

Politehnica University of Bucharest
Faculty of Automatic Control and Computers
Department of Automatic Control and Systems Engineering

PhD Thesis Abstract

*Adaptive Strategies for Separable Dictionary
Learning and Applications to Anomaly
Detection*

Andra-Elena Băltoiu

Advisor:

Prof. dr. ing. Bogdan Dumitrescu

2021

1 Introduction

Dictionary Learning for sparse representations (DL) is a class of signal processing techniques that approximate signals by a linear combination of few elements of a basis, called the dictionary. The learning problem requires that both the dictionary and the sparse representation are learned from the data. The model is used for solving general signal reconstruction problems, for classification, various image-related tasks and in the field of compressed sensing. When it comes to 2D signals, such as images, however, the model requires that the data is vectorized, an operation that potentially breaks the correlations present in the second dimension. A solution to this drawback is the so called Separable Dictionary Learning (SDL) model, in which the dictionary is structured as the Kronecker product of two smaller (therefore more computationally efficient) dictionaries.

The primary focus of the thesis concerns the SDL model and, in particular, it deals with determining the optimal values of the two underlying heuristics: the sparsity of the representation and the dictionary size. Signal sparsity is the main assumption of the DL model and while it is rigorously proven only for some types of signals, practice has shown the suitability of such a representation in numerous (other) applications. The question of choosing the target sparsity of the model therefore remains open: there are limited possibilities of inferring the real sparsity of the data. The second heuristic is related to the number of elements in the dictionary (called the atoms), of which few are chosen to represent each signal. Numerical experiments [1] have shown that while larger dictionaries produce better signal approximations, this improvement levels out beyond a point and does not justify the additional computational demand.

In this work, we also adapt the SDL model to the problem of identifying graph anomalies. We are concerned with detecting abnormal topologies, i.e. (sub-)graph structures that stand out from the neighboring connectivity patterns. In such a case, the graph signals can be represented by their corresponding Laplacian matrices and we expect the dictionary atoms to describe elementary relational configurations. Our contribution also includes two other solutions for the anomaly detection problem, that take into consideration the requirements of most applications: the appropriateness of the method for unsupervised and online scenarios.

The dictionary learning problem is formalized as

$$\mathbf{Y} = \mathbf{D}\mathbf{X} + \mathbf{V}. \tag{1}$$

where $\mathbf{Y} \in \mathbb{R}^{m \times N}$ are the N signals, $\mathbf{D} \in \mathbb{R}^{m \times n}$ is the dictionary, having n atoms, $\mathbf{X} \in \mathbb{R}^{n \times N}$ is the representation matrix and $\mathbf{V} \in \mathbb{R}^{m \times N}$ represents the noise, which in most cases is considered to be Gaussian.

Both the dictionary and the representation are to be learned from the signals, leading to the following optimization problem

$$\begin{aligned}
& \min_{\mathbf{D}, \mathbf{X}} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2 \\
& \text{s.t. } \|\mathbf{x}_l\|_0 \leq s, l = 1 : N \\
& \quad \|\mathbf{d}_j\|_2 = 1.
\end{aligned} \tag{2}$$

Note that while \mathbf{D} can also be designed in advance and kept fixed, we do not treat this case and instead consider only the case in which the dictionary is learned. Such an approach has the advantage of producing a better suited dictionary for approximating the signals. The above problem is usually solved via an alternate minimization scheme, in which \mathbf{D} and \mathbf{X} are iteratively refined until a certain condition is met. The procedure starts by initializing a random dictionary and computing the representation, a step called sparse coding. In the next iteration, \mathbf{X} is kept fixed and the dictionary is updated. A common stopping condition is the number of learning iterations, as most algorithms for updating \mathbf{D} are known to achieve good approximation performance after sufficient iterations.

In separable dictionary learning, the dictionary is formed as $\mathbf{D} = \mathbf{D}_2 \otimes \mathbf{D}_1$, with $\mathbf{D}_1 \in \mathbb{R}^{m_1 \times n_1}$ and $\mathbf{D}_2 \in \mathbb{R}^{m_2 \times n_2}$. The model becomes

$$\mathbf{Y} = \mathbf{D}_1 \mathbf{X} \mathbf{D}_2^\top + \mathbf{V} \tag{3}$$

and it is known to be equivalent to (1), since

$$\text{vec}(\mathbf{D}_1 \mathbf{X} \mathbf{D}_2^\top) = (\mathbf{D}_2 \otimes \mathbf{D}_1) \text{vec}(\mathbf{X}). \tag{4}$$

2 Sparsity Bayesian Learning for Separable DL

The first constraint in (2) assumes the target sparsity s is known in advance, which is rarely the case in practice. Performance is highly dependent on the choice of sparsity, especially in some classes of applications, such as compressed sensing. We show that misestimating sparsity leads to either overfitting or underfitting. The experiment is run on a synthetic dataset generated where the real sparsity level is known and evaluates several models (i.e. models having different target sparsity values) by assessing the ratio between train and test errors. Results show that underestimating sparsity leads to similar test and train errors, suggesting improper model training, while overestimating s leads to overfitting.

