the lowest level of abstraction (time of the day) while the values of the space dimension correspond to a level of abstraction that has been changed and moved from the generic subarea level to a more specific level (namely, the specific subareas with fields).

## IV. PERFORMANCE EVALUATION

The approach we followed to evaluate the performances is focused on measuring, through a simulated environment, the module's ability to discover anomalies within the identified scenario. Here we list the main steps of the approach.
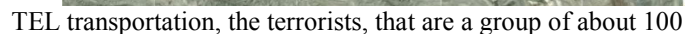
We prepared a data set composed by different test cases. The data represents the products collected during missions conducted by an unmanned airborne platform and then



Fig. 2.   Concept taxonomy from the domain ontology

processed by a Data Fusion system, that in this phase is simulated, having the scope to improve the attributes' accuracy of the collected products.

The collected data are given as input to the module that acts as if it was on board of the platform for discovering the anomalies embedded in the scenario.

### Test scenario

The scenario used for testing the effectiveness and the performance of the system is a typical large land area covering different subareas with varying population densities. The area that is selected for the analysis has s size of 2500 sq. km with the presence of hills and mountains.

It is a situation in which terrorist organizations have gained access to biological/chemical agents. These materials could be deployed as a Weapon of Mass Destruction and may pose a real threat to the International Community.

For launching the weapon, biochemical agents have to be obtained from research or production facilities, integrated into the launch missile, reloaded onto the Transporter Erector Launchers (TEL) and transported to a suitable launch site. With the scope to reduce the possibilities to be discovered during the



TEL transportation, the terrorists, that are a group of about 100

Fig. 3.   Map of the scenario

persons, prefer to assemble the launcher at the final launch site. The site has been selected in a inhabited area reachable only through the fields and not through the streets. Furthermore, the assembling activities are conducted at different times. Figure 3 represents the map of the area with the indicated zones and streets where the most part of the traffic is passing through.

### Results

Different kinds of tests have been performed on outcomes from 2 different missions (Mission 001 -conducted when activities were not present in zone 5 and 003 -conducted after starting of terrorists activities); they are discussed in the following.

### Module capability to detect characteristics of the traffic

Table I and Table II show the results obtained during the tests on Mission 001 and Mission 003 respectively and related to the capability to detect traffic characteristics. Not all the local areas are tested, but only a significant subset of them; in particular we tested the area where different traffic behaviour was included (zone 5) and another one where no significant changes occurred

(zone 3). In the first columns the Mission data (ID and date), the day of the week, the zones and its subarea components are reported. Then the ground truth data and the related results are placed side by side to allow an easy comparison. The following observations are in order:

**a)** The tracks or targets detected by the module reflect those present in the ground truth scenario, for both Missions 001 and 003. For the cases where a small difference exists ( 2 targets), the conducted analysis has shown that the reason is due to the used scenario, see next points d) and e).

**b)** Also the paths, constituted by the streets covered by tracks, have been correctly detected.

**c)** In general, the detected number of directions of movement is slightly higher than those of the ground truth scenario; the reason is due to a different mechanism of detection of the main direction. In fact the ground truth indicates only the main direction, while the module in the actual implementation indicates more precisely also the partial components of the directions.

**d)** In zone 5 and during the first hour (8:00 to 9:00) of Mission 003, the detected targets are 12 instead of 14. By analyzing the ground truth scenario, it has been noted that 2 targets transit trough zone 5 for a small time period (less than 2 minutes). The reason of the 2 missing detections is likely due to the discretization time used for the tests, equal to 10 minutes.

TABLE I.    RESULTS FROM MISSION 001

| Zone | Subareas | Ground Truth Tracks/hour | Direction | Street | Found Results |
|------|----------|--------------------------|-----------|--------|---------------|
| $z_1$ | $a_{25}, a_{34}, a_{35}$ | 1 (h:8:00-9:00) | SN | sec.ary,fields | |
| $z_3$ | $a_{51}, a_{52}, a_{42}, a_{43}$ | 7 (h:8:00-9:00) | WE | main | |
| | $a_{33}, a_{34}, a_{35}, a_{45}$ | 5 | EW | | |
| $z_4$ | $a_{42}, a_{43}$ | 4 (h:8:00-9:00) | WE, EW | main,sec.ary | |
| $z_5$ | $a_{33}, a_{34}, a_{35}$ | 1 (h:8-9:00) | | main,sec.ary | none |
| | $a_{44}, a_{45}, a_{54}, a_{55}$ | 0 | | fields | none |
| $z_3$ | $a_{51}, a_{52}, a_{42}, a_{43}$ | 7 (h:9-10:00) | WE | main | 16 tracks in 119 rules in all s_areas |
| | $a_{33}, a_{34}, a_{35}, a_{45}$ | 8 | EW | | |
| $z_5$ | $a_{33}, a_{34}, a_{35}$ | 1 (h:10-11:00) | SN | sec.ary,field | none |

**e)** In zone 3, for both the Missions 001 and 003, the detected targets number is higher than the targets present inside zone 3. The reason is that some subareas in zone 3 also belong to other zones of the scenario (zone 1 and 5), so the module can also capture other targets.

