

MIT Open Access Articles

Kernel dependence analysis and graph structure morphing for novelty detection with high-dimensional small size data set

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Mohammadi-Ghazi, Reza et al. "Kernel dependence analysis and graph structure morphing for novelty detection with high-dimensional small size data set." *Mechanical Systems and Signal Processing* 143 (September 2020): 106775 © 2020 Elsevier Ltd

As Published: <http://dx.doi.org/10.1016/j.ymssp.2020.106775>

Publisher: Elsevier BV

Persistent URL: <https://hdl.handle.net/1721.1/126504>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of use: Creative Commons Attribution-NonCommercial-NoDerivs License



Kernel dependence analysis and graph structure morphing for novelty detection with high-dimensional small size data set

Reza Mohammadi-Ghazi^a, Roy E. Welsch^b, Oral Büyüköztürk^{a,*}

^a*Department of Civil and Environmental Engineering, Massachusetts Institute of Technology, Cambridge, MA, 02139 USA*

^b*Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA, 02139 USA*

Abstract

In this study, we propose a new approach for novelty detection that uses kernel dependence techniques for characterizing the statistical dependencies of random variables (RV) and use this characterization as a basis for making inference. Considering the statistical dependencies of the RVs in multivariate problems is an important challenge in novelty detection. Ignoring these dependencies, when they are strong, may result in inaccurate inference, usually in the form of high false positive rates. Previously studied methods, such as graphical models or conditional classifiers, mainly use density estimation techniques as their main learning element to characterize the dependencies of the relevant RVs. Therefore, they suffer from the curse of dimensionality which makes them unable to handle high-dimensional problems. The proposed method, however, avoids using density estimation methods, and rather, employs a kernel method, which is robust with respect to dimensionality, to encode the dependencies and hence, it can handle problems with arbitrarily high-dimensional data. Furthermore, the proposed method does not need any prior information about the dependence structure of the RVs; thus, it is applicable to general novelty detection problems with no simplifying assumption. To test the performance of the proposed method, we apply it to realistic application problems for analyzing sensor networks and compare the results to those obtained by peer methods.

Keywords: Novelty detection, kernel independence, reproducing kernel Hilbert space, two-sample test, Hilbert-Schmidt independence criterion, graphical model, structural health monitoring, sensor network, video camera, camera-based measurement.

1. Introduction

From the pattern recognition point of view, novelty detection can be viewed as a one-class classification that aims to distinguish one well sampled class from all other possible classes for which the available data is insufficient to build an explicit model for the latter [1]. Important methods in this regard are the k-nearest neighbors [2], one-class support vector machines [3], neural networks [4], density estimation and clustering [5], and decision tree based techniques such as one-class random forests [6]. The data from the observed class are usually represented in terms of certain features which can be modeled as random variables (RV); and the

*Corresponding author

Email address: rezamg@mit.edu, rwelsch@mit.edu, obuyuk@mit.edu (Oral Büyüköztürk)

8 above-mentioned methods are usually most effective when the statistical dependencies of the relevant RVs
9 are weak as these techniques ignore such dependencies. As a result these methods may provide inaccurate
10 predictions, especially in the form of high false positive rates, in applications such as saliency detection
11 in image processing [7], analyzing sensor networks [8], structural health monitoring and damage detection
12 [8, 9, 10], where the dependencies of RVs are not negligible.

13 To consider the statistical dependencies of RVs in novelty detection, previous studies applied methods such
14 as statistical graphical models [7, 8, 11] and conditional classifiers [10]. These techniques mainly use density
15 estimation for characterizing the dependencies between the RVs of the problem. Noting that the density
16 estimation techniques usually suffer from the curse of dimensionality [12], the novelty detection techniques
17 which use density estimation may not be able to handle high-dimensional problems. Moreover, some of these
18 techniques have limited applications due to simplifying assumptions they make or specific prior information
19 they need about the dependence structure of RVs [8]. These issues motivate the objective of this study which
20 is to develop suitable novelty detection algorithm with the capability of considering statistical dependencies
21 of relevant RVs in high-dimensional problems.

22 To address the objective of this study we develop a kernel dependence novelty detection (KDND) algorithm
23 that uses kernel two-sample tests [13, 14, 15, 16] and kernel independence analysis [17, 18, 19] as the basis
24 for making inference. Our proposed KDND method aims to detect novel realizations of RVs with respect to
25 a baseline by tracking the changes in the pairwise dependence structure of RVs. This dependence structure
26 is learned by using a kernel dependence criterion [17, 18] that is robust with respect to dimensionality
27 of data. By doing so, the contributions of our study can be summarized as follows: (1) We formulate
28 a KDND classifier that is capable of considering the dependencies of RVs in arbitrarily high-dimensional
29 novelty detection problems without any prior information about these dependencies; (2) We experimentally
30 evaluate the proposed method in structural health monitoring (SHM) problems, one on a full scale steel
31 structure, and compare the results with other techniques.

32 The paper is organized as follows. First, a review of kernel two-sample tests and kernel dependence
33 analysis is presented in Section 3. Then, in Section 4 we describe the problem and formulate our proposed
34 KDND method followed by an implementation of this method in Section 5. Section 6 provides a discussion
35 about how to relax some of the assumption we made for formulating the proposed method as well as more
36 details about the proposed formulation of the classifier. The results of the experimental evaluation of the
37 proposed algorithm and its comparison with other methods are presented in Section 7. Finally, we conclude
38 with a summary of our findings and a discussion of future research directions.

39 **2. Notation**

40 Through out this paper, the RVs are denoted by sans-serif fonts, e.g., x , and deterministic quantities such
41 as the realizations of RVs are denoted using serified fonts, e.g., x .

42 To exclude particular entries from a set we use " \setminus " followed by the indices of the entries to be excluded.
 43 For instance, if $\mathbf{Z} = \{\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3\}$, then $\mathbf{Z} \setminus \{\mathbf{z}_2\} = \{\mathbf{z}_1, \mathbf{z}_3\}$. $|\cdot|$ is also the cardinality of a set.

44 Probability distributions are denoted by p with a subscript denoting the RV; e.g., p_y is the probability
 45 density of the RV y . The probability of an event e is denoted by $\mathbb{P}(e)$.

46 Finally, through out this paper, we use the *standard* as the opposite of *novel*, e.g., a standard data point
 47 with respect to a baseline distribution is the one that is drawn from that distribution. For the Gaussian
 48 distribution, we always use the term *Gaussian*.