Existing solutions alleviate these problems by either estimating s in a preliminary stage, by using a small set of signals to infer sparsity bounds, or employ adaptive strategies for inferring optimal sparsity as learning progresses. A third approach involves casting (2) as a Sparse Bayesian Learning (SBL) problem, which, in its hierarchical formulation, leads to a sparse solution without the need for explicitly setting a sparsifying prior. Consider the case where we are dealing with a single signal \mathbf{y} , instead of the entire signal matrix \mathbf{Y} , and its

corresponding representation \mathbf{x} . The basic, non-hierarchical, SBL approach considers the problem of learning \mathbf{x} as

$$\hat{\mathbf{x}} = \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}) = \arg \max_{\mathbf{x}} p(\mathbf{y}|\mathbf{x})p(\mathbf{x}). \quad (5)$$

This formulation leads to the Maximum a Posteriori (MAP) solution, which can be improved if instead of working directly on the elements of \mathbf{x} , a hyperparameter γ controlling the variance of each element is introduced

$$\hat{\gamma} = \arg \max_{\gamma} p(\gamma|\mathbf{y}, \sigma^2) = \arg \max_{\gamma} p(\mathbf{y}|\gamma, \sigma^2)p(\gamma). \quad (6)$$

As [2] has shown, a non-informative prior $p(\gamma)$ results in a sparse solution, therefore knowledge on the real sparsity is no longer required. Optimizing for γ leads to an estimation of the (Gaussian) posterior $p(\mathbf{x}|\mathbf{y}, \gamma, \sigma^2)$, having mean $\boldsymbol{\mu} = \sigma^{-2}\boldsymbol{\Sigma}_x\mathbf{D}^\top\mathbf{y}$ and variance

$$\boldsymbol{\Sigma}_x = (\sigma^2\mathbf{D}^\top\mathbf{D} + \boldsymbol{\Gamma}^{-1})^{-1}. \quad (7)$$

where $\boldsymbol{\Gamma}$ is a diagonal matrix containing the hyperparameters. The covariance matrix above is computationally demanding, because of its dimension. An alternative formulation, in terms of the smaller $\boldsymbol{\Sigma}_y$ is given in [3]

$$\boldsymbol{\Sigma}_y = \sigma^2\mathbf{I} + \mathbf{D}\boldsymbol{\Gamma}\mathbf{D}^\top. \quad (8)$$

The mean now becomes

$$\boldsymbol{\mu} = \boldsymbol{\Gamma}\mathbf{D}^\top\boldsymbol{\Sigma}_y^{-1}\mathbf{y}. \quad (9)$$

Estimation of γ can be performed following different approaches [2, 3]. We use the Expectation-Maximization (EM) formulation (k denotes the current iteration number):

$$\text{E step: } \mathbb{E}_{\mathbf{x}|\mathbf{y}, \gamma^{(k)}}[x_i^2] = (\boldsymbol{\Sigma}_x)_{i,i} + \mu_i^2, \quad (10)$$

$$\text{M step: } \gamma_i^{(k+1)} = \mathbb{E}_{\mathbf{x}|\mathbf{y}, \gamma^{(k)}}[x_i^2], \quad (11)$$

$$(\sigma^2)^{(k+1)} = \frac{1}{m} \left(\|\mathbf{y} - \mathbf{D}\boldsymbol{\mu}\|^2 + (\sigma^2)^{(k)} \sum_{i=1}^n \left(1 - \frac{\boldsymbol{\Sigma}_{x_{i,i}}}{\gamma_i} \right) \right)^{1/2}. \quad (12)$$

Our contribution (SBL-2D) [4] adapts the SBL framework to the separable DL problem and proposes a two-stage approach. First, the representation support is learned in a computationally efficient way, by changing the way the hyperparameters control the representation elements. Since we are dealing with two dictionaries, an element of \mathbf{X} is influenced by an atom from \mathbf{D}_1 (namely a column of the dictionary) and an atom from \mathbf{D}_2 (a row of \mathbf{D}_2^\top). We propose that the variance of the rows and columns of \mathbf{X} are controlled independently by two hyperparameters, $\boldsymbol{\beta}^{(1)} \in \mathbb{R}^{n_1}$ for rows and $\boldsymbol{\beta}^{(2)} \in \mathbb{R}^{n_2}$ for columns.