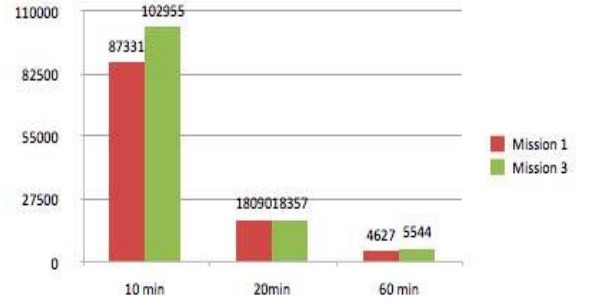TABLE II.    RESULTS FROM MISSION 003

| Zone | Subareas | Ground Truth Tracks/hour | Direction | Street | Found Results |
|------|----------|--------------------------|-----------|--------|---------------|
| $z_1$ | $a_{25}, a_{34}, a_{35}$ | 1 (h:8:00-9:00) | NS, SN | main,sec.ary | none |
| $z_3$ | $a_{51}, a_{52}, a_{42}, a_{43}$ | 8 (h:8:00-9:00) | WE | main | 12 tracks in >3 s_areas |
| | $a_{33}, a_{34}, a_{35}, a_{45}$ | 2 | EW | | |
| $z_4$ | $a_{42}, a_{43}$ | 3 (h:8:00-9:00) | WE, EW | main,sec.ary | none |
| $z_5$ | $a_{33}, a_{34}, a_{35}$ | 1 (h:8-9:00) | N | main,sec.ary | 12 tracks in >3 s_areas |
| | $a_{44}, a_{45}, a_{54}, a_{55}$ | 13 | SN,WE,EW | fields | |
| $z_5$ | $a_{33}, a_{34}, a_{35}$ | 3 (h:9:00-10:00) | NS | main | 2 tracks in >3 s_areas |

After the conduction of Mission 001 and Mission 003 and the processing of data by the module, zone 5 appears an area where an anomalous traffic behaviour is present. It corresponds to traffic that moved on fields (as opposite to streets). We discovered this fact by construction of the rule set that is composed by rules that appear in Mission 003 but not in Mission 001 (we used the rule comparator component). Next, by applying the Rule Matcher component we found the targets that satisfy at least one of those rules. The set of those targets constitutes a new source data set (that we call here AnomalousTargets) that is further used by a new, more specific query (generated autonomously; see next section).

**Performance of the module as a function of the time discretization, support and confidence values**



Number of rules obtained from samples with different discretization times

**f)** Figure 4 plots the number of rules extracted on data

Fig. 4.   Number of extracted rules in Mission 001 and Mission 003 with different time discretization.

sampled with different time discretization. It shows that the number of rules increases when the time granularity is decreased. However, such rules generally have a lower support value and a reasonable low number of rules is necessary to avoid a high computational effort and to focus on real interesting information.

**g)** Figure 5 plots the number of extracted rules when the support threshold value indicated in the query is lowered (re-ported as the value of the X axis). The support and confidence used in the tests (example of Mission 003) were set to very low values allowing the detection of almost any target so that we could verify the ability of the module to recognize the ground truth.



Fig. 5.   Number of extracted rules in Mission 003 with a changing support threshold.

**h)** Figure 6 plots the execution times for extracting the rules in Mission 003 when the support threshold value indicated in the query is lowered (reported as the value of the X axis). We observe that the computation time increases with the number of rules. This is most apparent when the number of rules grows in

the tens of thousands. In fact, in other cases the computation time is dominated by the database I/O operations.



Fig. 6. Computational times for the rules extraction component with a changing support threshold.

**Autonomous capability of subsequent queries preparation)**

After discovering a traffic anomaly in zone 5, by exploiting the concept taxonomy, the module is able to prepare other queries, in order to better understand the anomaly characteristics. Some possibilities are further investigated:

1) the hours of the day in which the anomalies occur,
2) the days of the week in which the anomalies occur,
3) the target types that compose the anomalous traffic.

Let us consider the specialization of the query described at point 3. The following query requests the list of the visited subareas (in the rule antecedent) associated to the vehicle type (the rule consequent) common to a consistent percentage of traffic paths (10%) that has been inferred as an anomaly with respect to the reference. Notice that the attribute moving vehicle comes from the event taxonomy and describes the vehicle type.

```
MINE RULE Anomalies with vehicle types AS
SELECT DISINCT 1..n subarea id AS BODY, moving vehicle AS HEAD, SUPPORT,
CONFIDENCE
FROM AnomalousTargets
GROUP BY track id
EXTRACTING RULES WITH SUPPORT: 0.1, CONFIDENCE: 0.5
```

Among the results it returns the following rule:

