



Anomaly Detection and Localisation using Mixed Graphical Models

Romain Laby, François Roueff, Alexandre Gramfort

► **To cite this version:**

Romain Laby, François Roueff, Alexandre Gramfort. Anomaly Detection and Localisation using Mixed Graphical Models. ICML 2016 Anomaly Detection Workshop, Jun 2016, New York, United States. <hal-01347167>

HAL Id: hal-01347167

<https://hal-imt.archives-ouvertes.fr/hal-01347167>

Submitted on 20 Jul 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Anomaly Detection and Localisation using Mixed Graphical Models

Romain Laby

ROMAIN.LABY@TELECOM-PARISTECH.FR

CNRS LTCI, Télécom ParisTech, Université Paris-Saclay, 46 Rue Barrault, 75013 Paris
Thales Airborne Systems, 2 Avenue Gay Lussac, 78990 Élan-court

François Roueff

FRANCOIS.ROUEFF@TELECOM-PARISTECH.FR

CNRS LTCI, Télécom ParisTech, Université Paris-Saclay, 46 Rue Barrault, 75013 Paris

Alexandre Gramfort

ALEXANDRE.GRAMFORT@TELECOM-PARISTECH.FR

CNRS LTCI, Télécom ParisTech, Université Paris-Saclay, 46 Rue Barrault, 75013 Paris

Abstract

We propose a method that performs anomaly detection and localisation within heterogeneous data using a pairwise undirected mixed graphical model. The data are a mixture of categorical and quantitative variables, and the model is learned over a dataset that is supposed not to contain any anomaly. We then use the model over temporal data, potentially a data stream, using a version of the two-sided CUSUM algorithm. The proposed decision statistic is based on a conditional likelihood ratio computed for each variable given the others. Our results show that this function allows to detect anomalies variable by variable, and thus to localise the variables involved in the anomalies more precisely than univariate methods based on simple marginals.

1. Introduction

Anomaly detection refers to the task of detecting anomalous samples within a dataset described by N variables, also called features. The localisation is the task that aims at identifying the subset of variables that are at the origin of the detected anomalies. While the problem of detection has been extensively studied in the machine learning literature (see (Hodge & Austin, 2004)), the problem of localisation in the presence of dependant variables remains a challenge.

In this paper, we propose to address this question using undirected probabilistic graphical models. Such models are

particularly useful to represent the joint distribution over a set of N random variables X_1, \dots, X_N , in an efficient and compact way. Undirected graphical models are commonly tied to Gaussian random variables yet recent works have studied the possibility of building models over heterogeneous variables : (Yang et al., 2014) proposes a general class of graphical models where each node-conditional distribution is a member of a univariate exponential distribution and (Lee & Hastie, 2015; Laby et al., 2015) investigate the problem of learning the structure of pairwise graphical model over both discrete and continuous variables. This is done by optimizing the likelihood or the pseudo-likelihood, penalized with a Lasso or group Lasso regularisation.

A standard approach to perform online anomaly detection on temporal data such as signals, is to use the CUSUM algorithm (see (Page, 1954) and (Basseville et al., 1993)). In this work we propose a two-sided test with an adapted CUSUM algorithm to detect anomalies that occur in the conditional distributions rather than in the marginal. The resulting algorithm allows to perform change-point detection and to detect, variable by variable, continuous or categorical, the time when the distribution of the data changes from the “normal” distribution.