To adapt relations (7-8), we note that the matrix $\mathbf{\Gamma}$ is replaced by two matrices $\mathbf{B}^{(d)}$, one for each SBL process. Each $\mathbf{B}^{(d)}$ can be expressed as a Kronecker product, denoted generically $\mathbf{B}^{(d)} = \mathbf{B}_2^{(d)} \otimes \mathbf{B}_1^{(d)}$, with the particular forms

$$\begin{aligned}\mathbf{B}^{(1)} &= \mathbf{I}_{n_2} \otimes \text{diag}(\boldsymbol{\beta}^{(1)}), \\ \mathbf{B}^{(2)} &= \text{diag}(\boldsymbol{\beta}^{(2)}) \otimes \mathbf{I}_{n_1}.\end{aligned}\tag{13}$$

Here \mathbf{I}_{n_d} denotes the identity matrix of size n_d ; the first above relation simply says that all elements of row i of \mathbf{X} are associated with $\beta_i^{(1)}$. We introduce the notation $\boldsymbol{\beta}^{(d)}$, $d \in \{1, 2\}$ for referencing both dimensions and extend it to all other variables that have both row and column values.

Using known Kronecker product properties, adapting (8) for the 2D case reads

$$\boldsymbol{\Sigma}_y^{(d)} = (\sigma^{(d)})^2 \mathbf{I} + (\mathbf{D}_2 \mathbf{B}_2^{(d)} \mathbf{D}_2^\top) \otimes (\mathbf{D}_1 \mathbf{B}_1^{(d)} \mathbf{D}_1^\top).\tag{14}$$

Structured covariances arise in several multivariate contexts. Consequently, computing the inverse of covariances such as (14) is known [5] to be efficiently obtained using the SVD decomposition, in our case

$$\begin{aligned}\mathbf{D}_1 \mathbf{B}_1^{(d)} \mathbf{D}_1^\top &= \mathbf{U}_1^{(d)} \mathbf{S}_1^{(d)} \mathbf{U}_1^{(d)\top}, \\ \mathbf{D}_2 \mathbf{B}_2^{(d)} \mathbf{D}_2^\top &= \mathbf{U}_2^{(d)} \mathbf{S}_2^{(d)} \mathbf{U}_2^{(d)\top}.\end{aligned}\tag{15}$$

By simply using the orthogonality of the matrices from the above relations, it results that

$$\left(\boldsymbol{\Sigma}_y^{(d)}\right)^{-1} = (\mathbf{U}_2^{(d)} \otimes \mathbf{U}_1^{(d)}) \left((\sigma^{(d)})^2 \mathbf{I} + \mathbf{S}_2^{(d)} \otimes \mathbf{S}_1^{(d)}\right)^{-1} (\mathbf{U}_2^{(d)} \otimes \mathbf{U}_1^{(d)})^\top,\tag{16}$$

where now the matrix to be inverted is diagonal. We denote

$$\left((\sigma^{(d)})^2 \mathbf{I} + \mathbf{S}_2^{(d)} \otimes \mathbf{S}_1^{(d)}\right)^{-1} = \text{diag}(\mathbf{t}^{(d)}),$$

with $\mathbf{t}^{(d)} \in \mathbb{R}^m$. Using the above expression and denoting

$$\mathbf{E}_1^{(d)} = \mathbf{B}_1^{(d)} \mathbf{D}_1^\top \mathbf{U}_1^{(d)}, \quad \mathbf{E}_2^{(d)} = \mathbf{B}_2^{(d)} \mathbf{D}_2^\top \mathbf{U}_2^{(d)},\tag{17}$$

the solution estimation (9) becomes

$$\boldsymbol{\mu}^{(d)} = (\mathbf{E}_2^{(d)} \otimes \mathbf{E}_1^{(d)}) \cdot \text{diag}(\mathbf{t}^{(d)}) \cdot (\mathbf{U}_2^{(d)} \otimes \mathbf{U}_1^{(d)})^\top \mathbf{y}.\tag{18}$$

The E step is completed by transforming the expression of the posterior covariances (7) into

$$\boldsymbol{\Sigma}_x^{(d)} = \mathbf{B}_2^{(d)} \otimes \mathbf{B}_1^{(d)} - (\mathbf{E}_2^{(d)} \otimes \mathbf{E}_1^{(d)}) \cdot \text{diag}(\mathbf{t}^{(d)}) \cdot (\mathbf{E}_2^{(d)} \otimes \mathbf{E}_1^{(d)})^\top.\tag{19}$$

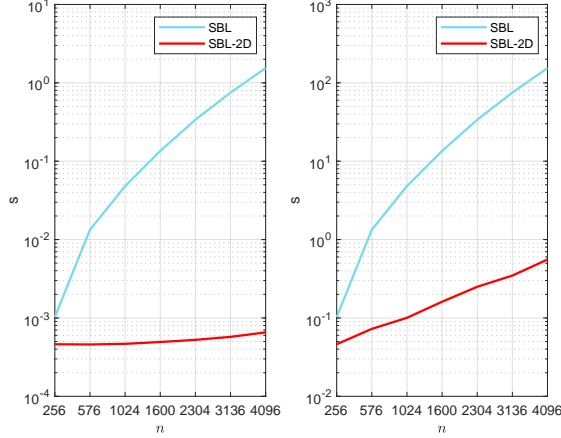


Figure 1: Execution time (in seconds) for SBL and SBL-2D. Left: time per iteration. Right: total time.