49 3. Review of kernel two-sample test and kernel independence analysis

50 In this study, we use the kernel two-sample test and kernel independence analysis as the building blocks
 51 for formulating our proposed KDND algorithm. Therefore, we first provide a review of these methods to be
 52 used in the following sections.

53 A two-sample test is a statistical test to infer whether two data sets are drawn from a same probability
 54 distribution [15]. One approach to perform such a test is to compare the estimated probability density models
 55 of the two data sets. The main issue of this approach is that density estimation techniques suffer from the
 56 curse of dimensionality and hence, density-based methods are not usually robust when the size of data set is
 57 relatively small compared to its dimensionality.

An alternative approach for two-sample test is kernel-based method that maps data sets into the re-
 producing kernel Hilbert space (RKHS) and uses appropriate similarity measures in this space to compare
 the two samples [13, 14]. To further clarify this procedure, consider two RVs \mathbf{y}_1 and \mathbf{y}_2 with the fixed but
 unknown probability distributions $p_{\mathbf{y}_1}$ and $p_{\mathbf{y}_2}$. Letting $k(\cdot, \cdot)$ to be a universal kernel associated with the
 RKHS, an appropriate distance between these distributions can be defined as [13]

$$D(p_{\mathbf{y}_1}, p_{\mathbf{y}_2}) = \mathbf{E}_{\mathbf{y}_1, \mathbf{y}'_1}[k(\mathbf{y}_1, \mathbf{y}'_1)] - 2\mathbf{E}_{\mathbf{y}_1, \mathbf{y}_2}[k(\mathbf{y}_1, \mathbf{y}_2)] + \mathbf{E}_{\mathbf{y}_2, \mathbf{y}'_2}[k(\mathbf{y}_2, \mathbf{y}'_2)], \quad (1)$$

where \mathbf{y}'_1 and \mathbf{y}'_2 are independent copies of \mathbf{y}_1 , and \mathbf{y}_2 , respectively, and \mathbf{E} is the expectation operator.
 $D(\cdot, \cdot)$, which is called the maximum mean discrepancy (MMD), basically measures the similarity of two
 distributions by comparing the expectation operators that are defined with respect to their distribution (See
 [13, 19] for more information about this similarity measure). In practice, however, we do not have access
 to RVs' distributions and we need to estimate the MMD using empirical data with finite samples. For such
 estimation, let $\mathbf{Y}_1 = \{\mathbf{y}_{11}, \dots, \mathbf{y}_{1m}\}$ and $\mathbf{Y}_2 = \{\mathbf{y}_{21}, \dots, \mathbf{y}_{2m}\}$ be two sets of m i.i.d. realizations drawn from
 $p_{\mathbf{y}_1}$ and $p_{\mathbf{y}_2}$, respectively. It can be shown that an unbiased estimate of $D(\cdot, \cdot)$ is [14]

$$D(p_{\mathbf{y}_1}, p_{\mathbf{y}_2}) \approx \widehat{D}(\mathbf{Y}_1, \mathbf{Y}_2) = \frac{1}{m(m-1)} \sum_{j \neq j'}^m k(\mathbf{y}_{1j}, \mathbf{y}_{1j'}) - 2k(\mathbf{y}_{1j}, \mathbf{y}_{2j'}) + k(\mathbf{y}_{2j}, \mathbf{y}_{2j'}). \quad (2)$$

58 where $j, j' \in \{1, \dots, m\}$.

The MMD, as it was mentioned before, is a measure of similarity between two arbitrary distributions; therefore, it can be extended to obtain a measure of dependency of RVs. To do that, first remember that two RVs \mathbf{y}_1 and \mathbf{y}_2 are statistically independent if their joint distribution equals the product of their marginal distributions, i.e., $p_{\mathbf{y}_1\mathbf{y}_2} = p_{\mathbf{y}_1}p_{\mathbf{y}_2}$. Therefore, comparing $p_{\mathbf{y}_1\mathbf{y}_2}$ and $p_{\mathbf{y}_1}p_{\mathbf{y}_2}$ can provide information about the dependency of \mathbf{y}_1 and \mathbf{y}_2 . Performing this comparison via the MMD, a measure of independence between the two RVs is obtained. This independence measure is called Hilbert-Schmidt independence criterion (HSIC) and estimated as follows [13, 18]

$$\eta_{\mathbf{y}_1\mathbf{y}_2} = \frac{1}{m^2} \text{tr}(\mathbf{H} \cdot \mathbf{K}_{\mathbf{y}_1\mathbf{y}_1} \cdot \mathbf{H} \cdot \mathbf{K}_{\mathbf{y}_2\mathbf{y}_2}), \quad (3)$$

where $\text{tr}(\cdot)$ is the matrix trace; $\mathbf{K}_{\mathbf{y}_1\mathbf{y}_1} = [k(\mathbf{y}_{1j}, \mathbf{y}_{1j'})]$ and $\mathbf{K}_{\mathbf{y}_2\mathbf{y}_2} = [k(\mathbf{y}_{2j}, \mathbf{y}_{2j'})]$; and $\mathbf{H} = \mathbb{I}_m - 1/m$ with \mathbb{I}_m to be the identity matrix of size m . The HSIC defined in (3) can be normalized as [18]

$$h_{\mathbf{y}_1\mathbf{y}_2} = \frac{\eta_{\mathbf{y}_1\mathbf{y}_2}}{\sqrt{\eta_{\mathbf{y}_1\mathbf{y}_1}\eta_{\mathbf{y}_2\mathbf{y}_2}}}. \quad (4)$$

This normalization maps $\eta_{\mathbf{y}_1\mathbf{y}_2}$ into the interval of $[0, 1]$, which makes it easier to interpret this independence criterion.

The MMD and normalized HSIC are used in the next section as a basis for formulating our KDND classifier.

4. Kernel dependence novelty detection classifier

Assume we have a system with n components that each can have two possible states, *standard* (intact/healthy) and *novel* (damaged), with respect to a baseline data set. Let x_i , $i \in \{1, \dots, n\}$, be a Bernoulli RV which can take on values in $\{-1, +1\}$ such that $x_i = +1$ if the i^{th} component of the system is intact and $x_i = -1$ if the i^{th} is novel. Assume we can observe all components of the system and represent the observations from the i^{th} component via a feature vector $\mathbf{y}_i \in \mathbb{R}^d$. Note that this feature vector is modeled as a d -dimensional RV, i.e., $\mathbf{y}_i = [y_{i1}, \dots, y_{id}]$, where y_{ik} , $k \in \{1, \dots, d\}$, is a univariate RV that represents a single feature that is extracted from the observations of the i^{th} component.