2. Mixed Model Presentation

The definition of graphical models relies on the factorisation of the joint distribution. Pairwise models form a particular class of models where the features are grouped in sets of one or two variables. Such models have been widely studied and have a number of practical advantages (Schmidt, 2010). In this paper, we focus on mixed models mixing binary variables $X_C = \{X_i, i \in C\}$ (called categorical thereafter), and continuous variables $X_Q = \{X_u, u \in Q\}$ (called quantitative thereafter). We have

$X = (X_C, X_Q)$, with values in $\{0, 1\}^{|\mathcal{C}|} \times \mathbb{R}^{|\mathcal{Q}|}$. We use the pairwise mixed model

$$p_\Omega(x) = \frac{1}{Z_\Omega} \exp \left(x_C^T \Theta x_C + \mu^T x_Q - \frac{1}{2} x_Q^T \Delta x_Q + x_C^T \Phi x_Q \right), \quad (1)$$

where $\Omega = (\Theta, \mu, \Delta, \Phi)$ contains all the parameters of the model. Here, $\Theta = (\theta_{ij})_{i,j \in \mathcal{C}}$ is a symmetric matrix, $\mu = (\mu_i)_{i \in \mathcal{Q}} \in \mathbb{R}^{|\mathcal{Q}|}$, $\Delta = (\delta_{uv})_{u,v \in \mathcal{Q}}$ is a positive definite symmetric matrix and $\Phi = (\phi_{iu})_{i,u \in \mathcal{C} \times \mathcal{Q}}$ is a general matrix.

The model (1) is a mixture between the classic Ising Graphical Model (IGM) and Gaussian Graphical Model (GGM). In the Gaussian model, that is, when p takes the form of a Gaussian density, the partition function Z is easy to calculate and only requires the calculation of the determinant of a $N \times N$ matrix. The Ising model is one of the earliest studied undirected model for modeling energy of a physical system involving interactions between atoms (see (Ising, 1925)). The Ising model has binary variables, *i.e.* each x_i takes values in $\{-1, 1\}$ or $\{0, 1\}$, depending on the authors. Here we use the state space $\{0, 1\}$. The Ising model can be generalized for discrete variables, for example with the Potts model (Potts, 1953), but this one can be reparametrized as an IGM using 1-of- K encoding, as explained in (Bishop, 2006), §4.3.4. In the following, we will therefore only consider binary categorical variables.

To illustrate the model (1), we show some simulations made with 2 quantitative and 3 categorical variables. Figure 1 shows simulations of X when $\Phi = 0$ (the quantitative variables X_Q are independent of the categorical variables X_C and thus have a Gaussian distribution) and when $\Phi \neq 0$ (X_Q is not independent of X_C and its distribution is a mixture of Gaussian distribution). Given X_C , the conditional distribution of X_Q is always Gaussian, namely

$$X_Q | X_C \sim \mathcal{N}(\Delta^{-1}(\mu + \Phi^T X_C), \Delta^{-1}). \quad (2)$$

While, except when $\Phi = 0$, the unconditional law of X_Q is not Gaussian but is a mixture of Gaussian distributions, the unconditional law of X_C is again an Ising model with density

$$p_\Omega(x_C) \propto \exp \left(x_C^T (\Theta + \Phi \Delta^{-1} \Phi^T / 2 + \text{Diag}(\Phi \Delta^{-1} \mu)) x_C \right). \quad (3)$$

With these two properties, we can design an algorithm to efficiently sample from the distribution (1). Since $p_\Omega(X_Q, X_C) = p_\Omega(X_C) p_\Omega(X_Q | X_C)$, one just need to first sample X_C from (3), using for instance Wolff's algorithm

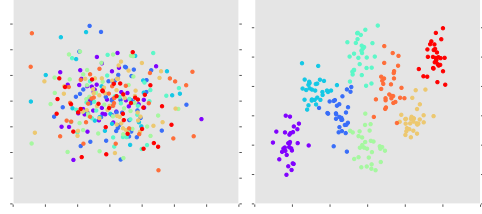


Figure 1. I.i.d. samplings of X_Q in two dimensions. The values of X_C are represented by 2^3 different colors. On the left, $\Phi = 0$ and on the right, $\Phi \neq 0$.

(Wolff, 1989), and then to sample X_Q from the conditional Gaussian density (2). This procedure will be used in our numerical experiments below.

In this paper, we do not aim at learning a graphical model but rather at exploiting one for anomaly detection and localisation. See (Yang et al., 2014), (Lee & Hastie, 2015) and (Laby et al., 2015) for recent works that investigate the task of learning the parameters of a mixed undirected graphical model.