Until now, we have presented the ingredients for solving the E step of the EM algorithm. The M step amounts to updating the hyperparameter values and noise estimation. We adapt the step to our separable problem as

$$\beta_i^{(d)} \leftarrow \frac{1}{n_{3-d}} \sum_{\substack{\ell=1 \\ c_d(\ell)=i}}^n \left((\Sigma_x)_{\ell,\ell}^{(d)} + (\mu_\ell^{(d)})^2 \right). \quad (20)$$

Finally, the noise estimation (12) is performed separately for the row and column cases with

$$(\sigma^{(d)})_{(k+1)}^2 = \frac{1}{m_1 m_2} \left(\|\mathbf{y} - (\mathbf{D}_2 \otimes \mathbf{D}_1) \boldsymbol{\mu}^{(d)}\|^2 + (\sigma^{(d)})_{(k)}^2 \sum_{\ell=1}^n \left(1 - \frac{(\Sigma_x)_{\ell,\ell}^{(d)}}{\beta_{c_d(\ell)}^{(d)}} \right) \right)^{1/2}. \quad (21)$$

With identifying the representation support in this manner, we complete the first stage of our SBL-2D algorithm. The second step involves using the regular SBL method on the reduced support in order to estimate the representation. The approach improves the computational costs, since the problem is now considerably smaller. Figure 1 shows execution times for both algorithms. Testing is done using a synthetic dataset, where the real sparsity is $s = 5$ and signals are corrupted by noise having signal-noise ratio $\text{SNR} = 40$.

Results show that accuracy is comparable to regular SBL and often considerably higher than other sparse representation solutions adapted for the 2D case. We test performance in a set of experiments on synthetic data as well as real images. We vary the real sparsity level and noise values and measure the signal approximation error (RMSE) and representation error (by comparing the

true representation with the SBL-2D estimate). The contribution also includes convergence tests - showing the reduction of the representation support at each iteration, the evolution of the two error metrics, as well as false negative and false positive counts for the identified support.

3 Dictionary Size Adaptation for Separable DL

When sparsity is unknown and not estimated by one of the strategies mentioned above, it is usually set with respect to dictionary size. However, while at least for some types of signals, an informed choice for sparsity is possible, in what the dimension of the dictionary is concerned, such prior information hardly exists. In turn, common practice is to set the size of \mathbf{D} with respect to signal dimension. The alternative is to apply adaptive strategies for determining the optimal size.

Regardless of the optimality criterion used, all strategies involve starting with either a small dictionary or a large one and adding and/or trimming atoms as learning progresses, in order to obtain a better suited model (i.e. dictionary). More precisely, size optimality is evaluated once every few iterations and the dimension of the dictionary is modified accordingly. However, the size adaptation strategy can interfere with learning, since newly added atoms require some iterations to align to meaningful directions. In an experiment on synthetic data, we investigate this effect by evaluating signal approximation error at each model evaluation step and the corresponding optimal size, which we in turn compare to the true dictionary size.

One solution for evaluating the model is the Minimum Description Length (MDL) principle of model selection, which assumes the correlations present in the signals allow for a parsimonious representation. As such, an optimal model is one that can compress the data, while not hindering on signal approximation. The framework includes several Information-theoretic Criteria (ITC) that seek a compromise between model performance and model complexity. In particular, the Renormalized Maximum Likelihood (RNML) criterion is known to permit a formulation (called Extended-RNML) for the problem of determining the optimal size of the dictionary [6].

Our contribution involves the adaptation of ERNML to the separable model and corresponding size adjustment scheme [7]. Because of the equivalence of the original DL model (1) and the separable counterpart (3), the ERNML criterion does not require any essential reformulation in order to fit the separable case. The only necessary adaptations concern dimension adjustments and, given we are working with two dictionaries instead of one, parameter recount. In particular, signal dimension is $T_{2D} = m_1 m_2 N$ and the number of parameters is $\text{NoP}_{2D} = sN + (m_1 - 1)n_1 + (m_2 - 1)n_2$, since we are counting the degrees of freedom for each dictionary separately. Adapted to the 2D case, the criterion reads

$$\begin{aligned}
\text{ERNML}_{2D} &= (T_{2D} - \text{NoP}_{2D}) \log \frac{\text{RMSE}^2}{T_{2D} - \text{NoP}_{2D}} \\
&+ \text{NoP}_{2D} \log \frac{\|\mathbf{D}_1 \mathbf{X} \mathbf{D}_2^\top\|_F^2}{T_{2D} \cdot \text{NoP}_{2D}} \\
&+ \log [\text{NoP}_{2D}(T_{2D} - \text{NoP}_{2D})] + 2N \log \binom{n_1 n_2}{s}.
\end{aligned} \tag{22}$$