[a44],[a45],[a55]->truck

Figure 7 plots the number of rules that denote the presence of a difference in the traffic behaviour in Mission 003 w.r.t. Mission 001. The data that match these rules, denote an inferred anomalous behaviour with respect to the reference. These data form a new, specific data source that could be analysed in more detail by submission of a new set of queries generated autonomously in the data mining process engine. Notice that the number of rules increases when the time discretization decreases the depth of the time interval (in the case of the example, it decreases the time interval from 20 minutes to 10

minutes). The number of corresponding rules goes from 354 to 552. This effect is due to the increased probability that some targets are observed in the same spatial area in a different time interval and it provokes the verification of an increased number of rules.
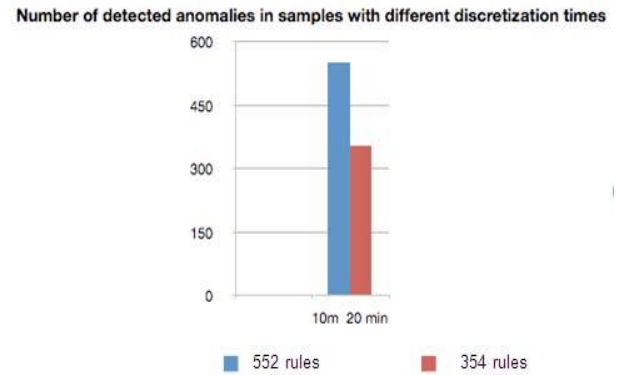


Fig. 7. plots the number of rules that denote the presence of a difference in the traffic behaviour in Mission 003 w.r.t. Mission

**Robustness of the module versus the errors present on the inputs**

We finally tested the robustness of the module with respect to the presence of errors in the detection of the spatial coordinates of the targets.

We tested the same mining query in two kinds of data: the first is the data coming from the Data Fusion system, that is characterized by a reduced error amount; the second is the data coming directly from the sensors which are affected by a larger error. Then, we compared the number of extracted rules, the number of targets that satisfy any of these rules and the main directions of target movement.

We note that the number of rules, as well as the number of targets, is the same in the two cases. A slight difference exists in the detection of the main directions of movement. This is due to the mechanism adopted in the detection of the main direction that contains a certain aleatory component in the cases in which a set of adjacent subareas is visited by the targets in the same time interval.

## V. CONCLUSION

The module has shown promising results that will be very useful in a persistent surveillance, a very important approach for the ISR applications. Results have proven that this module is able to detect traffic anomalies utilizing data detected by the Mission system. The autonomous capability of the module will also be essential to reduce the workload of the human Operator allowing him to concentrate more on the high level decision making tasks.

REFERENCES

[1] S. Thrun, C. Faloutsos, T. Mitchell and L. Wasserman, Automated Learning and Discovery: State Of The Art and Research Topics in a Rapidly Growing Field, *AI Magazine*, vol. 20, no. 3, 1999.

[2] F. Provost and V. Kolluri, A survey of methods for scaling up inductive algorithms. Data Mining and Knowledge Discovery, *Data Mining and Knowledge Discovery* , vol. 3, no. 2, pp. 131-169, 1999.

[3]  A.W. Moore and M.S. Lee, Cached Sufficient Statistics for Efficient Machine Learning with Large Datasets,' *Journal of Arti ficial Intelligence Research*, vol. 8, 1997.

[4]  L. Miller and J. Wong, Distributed Knowledge Networks, *Proceedings of the IEEE Information Technology Conference*, 1998.

[5]  L. Dehaspe and H. Toivonen, Discovery of frequent Datalog patterns, *Data Mining and Knowledge Discovery*, vol. 3, no. 1, pp.7-36, 1999.

[6]  F. Giannotti, M. Manco, M. Nanni, and D. Pedreschi, Datalog++: a Basis for Active Object-Oriented Databases, *A Basis for Active Object-Oriented Databases*, 1997.

[7]  N. Arni, K. Ong, S. Tsur, and C. Zaniolo, LDL++: A Second Generation Deductive Databases Systems, *ACM Transactions on Computational Logic*, 1993.

[8]  A. Appice, M. Ceci, A. Lanza, F.A. Lisi and D. Malerba, Discovery of Spatial Association Rules in Georeferenced Census Data: A Relational Mining Approach, *Intelligent Data Analysis*, vol. 7, 2003.

[9]  M. Ceci, A. Appice, and D. Malerba, Spatial Associative Classification at Different Levels of Granularity: A Probabilistic Approach, Knowledge Discovery in Databases: PKDD 2004, Lecture Notes in Artificial Intelligence, 3202, 99-111, Springer, Berlin, Germany, 2004.

[10]  R. Trasarti, F. Giannotti, M. Nanni, and D. Pedreschi, C.Renso, A Query Language for Mobility Data Mining, IJDWM vol. 7, no. 1, pp. 24-45, 2011.

[11]  S. Kullback, R.A. Leibler, On information and sufficiency, Annals of Mathematical Statistics vol. 22, no. 1, pp. 7986, 1951 doi:10.1214/aoms/1177729694. MR 39968.

[12]  L. L. Scharf, Statistical Signal Processing: Detection, Estimation, and Time Series Analysis, Pearson, 1991, ISBN-10: 0201190389.