3. Anomaly detection and localisation

In this section, we present a method to detect and localise anomalies from a sequence of new data $(X_C^{(t)}, X_Q^{(t)})$, $t = 1, 2, \dots$, assuming a reference model Ω that has already been learned using normal data.

The idea to localise anomalies is to monitor each term of the log-pseudo-likelihood (Besag, 1975) as a function of time. The CUSUM algorithm (Page, 1954) has been introduced to sequentially detect a change in the mean of a random variable. Since we want to detect an increase or a decrease, we use the two-sided CUSUM algorithm as proposed in (Basseville et al., 1993). For each t and each variable X_i , we define the instantaneous conditional log-likelihood ratio

$$s_i^{(t)} = \log \left(\frac{p(X_i^{(t)} | X_{-i}^{(t)})}{p_\Omega(X_i^{(t)} | X_{-i}^{(t)})} \right) \quad (4)$$

where $X_{-i} = \{X_j, j \in \mathcal{C} \cup \mathcal{Q}, j \neq i\}$, and a decision statistic defined recursively by $S_i^{(0)} = 0$ and

$$S_i^{(t)} = \left(S_i^{(t-1)} + s_i^{(t)} \right)^+, \quad t = 1, 2, \dots, \quad (5)$$

where $(z)^+ = \max(z, 0)$. Here p denotes the density of the alternative hypothesis, that is, the conditional density of the targeted anomalous behaviour.

We focus first on the quantitative variables. By (2), the conditional distribution of $X_Q^{(t)}$ given $X_C^{(t)}$ is the multivariate

Gaussian $\mathcal{N}(\nu^{(t)}, \Delta^{-1})$, with $\nu^{(t)} = \Delta^{-1}(\mu + \Phi^T X_c^{(t)})$. It follows that, for all $i \in \mathcal{Q}$, the conditional distribution of $X_i^{(t)}$ given $X_{-i}^{(t)}$ is Gaussian univariate with mean

$$e_i^{(t)} = \mathbb{E}_\Omega[X_i^{(t)} | X_{-i}^{(t)}] = \Delta_{i,-i}^{-1} \Delta_{-i,-i} \left(X_{\mathcal{Q}-i}^{(t)} - \nu_{-i}^{(t)} \right) + \nu_i^{(t)}$$

and variance

$$\sigma_i^2 = \text{Var}_\Omega(X_i^{(t)} | X_{-i}^{(t)}) = \Delta_{ii}^{-1} - \Delta_{i,-i}^{-1} \Delta_{-i,-i} \Delta_{-i,i}^{-1}.$$

We actually see from (2) that $e_i^{(t)}$ depends on $X_c^{(t)}$ and a fortiori on t , whereas it is not the case for σ_i .

For each quantitative variable X_i , $i \in \mathcal{Q}$, we want to detect a change in $p_\Omega(X_i^{(t)} | X_{-i}^{(t)})$. We define the conditional density $p(X_i^{(t)} | X_{-i}^{(t)})$ of the alternative hypothesis as a Gaussian density with same variance σ_i^2 and a modified mean $e_i^{(t)} + \delta\sigma_i$. The ratio (4) then becomes, for $i \in \mathcal{Q}$,

$$\begin{aligned} s_i^{(t)} &= \frac{1}{2} \left(\frac{X_i^{(t)} - e_i^{(t)}}{\sigma_i} \right)^2 - \left(\frac{X_i^{(t)} - (e_i^{(t)} + \delta\sigma_i)}{\sigma_i} \right)^2 \\ &= \frac{(X_i^{(t)} - e_i^{(t)})}{\sigma_i} \delta - \frac{1}{2} \delta^2. \end{aligned} \quad (6)$$

Setting $\delta > 0$ or $\delta < 0$ defines two statistics $S_i^{(t)\uparrow}$ and $S_i^{(t)\downarrow}$ in (5), for detecting respectively increase and decrease of the conditional mean $e_i^{(t)}$. In our experiments in Section 4, we will consider the sum $\bar{S}_i^{(t)} = S_i^{(t)\uparrow} + S_i^{(t)\downarrow}$ in order to detect a change in both possible directions.