Consider the two sets $\mathbf{x} = \{x_1, \dots, x_n\}$ and $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$, and assume that the relationship between \mathbf{x} and \mathbf{Y} can be captured by a fixed but unknown joint probability distribution denoted by $p_{\mathbf{Y}, \mathbf{x}}$. Let us denote the baseline training set by \mathcal{T} which contains m data points, i.e., $\mathcal{T} = \{(\mathbf{Y}_j^{(b)}, \mathbf{x}_j^{(b)}); j \in \{1, \dots, m\}\}$ where $\mathbf{Y}_j^{(b)} = \{\mathbf{y}_{1j}^{(b)}, \dots, \mathbf{y}_{nj}^{(b)}\}$ is the set of feature vectors associated with the j^{th} baseline observation of the system, $\mathbf{x}_j^{(b)} = \{x_{1j}^{(b)}, \dots, x_{nj}^{(b)}\}$ are the class labels for these observations, and the superscript (b) denotes *baseline*. Note that in novelty detection problems, the baseline data belongs to only one state of the system [1, 20] in which all components are intact and hence, $x_{ij}^{(b)} = +1, \forall(i, j)$. Given \mathcal{T} and a set of m' new realizations of \mathbf{Y} , denoted by $\mathcal{Y} = \{\mathbf{Y}_1, \dots, \mathbf{Y}_{m'}\}$, the objective is to distinguish the standard and novel components of the system based on \mathcal{Y} , i.e., to predict $x_i, \forall i$.

80 Our proposed approach to satisfy the above-mentioned objective is to track the changes in the dependence
 81 structure of the RVs of the problem using kernel dependence methods. In doing so, first we consider a pairwise
 82 graphical model $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where \mathcal{V} is the set of vertices such that vertex $i \in \mathcal{V}$ is associated with \mathbf{y}_i , and
 83 the edges $\mathcal{E} = \mathcal{V} \times \mathcal{V}$ represent the pairwise dependencies between each pair of RVs. For the moment, let us
 84 assume that the edge potentials of this graph are fixed and unique for each possible realization of \mathbf{x} . Thus,
 85 we expect that the state transition of a subset of the system’s components from intact to novel is reflected in
 86 the properties of \mathcal{G} . Therefore, finding novel components of the system can be accomplished by comparing
 87 the dependencies encoded by two such graphical models which are respectively learned from \mathcal{T} and \mathcal{Y} . Note
 88 that the assumption about the uniqueness of the edge potential of \mathcal{G} for each realization of \mathbf{x} may not always
 89 hold. In Section 6 we discuss how to deal with the cases where this assumption is not true.

To proceed with formulating our KDND classifier, let N_i be the set of neighboring vertices of the i^{th} vertex
 of \mathcal{G} . For any edge $(i, i') \in \mathcal{E}, i' \in N_i$, we use $h_{\mathbf{y}_i \mathbf{y}_{i'}}$ as an estimate of the dependence strength associated
 with this edge. We compute $h_{\mathbf{y}_i \mathbf{y}_{i'}}$ using (4) and denote it by $h_{ii'}$ for ease of notation. Given the above
 assumptions and definitions, we define classifier $C(h_{ii'})$ as

$$x_i = C(h_{ii'}) = \text{sign}(f(h_{ii'})), \quad (5a)$$

$$f(h_{ii'}) = \sum_{i' \in N_i} w_{ii'} \ell_{ii'}, \quad (5b)$$

$$\ell_{ii'} = \text{sign}\left(p_{h_{ii'}}^{(b)}(h_{ii'}) - p_{ii'}^*\right), \quad (5c)$$

90 where $w_{ii'}$ is the weight that is (pre-)assigned to the edge $(i, i') \in \mathcal{E}$ and $\sum_{i' \in N_i} w_{ii'} = 1$. $w_{ii'}$ can be viewed
 91 as an importance factor of the (i, i') edge in making inferences about the i^{th} component of the system. The
 92 reason for considering these weights is to make our formulation applicable to problems where prior information
 93 is available on the dependencies of a subset of RVs. In case there is no such information, the weights can
 94 be set to $w_{ii'} = 1/|N_i|$ for $(i, i') \in \mathcal{E}$ and $\forall i' \in N_i$. $p_{h_{ii'}}^{(b)}(\cdot)$ is the baseline distribution of $h_{ii'}$, which is the
 95 dependence strength associated with the (i, i') edge; and $p_{ii'}^*$ is an appropriate likelihood threshold. Note
 96 that in this formulation, the edge potentials are modeled as RVs and shown using sans-serif font, i.e., $h_{ii'}$.
 97 The realization of this RV, that is obtained by using \mathcal{Y} along with (4), is shown by serif font, $h_{ii'}$.

98 To find the novel components of the system, the proposed classifier starts with classifying each edge of
 99 \mathcal{G} via (5c). The aim of this preliminary classification is to find the edges that encode significantly different
 100 dependence strengths, in a statistical sense, for \mathcal{T} and \mathcal{Y} . In doing so, the new edge potentials of the graph,
 101 $h_{ii'}$, are calculated using \mathcal{Y} along with (4). Then the likelihood of these potentials with respect to their
 102 baseline distribution, $p_{h_{ii'}}^{(b)}$, is compared with a likelihood threshold $p_{ii'}^*$ as stated in (5c). The result of this
 103 preliminary classification is $\ell_{ii'}$ which is a Bernoulli RV that can take on values in $\{-1, +1\}$. $\ell_{ii'} = +1$ means
 104 that the dependencies between \mathbf{y}_i and $\mathbf{y}_{i'}$ for the two data sets \mathcal{T} and \mathcal{Y} are not significantly different. The
 105 converse is true for $\ell_{ii'} = -1$ which indicates that the edge potential that is obtained from \mathcal{Y} for the (i, i')
 106 edge is statistically different from its baseline. Note that $h_{ii'}$ in (5) is obtained by using \mathcal{Y} along with (4),
 107 and $p_{h_{ii'}}^{(b)}$ is learned from the training set \mathcal{T} . We will address the details of learning this classifier in the next

108 section. Note that performing a likelihood test in (5c) is not the only way of classifying the edges. One can
 109 use other methods, such as off-the-shelf classifiers such as the SVM to do this task.