Adapting dictionary size via ITC in [6] involves routinely evaluating the model dimension-performance trade-off for a number n_{cand} of candidate models, more precisely for different sized dictionaries. Following Occam’s razor, the strategy is biased towards examining smaller dictionaries, which also has the benefit of keeping the complexity of computing the criterion in check. It does, however, account for cases where the optimal size is larger. The idea is, briefly, to test whether the current dictionary is not over-sized with useless atoms. Accordingly, every $iter_{adapt}$ iterations, the (current) dictionary atoms are ordered based on their representation power. The power of an atom j is defined as [1]

$$P_j = \|\mathbf{x}_j^\top\|_2^2. \tag{23}$$

The candidate models are formed by excluding the least powerful atoms, the smallest candidate having $n_{current} - n_{cand} + 1$ atoms (where, clearly, $n_{current}$ is the size of the dictionary at the present iteration). A representation is computed for each model and the corresponding ERNML criteria are evaluated. The best dictionary is that for which ERNML is smallest. The case where the current dictionary is the one yielding the lowest value is interpreted as an indication that the model is not complex enough and hence atom addition is performed.

When working with separable dictionaries, however, we refer to the representation power of an atom combination $\mathbf{D}_{1,i}$ and $\mathbf{D}_{2,j}$. More precisely, the measure now takes the form

$$P_{2D} = \sum_{k=1}^N X_{i,j,k}^2. \tag{24}$$

We therefore sort the atom combinations based on power and produce the corresponding ranked atom lists for D_1 and D_2 respectively. There are now n_{cand}^2 candidates, since at most n_{cand} atoms can be excluded from each dictionary. The criterion is then computed for every pair combination of dictionary dimensions. Size adjustment is done separately on the two dictionaries, by performing the necessary atom additions/removals that are indicated by the lowest ERNML value. The resulting best \mathbf{D}_d s may have different sizes, which lays the ground for a well adapted model, since the signal patterns in one dimension (say the rows of \mathbf{y}) may differ in size and complexity from the ones in the other (columns of \mathbf{y}), thus requiring a different number of atoms for their representation.

We also propose a solution that minimizes interference of dimension modification with learning, which involves relying on ERNML_{2D} only as an indicator of the change. The actual optimal size is computed by smoothing the optimal size values provided by the ITC criterion over the last ws iterations, in order to avoid abrupt size changes that may lead to under-trained atoms. The solution uses a moving average filter of window ws . The complete size adaptation scheme is presented in Algorithm 1.

Results show good performance in determining the true dictionary sizes for several dimension configurations, different sparsity levels and noise values. Most encouragingly, dictionary size is hardly underestimated, especially when the true size is small.

4 Anomaly Detection

The anomaly detection (AD) problem can be seen as binary classification, where one class represents the normal signals, while the second the anomalies. Our contributions use dictionary learning to solve the AD problem for three applications: malware identification, financial fraud detection and the more general problem of graph anomaly detection.

Malware detection is essentially a large scale problem, given the multitude of software applications. The computational costs involved in learning such a large set of signals can be alleviated by resorting to online algorithms that can be successively trained on parts of the original dataset, without significantly compromising on accuracy. Another aspect concerns the dynamics of the field, where new types of malware are constantly produced, to evade the advances in anti-virus solutions. Models must accommodate this adaptability, which translates into the ability of identifying malware for which no previous example is available. Unsupervised methods are, thus, preferable.

The majority of the above observations also hold for the problem of detecting financial frauds, perhaps even more so. Financial fraud poses an additional problem for machine learning. A transaction is best described as a link between two nodes - the financial entities. Therefore, a suited data representation is that of a graph, which has the advantage of the representing inter-dependencies of transactions. Usually, both nodes and edges are attributed with information such as concerning identity, transaction amount and currency. Therefore, data consists in both numeric and relational information. This underlying graph is especially relevant with money laundering, where individual transactions in the scheme may in themselves seem legitimate; it is only in the context of the graph which links them that fraud becomes apparent.

Motivated by the problem of identifying money-laundering schemes and other financial frauds, our work [8] reviews anomaly detection methods that are aimed at graph data.