Note that, by (6), the conditional negative drift under the null hypothesis (when no changes occur) of the decision statistic (5) is given by $\mathbb{E}_\Omega[s_i^{(t)} | X_{-i}^{(t)}] = -\delta^2/2$.

We focus now on the categorical variables. Each variable X_i , $i \in \mathcal{C}$ has a conditional Bernoulli distribution with mean

$$p_i = \mathbb{E}_\Omega[X_i | X_{-i}] = \frac{e^{q_\Omega(X,i)}}{1 + e^{q_\Omega(X,i)}}, \quad (7)$$

where

$$q_\Omega(X, i_0) = \theta_{i_0 i_0} + 2\Theta_{i_0, -i_0} X_{-i_0} + \Phi_{i_0, \mathcal{Q}} X_{\mathcal{Q}}. \quad (8)$$

In the case of categorical variables, we define the conditional distribution of the alternative hypothesis as a Bernoulli distribution with mean $a_i^{(t)}$. The instantaneous log-likelihood ratio is then given by

$$s_i^{(t)} = X_i^{(t)} \log \frac{a_i^{(t)}}{p_i^{(t)}} + (1 - X_i^{(t)}) \log \left(\frac{1 - a_i^{(t)}}{1 - p_i^{(t)}} \right). \quad (9)$$

We choose $a_i^{(t)}$ such as the drift of the decision function (5) under the null hypothesis is set to the same value $-\delta^2/2$,

as for quantitative variables in (6). This drift is given by computing $\mathbb{E}_\Omega[s_i^{(t)} | X_{-i}^{(t)}]$ with $s_i^{(t)}$ as in (6), yielding the equation

$$p_i^{(t)} \log \frac{a_i^{(t)}}{p_i^{(t)}} + (1 - p_i^{(t)}) \log \left(\frac{1 - a_i^{(t)}}{1 - p_i^{(t)}} \right) = -\frac{\delta^2}{2}. \quad (10)$$

It is easy to show that this equation in $a_i^{(t)}$ (with δ and $p_i^{(t)}$ fixed) has two distinct solutions $a_i^{(t)\uparrow} \in [p_i^{(t)}, 1]$ (associated to the statistic $S_i^{(t)\uparrow}$) and $a_i^{(t)\downarrow} \in [0, p_i^{(t)}]$ (associated to $S_i^{(t)\downarrow}$), detecting respectively increase and decrease of the mean $p_i^{(t)}$, with a conditional negative drift $-\delta^2/2$ under the null hypothesis. For the same reasons as with quantitative variables, we will consider the sum $\bar{S}_i^{(t)} = S_i^{(t)\uparrow} + S_i^{(t)\downarrow}$ in the experiments.

Under the null hypothesis each decision statistic $S_i^{(t)\uparrow}$ or $S_i^{(t)\downarrow}$ evolves with a negative drift $-\delta^2/2$. Hence, because of the positive part in (5), it remains close to zero with high probability. In contrast, under the alternative, the conditional drift becomes positive and the decision statistic $\bar{S}_i^{(t)}$ eventually increase above any arbitrarily high threshold h . We thus label as a change time the first times t when $\bar{S}_i^{(t)} > h$. The choice of δ sets how sensitive the test is to a close alternative, while the choice of h is a compromise between the false alarm probability over a given horizon and the delay needed to raise an alarm after a change of distribution. Finally and most interestingly, the set of indices i for which the alarm is raised provides a way to identify the variables for which not only the marginal distribution has changed but also the conditional one, given all other available variables.