110 After classifying all edges of \mathcal{G} , we find the RVs which are most responsible for the discrepancies between
 111 the two graphs that are learned from \mathcal{T} and \mathcal{Y} . The idea to identify such RVs is to find the vertices whose
 112 dependencies with their neighbors have changed the most. This task is carried out in (5b) that counts the
 113 number of those incoming edges to a given vertex i which are significantly different from their baseline, i.e.,
 114 $\ell_{ii'} = -1, i' \in N_i$. This can be viewed as a voting approach where each edge can vote for the two vertices it
 115 connects. The advantage of using the voting strategy over other possible prediction methods will be discussed
 116 in Section 6. Note that the importance of the incoming edges to a given vertex may not be equal; thus, we
 117 consider $w_{ii'}$ to account for such importance and the generality of the formulation. The final step, which is
 118 carried out in (5a), is to classify the components of the system based on the majority of the votes that they
 119 have received.

120 Based on the above explanations, three tasks are needed for implementing the proposed KDND classifier.
 121 These tasks are: (1) choosing an appropriate kernel and learning its parameters, (2) learning $p_{h_{ii'}}^{(b)}$ from the
 122 training set, and (3) determining $p_{ii'}^*$. In what follows, we explain our solutions for each of these tasks.

123 5. Implementation of the proposed KDND classifier

124 5.1. Choice of kernel and determining its parameters

Using a universal kernel in the sense of [21] is a necessary condition for derivation of the MMD and HSIC
 [13, 14]. Based on the discussion provided in [22], we chose a Gaussian kernel to be used in this study. To
 use this kernel, consider $\mathbf{y}_{ij} = [y_{ij1}, \dots, y_{ijd}]$ as the j^{th} realization of \mathbf{y}_i , and $\mathbf{y}_{i'j'} = [y_{i'j'1}, \dots, y_{i'j'd}]$ as j'^{th}
 realization of $\mathbf{y}_{i'}$. Note that y_{ijk} is the k^{th} feature that is extracted from the j^{th} observation of the response
 of the i^{th} system's component. The dissimilarity between the two realizations \mathbf{y}_{ij} and $\mathbf{y}_{i'j'}$ with respect to
 an isotropic Gaussian kernel is

$$k(\mathbf{y}_{ij}, \mathbf{y}_{i'j'}) = \exp\left(-\frac{\|\mathbf{y}_{ij} - \mathbf{y}_{i'j'}\|^2}{2\sigma^2}\right), \quad (6)$$

125 where σ is the kernel width and $\|\cdot\|$ is vector norm. In practice, the median distance between the aggregate
 126 data points can be used as an estimate of σ [17, 18, 19].

Previous studies on kernel two-sample tests usually used the above-mentioned form of the Gaussian kernel.
 However, the isotropic property may not effectively capture the dissimilarities of data points in multivariate
 problem where the distributions of the feature y_{ik} , that were defined in Section 4, are different. To address
 this problem, we suggest using an anisotropic Gaussian kernel with the form of

$$k(\mathbf{y}_{ij}, \mathbf{y}_{i'j'}) = \exp\left(-\frac{1}{2}(\mathbf{y}_{ij} - \mathbf{y}_{i'j'})^T \mathbf{\Sigma}^2 (\mathbf{y}_{ij} - \mathbf{y}_{i'j'})\right), \quad (7)$$

127 where $\mathbf{\Sigma} = \text{diag}(1/\sigma_1, \dots, 1/\sigma_d)$ with $\sigma_k, k \in \{1, \dots, d\}$, is the kernel width along the k^{th} dimension which
 128 corresponds to the k^{th} feature. The superscript T is matrix/vector transpose. Note that the kernel matrix

129 that is obtained from (7) is strictly positive definite because Σ^2 is diagonal with positive diagonal elements;
 130 therefore, this anisotropic kernel is universal [23, 24].

131 The anisotropic kernel in (7) considers a unique kernel width for each dimension of the multivariate RVs
 132 \mathbf{y}_i and $\mathbf{y}_{i'}$ and hence, it is expected to better capture the shape of data. However, the trade-off of using such
 133 kernel is to learn d parameters instead of one. More information about anisotropic Gaussian kernels and
 134 learning its parameters for various applications can be found in [25, 26, 27]. In this study, we propose using
 135 the median distance of aggregate data points along the k^{th} direction as an estimate of σ_k . This proposition
 136 basically uses the conventional method suggested in [17, 18, 19] for each feature individually to estimate its
 137 corresponding kernel width.

138 5.2. Learning $p_{\mathbf{h}_{ii'}}^{(b)}$

139 Our approach here is to generate multiple realizations of $\mathbf{h}_{ii'}$ and learn its distribution accordingly. For
 140 a given $(i, i') \in \mathcal{E}$, consider $(\mathbf{y}_{ij}^{(b)}, \mathbf{y}_{i'j'}^{(b)}, \forall (j, j'))$, which are all permutations of training sample points for the
 141 two vertices i and i' . The dissimilarity between each pair of data points corresponding to a unique (j, j') with
 142 respect to an anisotropic Gaussian kernel can be obtained using (7). The calculated dissimilarities for all
 143 pairs of (j, j') can be assembled in an $m \times m$ kernel matrix $\mathbf{K}_{ii'}$. The procedure can be followed to assemble
 144 \mathbf{K}_{ii} and $\mathbf{K}_{i'i'}$ for (i, i) and (i', i') , respectively. Using these matrices along with (3) and (4) provides a single
 145 realization of $\mathbf{h}_{ii'}$, whereas we need multiple such realizations for learning $p_{\mathbf{h}_{ii'}}^{(b)}$.

146 In order to generate multiple samples for $\mathbf{h}_{ii'}$, we use bootstrap aggregation (bagging). For doing so, we
 147 randomly pick \tilde{m} number of training data points. The dissimilarities between the chosen data points can be
 148 determined using (7), and assembled to form new kernel matrices which are called sampled kernel matrices
 149 and denoted by $\tilde{\mathbf{K}}_{ii'}$, $\tilde{\mathbf{K}}_{ii}$, and $\tilde{\mathbf{K}}_{i'i'}$ for $\mathbf{K}_{ii'}$, \mathbf{K}_{ii} , and $\mathbf{K}_{i'i'}$, respectively. Note that, for computational
 150 efficiency, the sampled kernel matrices can be formed by finding and assembling the corresponding elements
 151 of $\mathbf{K}_{ii'}$, \mathbf{K}_{ii} , and $\mathbf{K}_{i'i'}$ to the randomly chosen data points. Using the sampled kernel matrices along with (3)
 152 and (4), a new realization of $\mathbf{h}_{ii'}$ is obtained. By running this procedure multiple times, a set of realizations
 153 for $\mathbf{h}_{ii'}$ can be generated. Finally, we use Gaussian mixture models (GMM) to learn $p_{\mathbf{h}_{ii'}}^{(b)}$ from the generated
 154 realizations of $\mathbf{h}_{ii'}$. Note that the curse of dimensionality is not a problem in using GMM for learning $p_{\mathbf{h}_{ii'}}^{(b)}$,
 155 because it is a univariate distribution and we can generate as many realizations of $\mathbf{h}_{ii'}$ as needed through the
 156 bagging procedure.