The general DL framework can be extended for the classification task, especially in the supervised setting. The goal is to learn a dictionary such that different sets of atoms describe samples from each class. Moreover, it is also

Algorithm 1: ITC-ADL-2D

Data:

signals, $\mathbf{Y} \in \mathbb{R}^{m \times N}$
sparsity, s
number of DL iterations, K
iterations step for size adaptation, K_{adapt}
number of candidate sizes, n_{cand}
minimum number of dictionary atoms n_{min}
moving average window, ws

Result: optimum-sized dictionaries, \mathbf{D}_1 and \mathbf{D}_2

- 1 Initialize dictionaries \mathbf{D}_1 and \mathbf{D}_2 of sizes n_d
 - 2 **for** $k = 1 : K_{adapt} : K$ **do**
 - 3 Perform K_{adapt} iterations of 2D DL and obtain updated dictionary, $\mathbf{D}_1, \mathbf{D}_2$ and representation, \mathbf{X} with s target sparsity
 - 4 Order atoms separately in \mathbf{D}_1 and \mathbf{D}_2 according to power
 - 5 Compute $ERNML_{2D}$ criterion for all n_{cand}^2 candidates using (22) and return sizes $n_{d,ITC}$ that yield smallest $ERNML_{2D}$ value
 - 6 Compute $n_{d,opt}$ as moving average of $n_{d,ITC}$ over a ws window
 - 7 Apply $\mathbf{D}_d = \text{AdjustSize}(\mathbf{D}_d, n_{d,opt}, n_d)$, where $d \in \{1, 2\}$ for both \mathbf{D}_1 and \mathbf{D}_2
 - 8 Perform K_{adapt} more DL iterations with $n_{d,opt}$ sized dictionaries
- Function** $\mathbf{D}_d = \text{AdjustSize}(\mathbf{D}_d, n_{d,opt}, n_d)$
- 9 **if** $n_{d,opt} - n_d > 0$ **then**
 - 10 Set current adjusted size $n_d = n_{d,opt}$
 - 11 Add $n_{d,opt} - n_d$ atoms to dictionary \mathbf{D}_d
 - 12 **else if** $n_{d,opt} - n_d < 0$ **then**
 - 13 Set current adjusted size $n_d = \min(n_{d,opt}, n_{min})$
 - 14 Trim \mathbf{D}_d to its most used n_d atoms
 - 15 **else**
 - 16 Set current adjusted size $n_d = n_{d,opt} + 1$
 - 17 Add 1 atom to dictionary \mathbf{D}_d
-

beneficial that the representations of the signals in each class differ significantly from the ones in other classes. These two constitute additional constraints imposed on the learning problem and can be achieved by forcing \mathbf{X} to represent class labels and respect class atom allocation, along with the regular signal approximation objective.

Adding the three objectives together: signal reconstruction, discrimination (between classes) and label consistency (atom-class allocation), results in the following optimization problem

$$\min_{D, \mathbf{W}, \mathbf{A}, \mathbf{X}} \|\mathbf{Y} - \mathbf{D}\mathbf{X}\|_F^2 + \alpha \|\mathbf{H} - \mathbf{W}\mathbf{X}\|_F^2 + \beta \|\mathbf{Q} - \mathbf{A}\mathbf{X}\|_F^2. \quad (25)$$

where $\mathbf{H} \in \mathbb{R}^{c \times N}$ denotes the label matrix, $\mathbf{W} \in \mathbb{R}^{c \times n}$ the classifier matrix, $\mathbf{Q} \in \mathbb{R}^{n \times N}$ the atom allocation matrix, $\mathbf{A} \in \mathbb{R}^{n \times n}$ is a dictionary that imposes label consistency on the representation, and c stands for the number of classes.

The solution is called Label Consistent K-SVD (LC-KSVD) [9] since it can be recast as a regular DL problem and solved by standard algorithms such as K-SVD. Specifically, (25) is equivalent to

$$\min_{D, \mathbf{W}, \mathbf{A}, \mathbf{X}} \left\| \begin{bmatrix} \mathbf{Y} \\ \sqrt{\alpha} \mathbf{H} \\ \sqrt{\beta} \mathbf{Q} \end{bmatrix} - \begin{bmatrix} \mathbf{D} \\ \sqrt{\alpha} \mathbf{W} \\ \sqrt{\beta} \mathbf{A} \end{bmatrix} \mathbf{X} \right\|_F^2. \quad (26)$$

Online DL approaches to anomaly detection use several heuristics to control the confidence of classification. While this has indeed a positive impact on accuracy, it is not suited for the two applications we described. Both financial frauds and malware entail a high degree of novelty: it is to be expected that new fraudulent schemes are developed as old ones are exposed by anti-money laundering efforts and new types of malware are created as antivirus programs catch existing ones.

We propose a semi-supervised approach (called Tolerant Online Discriminative DL with Regularization - TODDLer), which provides a solution to this compromise. The setting involves an offline pre-training stage, where a dictionary is learned using a small labeled sample set. Regular classification methods for DL can be performed for this task, such as LC-KSVD described above. The dictionary is used to initialize the unsupervised online stage, where \mathbf{D} is updated by feeding the bulk of the dataset one sample at a time.