4. Applications on synthetic data

In this section, we present results of anomaly detection and localisation with synthetic data. We suppose here that we have already learned the model parameters Ω from normal data. The data are composed of 50 normal observations sampled from the model using the algorithm explained at the end of section 2, and 50 anomalous observations sampled from an altered model where one parameter value in Ω has been modified.

We use the same model structure as in (Lee & Hastie, 2015), with 4 categorical and 4 quantitative variables. The model is represented in Fig. 2, with a colormap that will be kept for all experiments. The parameters have been chosen as follows: upper and lower diagonal of θ are filled with .5, $\theta_{i,i} = -\sum_{j \neq i} \theta_{i,j}$ and 0 elsewhere; $\mu_i = 0$, $\Delta_{i,i} = 1$, lower and upper diagonal of Δ is filled with .25 and 0 elsewhere; $\phi_{i,i} = .5$ and 0 elsewhere. We have tested three different modifications on the parameters of Ω : 1) the conditional distribution of the second (green) quantitative vari-

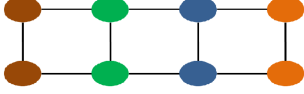


Figure 2. Structure of mixed graphical model used in the experiment. The upper and lower layer represents quantitative and categorical variables. Each i -th quantitative and i -th categorical variable have the same color.

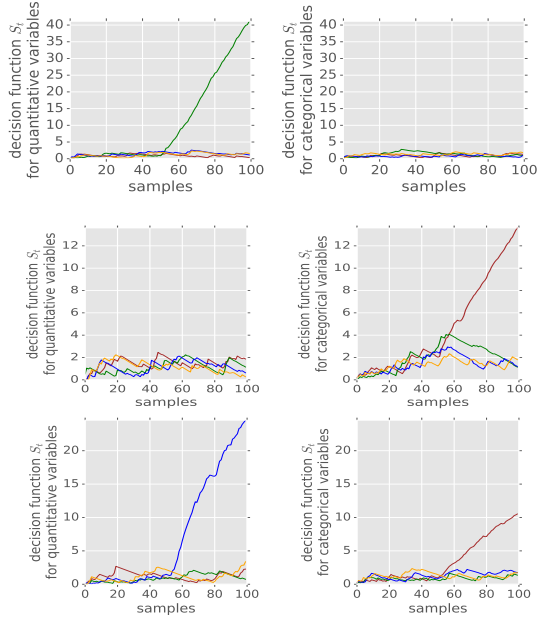


Figure 3. Time evolution of $\bar{S}_i^{(t)}$ for quantitative variables on the left and categorical variables on the right. The colors of the plots correspond to the colors of the variables in the graph of Figure 2. Top row : change on μ_1 . Middle row : change on $\theta_{0,0}$. Bottom row : change on $\phi_{0,2}$. For each experiment, the first 50 samples are sampled with parameter Ω , and the last 50 samples are sampled with the modified parameter.

able is changed by moving μ_1 from 0 to 3, 2) the conditional distribution of the first (red) categorical variable is changed by moving $\theta_{0,0}$ from -1 to -4 and 3) the conditional distributions of the first (red) categorical and third (blue) quantitative variable are changed by moving $\Phi_{0,2}$ from 0.5 to 2. Figure 3 shows the temporal evolution of the statistic $\bar{S}_i^{(t)}$ computed for every variable and for the three kinds of anomalies. As expected, the plots on the top row show that when changing μ_1 , only the statistic of the green quantitative variable is increasing, indicating that the green variable is carrying alone the change of conditional distribution. The same thing can be concluded for the two others modifications on $\theta_{0,0}$ and $\phi_{0,2}$. These results show that our method correctly detects and localises the changes in the conditional distributions.