157 5.3. Determining $p_{ii'}^*$

158 If there exists a prior distribution for the novel realizations of \mathbf{Y} , the parameter $p_{ii'}^*$ can be determined
 159 accordingly. Otherwise, in the absence of such prior information, hypothesis testing can be used. In this
 160 case, $p_{ii'}^*$ is the likelihood threshold associated with a $100(1 - \alpha)\%$ confidence bound for $p_{\mathbf{h}_{ii'}}^{(b)}$, where α is a
 161 predefined significance level. Note that the notion of confidence interval for non-symmetric or multi-modal
 162 distributions is controversial [28]. Thus, we propose determining $p_{ii'}^*$ by using the notion of a high density
 163 region (HDR), which is described as follows.

For a probability distribution $p_z(z)$ of RV z , the $100(1 - \alpha)\%$ HDR is defined as the subset $R(p_\alpha^*)$ of the sample space of z such that [29]

$$R(p_\alpha) = \{z : p_z(z) \geq p_\alpha^*\}, \quad (8)$$

164 where p_α^* is the largest value such that $\mathbb{P}(z \in R(p_\alpha^*)) \geq 1 - \alpha$. This is schematically shown in Figure 1 for
 165 a univariate distribution. By applying this concept to our problem, $p_{ii'}^*$ becomes the $100(1 - \alpha)\%$ HDR of
 $p_{h_{ii'}}^{(b)}$.

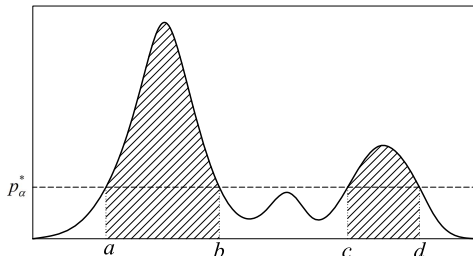


Figure 1: The $100(1 - \alpha)\%$ HDR for z is $R(p_\alpha^*) = [a, b] \cup [c, d]$ if $\mathbb{P}(z \in R(p_\alpha^*)) \geq 1 - \alpha$ [29]

166

167 6. Discussions

168 In this section, we clarify how to solve the problem when the main assumption we made in Section 4 do
 169 not hold. We also discuss alternative approaches for edge classification instead of density estimation, the
 170 reasons for using the averaging in (5b), and how to use this classifier when the data sets are unbalanced.

171 6.1. Using the KDND method when the main assumptions do not hold

172 To formulate the proposed KDND method we assumed that the changes in the state of some components
 173 of the system, from standard/intact to novel, affect the edge potentials of \mathcal{G} . This assumption may not
 174 always be true as it is theoretically possible that there exist special state transformations which change the
 175 distributions of the relevant RVs in such a way that their dependencies are retained unaffected. Although
 176 such special cases are unlikely in practice due to noise and complexities of the systems being monitored, we
 177 address a solution for this issue to ensure the robustness of our method.

178 For the situation that was described above, using the HSIC is no longer effective in detecting the novel
 179 realizations of the RVs. But, we can still use MMD as suggested in section 7.5 of [14] to detect any changes
 180 in the marginal distributions of the RVs. MMD, as described in this reference, is capable of detecting the
 181 novel realizations of RVs without suffering from the curse of dimensionality; however, this method may
 182 result in high false positive rates due to ignoring the dependencies of RVs. Thus, we suggest performing a
 183 preliminary classification using the MMD to ensure capturing the state changes in the system. Once some
 184 state transformations are detected in this preliminary analysis, we can use the proposed KDND algorithm to
 185 improve the accuracy of our predictions.

186 6.2. Alternative methods for learning $p_{h_{ii'}}^{(b)}$

187 As was explained in section 4, we use likelihood tests in (5c) to classify each edge of \mathcal{G} . This requires
188 learning $p_{h_{ii'}}^{(b)}$ for which we proposed using GMM; however, this is not the only approach for classifying the
189 edges. In fact, any novelty detection classifier can be trained using the realizations of $h_{ii'}$, that are obtained
190 via bagging (see section 5.2), and used instead of the likelihood test in (5c).

191 6.3. Reason for using the voting strategy in (5b)

192 In the formulation of our proposed KDND classifier, $\ell_{ii'}$ is a binary RV which is the vote of edge $(i, i') \in \mathcal{E}$
193 for the i^{th} and i'^{th} vertices. An alternative for this voting strategy would be the direct use of the likelihood
194 difference, $p_{h_{ii'}}^{(b)}(h_{ii'}) - p_{ii'}^*$, in our predictions. Due to the inverse relationship of the HSIC's magnitude with
195 its variance [14], making inferences based on the explicit value of $p_{h_{ii'}}^{(b)}(h_{ii'})$ results in the dominance of those
196 edges of \mathcal{G} whose potentials are small. This is equivalent to assigning more weight to the weakly dependent
197 RVs in the process of decision making. By doing so on the extreme case, we only consider the RVs which are
198 almost independent, and this is not consistent with the objective of this study. For more information about
199 the variance analysis of the MMD and HSIC the readers are referred to [14]. For an example showing the
200 behavior of these two measures in comparing Gaussian RVs see Appendix Appendix A.

201 6.4. Dealing with unbalanced data sets

202 It is common in practice that the baseline and test data sets are not equally sized, i.e., $m \neq m'$. This
203 results in different error rates in computing the dependence strength of RVs for \mathcal{T} and \mathcal{Y} [14]. Herein, we
204 provide two solutions to get around this problem. The first solution is applicable if $m > m'$ and can be
205 applied to the bagging procedure for learning the underlying distribution of $h_{ii'}$. To do that, we suggest
206 simply choosing $\tilde{m} = m'$ for each bagging iteration. Otherwise, if $m' > m$ or the sizes of data sets are too
207 small, we can use over sampling methods such as the synthetic minority over sampling technique [30].