The following solution and results constitute our contribution published in [10]. The method is based on an existing method [11], but in order to make the solution better suited for malware detection, where unseen malware types need to inform the model for future use, we aim at using all signals to train the model. Improper labeling can misguide learning, therefore we wish to prevent signals from drastically modifying the current model, since there is no guarantee the change is for the best. The solution is to add regularization terms for the classifier and label consistency matrix that control their rate of change. Modified in this way, the label consistent classification problem (26) becomes

$$\begin{aligned} \min_{D, \mathbf{W}, \mathbf{A}} & \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2 + \alpha \|\mathbf{h} - \mathbf{W}\mathbf{x}\|_2^2 + \beta \|\mathbf{q} - \mathbf{A}\mathbf{x}\|_2^2 \\ & + \lambda_1 \|\mathbf{W} - \mathbf{W}_0\|_F^2 + \lambda_2 \|\mathbf{A} - \mathbf{A}_0\|_F^2. \end{aligned} \quad (27)$$

In the formulation above we have also adapted the online approach: instead of dealing with the entire signal matrix \mathbf{Y} , the objective is minimized with every incoming signal \mathbf{y} , that has a corresponding label vector \mathbf{h} and an associated atom allocation vector \mathbf{q} . By \mathbf{W}_0 and \mathbf{A}_0 we denote the current values of the dictionaries, learned based on the previous signals.

The updated values of \mathbf{W} and \mathbf{A} can be obtained by fixing everything else in (25) and solving the corresponding minimization objectives

$$f(\mathbf{W}) = \|\mathbf{h} - \mathbf{W}\mathbf{x}\|_2^2 + \lambda_1 \|\mathbf{W} - \mathbf{W}_0\|_F^2, \quad (28)$$

$$g(\mathbf{A}) = \|\mathbf{q} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda_2 \|\mathbf{A} - \mathbf{A}_0\|_F^2. \quad (29)$$

Setting the gradients of the above functions to 0 leads to the following least squares solutions

$$\mathbf{W} = (\mathbf{h}\mathbf{x}^T + \lambda_1 \mathbf{W}_0)(\mathbf{x}\mathbf{x}^T + \lambda_1 \mathbf{I})^{-1}, \quad (30)$$

$$\mathbf{A} = (\mathbf{q}\mathbf{x}^T + \lambda_2 \mathbf{A}_0)(\mathbf{x}\mathbf{x}^T + \lambda_2 \mathbf{I})^{-1}. \quad (31)$$

We test the solution on two malware datasets and a financial fraud dataset. Results show comparable or improved classification accuracy over other online DL methods for classification.

Our second contribution to AD is an unsupervised method that progressively filters out the normal signals [12]. It is based on the assumption that normal samples, which generally outnumber the anomalies, are better represented by the model, and we thus used an error criteria to differentiate between the two. The solution also takes into consideration the fact that, as filtering leaves less regular signals to train on, this imbalance attenuates. We propose a composite dictionary structure, in which the new models from each iteration are combined with existing ones, since overtraining the dictionary on anomalies would corrupt the error criterion.

We also experiment with using atom properties to inform the class labels. An atom can be characterized in terms of the number of signals which are represented using the atom, a property that is called usefulness, $U_j = \|\mathbf{x}_j^\top\|_0$. Another approach to progressively filtering signals involves the restriction of the set \mathcal{A} of potential anomalies to signals represented by atoms \mathbf{d}_j with usefulness $U_j < N_a$, namely

$$\mathcal{A} = \{ \mathbf{y}_k \mid x_{j,k} \neq 0 \wedge U_j < N_a, \forall j, k \}. \quad (32)$$

In other words, we expect atoms being used by less than N_a signals to represent the features of anomalies.

Both methods yield good results in lowering the false positive count, while keeping the false negatives in check and can thus be used for obtaining a smaller dataset with better class balance, since only normal signals are eliminated by filtering, if proper stopping conditions are set.

We also employ atom usefulness, together with the atom power measure in a multiscale dictionary setup to investigate the separability of normal and abnormal signals. Intermediate results show that using the standard DL formulation, without any added classification constraints, can still be informative on class labels if we use atom properties in order to characterize the signals. The experiments show that the distributions of atom usefulness and atom power