We compare our method to the Wilcoxon test presented in (Lung-Yut-Fong et al., 2011), which is designed to detect changes in the distribution of a set of quantitative variables from batch data. In the following, we thus apply this approach to detect a change of distribution for each quantitative variable. Figure 4 displays the statistic of this test as a function of the possible change times. When only one change occurs in the data, this statistic is expected to approximately have a triangle shape with a maxima or a minima around the true change time. We use the same dataset as for the experiment with the anomalies localised on the second (green) quantitative variable, where μ_1 changes from 0 to 3 at time $t = 50$. Figure 4 should thus be compared with the top row of Figure 3. In contrast to online methods such as the one we propose, this Wilcoxon statistic cannot be computed recursively as it requires the whole set of data to be computed. Moreover it is not suited to localise the anomaly since a change of μ_1 , although it only modifies the conditional distribution of X_1 given X_{-1} , yields a change of all the marginal distributions. This is why in Figure 4, the Wilcoxon statistics display triangle shapes for all the quantitative variables with a more obvious change for the variables directly connected to X_1 .

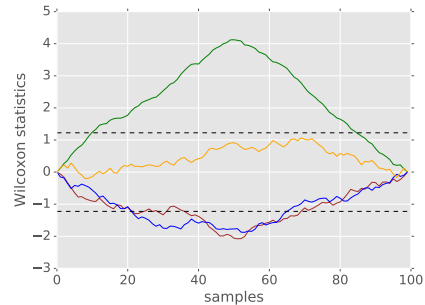


Figure 4. Evolution of the Wilcoxon statistic for 100 samples of 4 quantitative variables. After the 50th sample, we have modified Ω with $\mu_1 = 3$. The dashed lines indicate the thresholds for detecting a change with a 5% false detection probability.

5. Conclusion

In this paper, we proposed an online method that allows to detect anomalies in a data stream, but more importantly to localise which variables are at the origin of the problem. By using a mixed undirected graphical model learned over a set of normal data, we manage to track changes occurring in the conditional distributions which offers more specific detections than when studying only marginal distributions. This method is based on a two-sided CUSUM algorithm, where decision statistics are computed for every variable and involve the calculation of conditional likelihoods.

References

- Basseville, Michèle, Nikiforov, Igor V, et al. *Detection of abrupt changes: theory and application*, volume 104. Prentice Hall Englewood Cliffs, 1993.
- Besag, J. Statistical analysis of non-lattice data. *The Statistician*, 24, 1975.
- Bishop, C. *Pattern Recognition and Machine Learning*. 2006.
- Hodge, Victoria J and Austin, Jim. A survey of outlier detection methodologies. *Artificial Intelligence Review*, 22(2):85–126, 2004.
- Ising, E. Beitrag zur theorie des ferromagnetismus. *Zeitschrift fur Physik*, 31:253–258, 1925.
- Laby, Romain, Gramfort, Alexandre, Roueff, François, Enderli, Cyrille, and Alain, Larroque. Sparse pairwise Markov model learning for anomaly detection in heterogeneous data. June 2015. URL <https://hal-institut-mines-telecom.archives-ouvertes.fr/hal-01167391>.
- Lee, Jason D and Hastie, Trevor J. Learning the structure of mixed graphical models. *Journal of Computational and Graphical Statistics*, 24(1):230–253, 2015.
- Lung-Yut-Fong, Alexandre, Lévy-Leduc, Céline, and Cappé, Olivier. Homogeneity and change-point detection tests for multivariate data using rank statistics. *arXiv preprint arXiv:1107.1971*, 2011.
- Page, ES. Continuous inspection schemes. *Biometrika*, 41 (1/2):100–115, 1954.
- Potts, R.B. Some generalized order-disorder transformations. *Proc. Cambridge Philosophie Soc*, 1953.
- Schmidt, Mark. *Graphical model structure learning with l_1 -regularization*. PhD thesis, University Of British Columbia (Vancouver), 2010.
- Wolff, Ulli. Collective monte carlo updating for spin systems. *Physical Review Letters*, 62(4):361, 1989.
- Yang, Eunho, Baker, Yulia, Ravikumar, Pradeep D, Allen, Genevera I, and Liu, Zhandong. Mixed graphical models via exponential families. In *AISTATS*, pp. 1042–1050, 2014.