208 7. Experimental evaluation

209 To evaluate the efficacy of the proposed method, we applied it to realistic problems of analyzing sensor
210 network data in SHM applications. In vibration-based SHM, the dynamic behavior of structures is used as a
211 basis for evaluating their health and safety. For this purpose, their vibrational responses are measured and
212 the measurements are compared with a baseline that represents the intact state of the system. This baseline
213 can be a set of theoretical facts about the behavior of structures or an empirical data set that has been
214 provided from the same structure at the intact or reliable state. Noting that there may be an infinite number
215 of damaged states for a structure, the SHM problem clearly fits the novelty detection framework [31, 32].

216 One of the main goals of this comparison is to detect and localize damages in the structure [33]. In recent
217 years, statistical methods and machine learning techniques have assisted the researchers in achieving this goal
218 [31, 32, 34, 35]. SVM [36], clustering techniques [37], and deep learning [38] are examples of such methods.

219 However, the predictions from most of these techniques are inaccurate if the size of data set is small, the
220 relevant RVs are strongly dependent, or the dimensionality of the data is extremely large [12]. Our proposed
221 KDND classifier has been specifically designed for handling these issues and hence, can be considered as an
222 approach to deal with real-world SHM applications.

223 In the following two sections, we present the application of our proposed method in monitoring a plate
224 structure and a full scale steel structure. In each of these sections we first describe the experimental setup
225 and the data acquisition system we used to measure the response of the structures. Then, we show the
226 damage localization result of the proposed KDND algorithm, and its comparison with peer novelty detection
227 methods.

228 *7.1. Plate structure*

229 *7.1.1. Experimental Setup*

230 This structure was a steel plate with dimensions of 60 cm \times 5.08 cm \times 0.64 cm, fixed to a massive concrete
231 base using four bolts. The experimental setup for measuring the plate involved a shaker that is attached to
232 the top of the plate, high-speed camera, and extra lighting, as shown in Figure 2. The video camera and the
233 extra lighting are used along with the phased-based optical flow approach [39, 40] to extract the displacement
234 field of the plate from video.

235 A summary of how the displacement extraction algorithm works is as follows. The complex steerable
236 pyramid filters are used to obtain the local amplitude and phase at multiple orientations and physical length
237 scales [41]. As was shown in previous works [41], the motion of constant phase contours corresponds to motion
238 signals in the video. Therefore, we use the decomposed local phase signals to calculate the displacement signal
239 at every pixel in the video. More details about the video decomposition and the calculation of displacements
240 is contained in [42]. Note that there are few assumptions in calculating the displacement signal from objects
241 in video. The first assumption is that the motion must be small and on the order of one pixel or smaller;
242 otherwise, the local filters stop working. Secondly, displacements are only well defined at edges or textured
243 regions in the video. In order to satisfy this condition, we applied a speckled pattern on the steel plate by first
244 painting it white and then spraying a random pattern of black paint on the plate in a spotty manner. The
245 last assumption for the measurement is that the lighting is constant. Flickering lights, such as fluorescent
246 lighting, introduce an apparent motion signal into the video with the same variation in time as the lighting.
247 To prevent this from happening, the object under test is flood illuminated with several bright battery powered
248 lamps so that the lighting stays consistent.

249 To run the tests, the shaker excited the plate with a white Gaussian noise waveform in a horizontal
250 direction in the video. After allowing some time for the excitation to reach steady state, 3.5 seconds of video
251 were recorded at 2000 frames per second with a resolution of 1736 \times 244 pixels using the high speed camera.
252 This was repeated for the damaged plate with a machined crack towards its base. The details of the intact and
253 notched plates, and screenshots from the input videos are shown in Figure 3. To extract the displacements
254 of the plate, 100 pixels are chosen, on a grid of 5 \times 20, as the pseudo sensor locations on this structure.

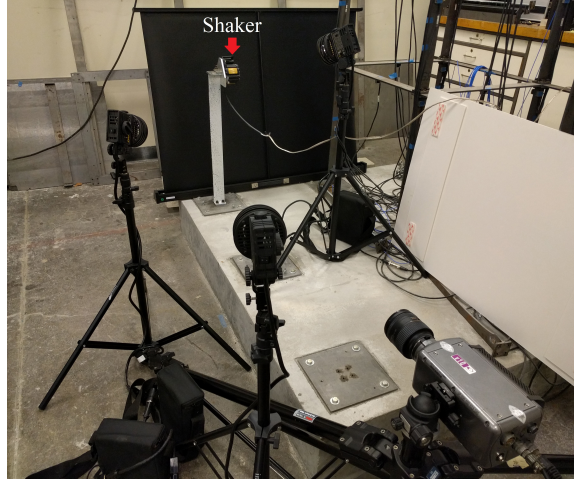


Figure 2: Picture of the experimental setup for the plate, showing the plate fixed to a concrete base, the shaker bolted to the top of the plate, high-speed camera, and extra lighting

255 The displacements of these locations were calculated from the recorded video using the phase-based method
 256 described above. Each pixel’s displacement time history was windowed in time to produce 22 pseudo tests
 257 that each consisted of 800 time points. The reason for dividing one measurement into multiple windows is
 258 that in the real world it is unlikely to have multiple extreme events in civil structures. However, we may
 259 have few extreme events that last for about 10 to 20 seconds or even more.

260 *7.1.2. Damage sensitive features and evaluation criteria*

261 The Fourier coefficients in the frequency range of 1 to 400 Hz are used as the damage sensitive features
 262 in this example. As a result, the dimensionality of the feature space becomes 400 while we have only 22
 263 pseudo tests. Many of the widely used novelty detection algorithms, such as the one-class SVM, cannot be
 264 learned for such a data set in which the number of features significantly exceeds the number of sample points.
 265 Therefore, to evaluate the efficacy of the proposed KDND method, its classification results were compared
 266 to a one-class gradient boosting algorithm [6, 12] which is robust with respect to dimensionality of data sets.