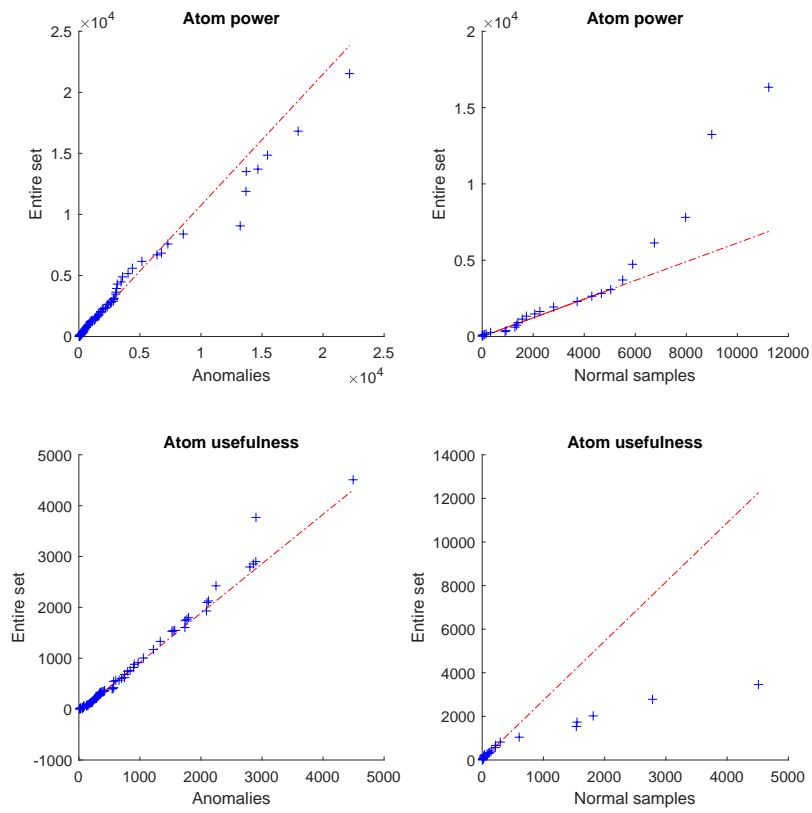


Figure 2: Quantile-Quantile plot of atom power and usefulness distributions in the normal and anomalies classes

differ significantly in anomalies from the normal samples. Figure 2 shows the quantile-quantile plot illustrating this comparison.

Our last contribution, called Separable Laplacian Classification [13] involves exploiting the structural information of graphs in order to identify anomalous patterns. We work directly with the structure of the graphs and adapt the separable dictionary learning problem, which takes into account vicinity patterns in 2D data. Our strategy is to exploit the two dimensional structure of a graph Laplacian in order to learn connectivity patterns that are specific to each class of graphs. The solution is tested on a synthetic dataset consisting of graph signals, in which the anomalies are structurally different from the normal samples. The investigated patterns are rings and cliques, since they are common in fraud-detection tasks. The training signals from each class are used separately to train one pair of dictionaries. Classifying a new, test signal, accounts to evaluating which pair of dictionaries is better at representing the signal. We also perform a test that investigates the effect of dictionary size on the ability of the Separable Laplacian Classification method in identifying the circular graph patterns.

References

- [1] B. Dumitrescu, P. Irofti, Dictionary Learning Algorithms and Applications, Springer International Publishing, 2018 (2018).
- [2] M. Tipping, Sparse Bayesian Learning and the Relevance Vector Machine, *Journal of Machine Learning Research* 1 (2001) 211–244 (2001).
- [3] D. P. Wipf, B. D. Rao, Sparse Bayesian learning for basis selection, *IEEE Transactions on Signal Processing* 52 (8) (2004) 2153–2164 (2004).
- [4] A. Băltoiu, B. Dumitrescu, Sparse bayesian learning algorithm for separable dictionaries, *Digital Signal Processing* 111 (2021) 102990 (2021).
- [5] O. Stegle, C. Lippert, J. Mooij, N. Lawrence, K. Borgwardt, Efficient inference in matrix-variate Gaussian models with iid observation noise, in: *Proceedings of the 24th Neural Information Processing Systems Conference*, 2011, pp. 630 – 638 (2011).
- [6] B. Dumitrescu, C. D. Giurcăneanu, Adaptive-Size Dictionary Learning Using Information Theoretic Criteria, unpublished document (2019).
- [7] A. Baltoiu, B. Dumitrescu, Size adaptation of separable dictionary learning with information-theoretic criteria, in: *2019 22nd International Conference on Control Systems and Computer Science (CSCS)*, 2019, pp. 7–11 (2019).
- [8] P. Irofti, A. Băltoiu, A. Pătrașcu, Fraud detection in networks, in: *Enabling AI applications in Data Science*, Springer, 2020, pp. 517–536 (2020).
- [9] Z. Jiang, Z. Lin, L. Davis, Learning A Discriminative Dictionary for Sparse Coding via Label Consistent K-SVD, in: *IEEE Conf. Computer Vision and Pattern Recognition*, 2011, pp. 1697–1704 (2011).

- [10] P. Irofti, A. Băltoiu, Malware identification with dictionary learning, in: 27th European Signal Processing Conference, 2019, pp. 1–5 (2019).
- [11] S. Matiz, K. Barner, Label consistent recursive least squares dictionary learning for image classification, in: 2016 IEEE International Conference on Image Processing (ICIP), IEEE, 2016, pp. 1888–1892 (2016).
- [12] P. Irofti, A. Băltoiu, Unsupervised dictionary learning for anomaly detection, in: International Traveling Workshop on Interactions Between Sparse Models and Technology, 2020, pp. 1–3 (2020). arXiv:2003.00293.
- [13] A. Băltoiu, A. Pătrașcu, P. Irofti, Graph anomaly detection using dictionary learning, in: The 21st World Congress of the International Federation of Automatic Control, 2020, pp. 3551–3558 (2020).