For comparing the results of the proposed method and peer methods we used the *false positive reduction* and *true positive improvement* criteria which were defined in [8]. The definition of false positive (FP) and true positive (TP) are controversial in SHM application as it is unlikely that the damage occurs exactly at the sensor locations. Therefore, the FP for SHM is defined in [8] as the detection of damage at sensor locations which may not coincide with the damage location nor the closest neighboring sensor locations. The TP is also defined as detecting the damage at the sensor location where the damage is located or its closest neighboring sensor locations. Assume $I = \{1, \dots, n\}$ with n to be the number of sensors in a network, D_a is the set of sensor locations which are either at the damage location or its closest neighborhood, and D_a^c is the complement of D_a with respect to I . Also assume that D_k and D_g are the set of damaged locations which are detected by the proposed KDND method and the alternative novelty detection method, respectively. Then,

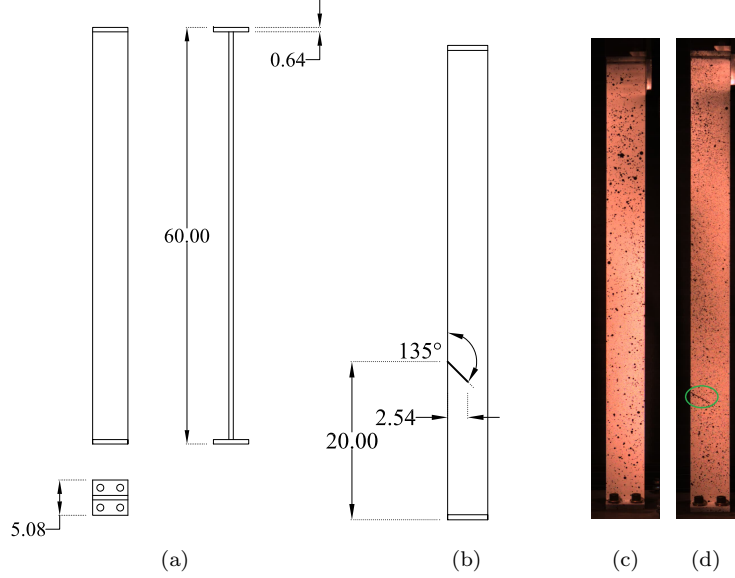


Figure 3: Details of the plate structures (dimensions are in centimeters): a) detail of the plates, b) location of the machined crack on the damaged plate, c) screenshot of input video for the intact plate, (d) screenshot of input video for the damaged plate with crack

for a new measurement the FP reduction (FPR) and TP improvement (TPI) criteria are defined as

$$FPR = \frac{|D_g \cap D_a^c| - |D_k \cap D_a^c|}{|D_a^c|} \quad (9a)$$

$$TPI = \mathbf{1}_{D_k \cap D_a \neq \emptyset} - \mathbf{1}_{D_g \cap D_a \neq \emptyset} \quad (9b)$$

267 where $|\cdot|$ is the cardinality of a set, and $\mathbf{1}_{(\cdot)}$ is the indicator function. *FPR* essentially shows how much
 268 the FP is reduced by using the proposed method compared to the gradient boosting algorithm. *TRI* shows
 269 which algorithm can/cannot localize the damage.

270 7.1.3. Fitting the KDND models and damage detection results

271 The graph \mathcal{G} in this example is considered as a fully connected pairwise graphical model in which, a
 272 generic edge (i, i') encodes the dependencies of the displacement responses at the i^{th} and i'^{th} pseudo sensor
 273 locations. Note $\mathbf{y}_i, i \in \{1, \dots, 100\}$, which is the feature vector associated with the response of the i^{th} pseudo
 274 sensor is a 22×400 matrix in this experiment. Using these vectors, we followed the procedure described in
 275 Section 5 to learn the KDND model. Due to the lack of prior information about the dependencies of the
 276 sensor measurements, the edge weights are considered to be $w_{ii'} = 1/|N_i|$ for $(i, i') \in \mathcal{E}, \forall i \in \mathcal{V}$ and $\forall i' \in N_i$.
 277 We also chose an anisotropic Gaussian kernel which results in computing 400 different kernel widths for each
 278 edge of the graph. Figure 4 shows the kernel width for four of the features (out of 400) and all edges. The
 279 kernel width between sensors i and j is colored at the intersection of the i^{th} row and the j^{th} column of these
 280 plots. Note that these graphs are symmetric because there is no difference between the edges (i, j) and (j, i)
 281 for computing the kernel width. Also, the diagonal entries are zero in these plots, because we do not consider

self-looping in \mathcal{G} and hence, we do not compute the associated kernel width. These plots imply that the kernel width can be quite different for different features and different subsets of the edges.

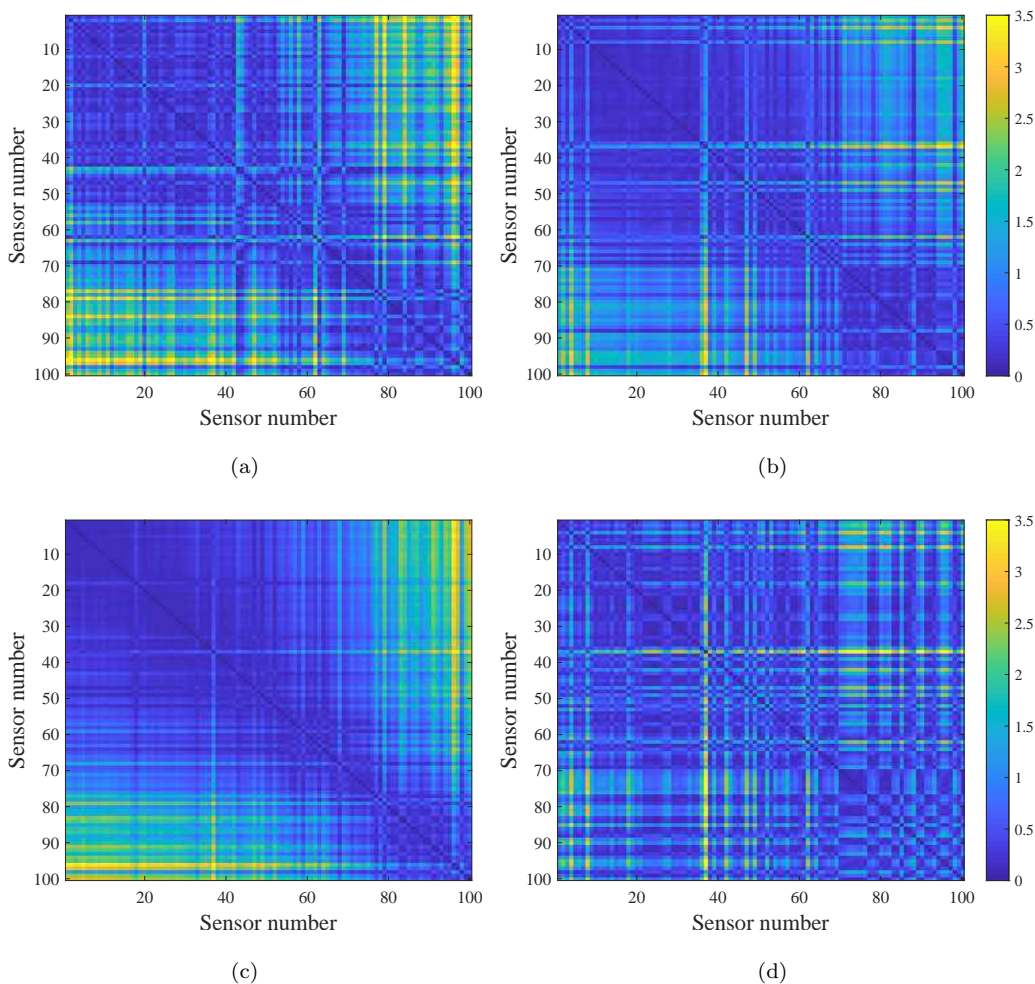


Figure 4: Kernel width for (a) first feature, (b) 10th feature, (c) 200th feature, (d) 100th feature.

For learning $p_{h_{i,i'}}^{(b)}$, we first set $\alpha = 0.05$ as the HDR significance level for decision making and then started fitting KDND models with various bagging trials starting from 10 and increased this number until the prediction error converges. The result of this procedure is shown in Figure 5(a). It follows from this plot that the FP rates does not change after the bagging trials reach 30.

In this case study we chose the one-class gradient boosting method as the alternative method to be compared to the KDND. The main reason for choosing the gradient boosting method in this part is the robustness of decision tree based algorithms with respect to the dimensionality of learning problems. Also, most of other widely used novelty detection algorithms, such as the one-class SVM, are not applicable to this application example due to the high dimensionality of the feature space. To learn the one-class gradient boosting classifier, we followed the approach that is explained in [6, 12].

294 The damage localization result of the proposed algorithm and its comparison with the predictions of
 295 the one-class gradient boosting algorithm are shown in Figures 5(b) and 5(c). In Figure 5(b), the average
 296 classification results of the gradient boosting algorithm for the 22 pseudo tests on the damaged structure
 297 are color coded. The red and green colors in these plots are, respectively, used to show whether a pseudo
 298 sensor location is predicted as damaged or intact. It follows from the results shown in Figure 5 that both
 299 algorithms can detect the damage, but their damage localization accuracies are significantly different. The
 300 localization results of the gradient boosting algorithm are almost inconclusive, while the proposed KDND
 301 algorithm can perfectly localize the damage by detecting two pseudo sensor locations right above the notch.
 302 The proposed algorithm also has some false positives, mainly at the top of the plate. These false alarms might
 303 have been due to the loss of the speckled pattern in that zone on the plate and hence inaccurate displacement
 304 extraction, especially since the lightening is weaker at the top of the plates compared to the other areas (see
 305 Figure 3(c)). This suggests that more attention needs to be paid when using video-based measurements in
 306 practice as this method is quite sensitive to lightening. Comparison of the two algorithms using the criteria
 307 defined in Section 7.1.2 shows that the proposed algorithm reduces the FP rate by 24% while keeping the
 308 TP rate the same. This means that by using the KDND method we localize the anomalies more accurately
 309 without sacrificing the TP.

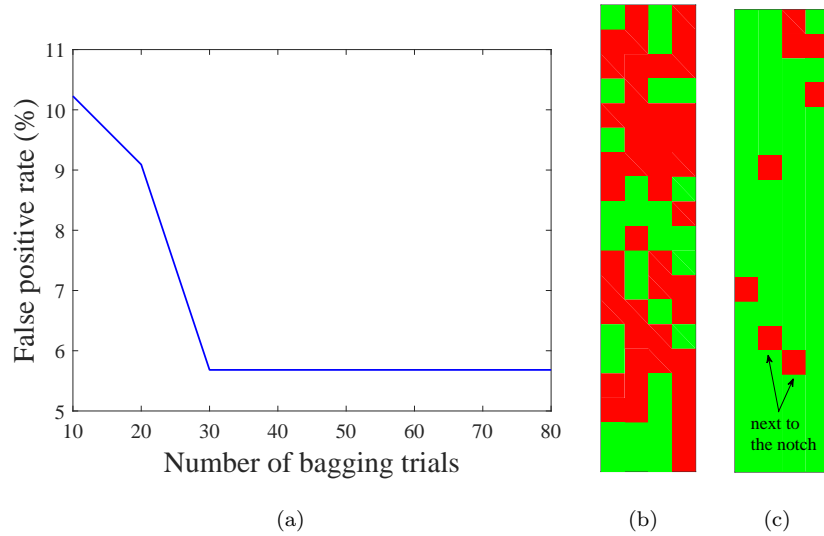


Figure 5: KDND training and damage localization results; (a) Variation of false positive rates versus the number of bagging trials, (b) damage detection results of gradient boosting, (c) damage detection results of the proposed KDND algorithm. In (a) and (b), the colors show the binary classification results where red and green colors at a pseudo sensor location, respectively, means that location is predicted as damaged or intact.

310 7.2. Full scale steel structure

311 7.2.1. Experimental Setup

312 The full scale structure that we study in this chapter is the upper half of a truncated telecommunication
313 tower. The dimensions of this tower are $83 \times 83 \times 239$ inches and it consists of L-shaped steel elements with
314 bolted connections as shown in Figure 6. The structure has three stories, which are all similar in shape
315 and configuration, and a cab floor which is shown in Figure 6(c). The beams are connected to the columns
316 using one bolt at each side. The bracing-column connections are also made by one bolt, but the connection
317 between a pair of bracings when they cross are made using three bolts as shown in Figures 6(d) and 6(e).
318 The four faces of the structure are identical; therefore, we have named them as A, B, C, and D to be used
319 when addressing damage scenarios. These faces and the global coordinate system are shown in Figure 6(c).

320 To measure the dynamic response of the structure it was instrumented by 48 triaxial MEMS accelerome-
321 ters, as well as two shakers at its top corner for generating excitation along the two perpendicular axis of the
322 structure (Figure 6). The maximum sampling rate of the MEMS sensors was 2 kHz, and the shakers could
323 generate white Gaussian noise excitation in the frequency range of 5 to 350 Hz. The sensor network and the
324 shakers were controlled via a central computer and a data acquisition system that were placed in a trailer
325 next to the structure. Note that the structure was symmetric, but torsional modes could be excited due to
326 the placement of the shakers.

327 7.2.2. Experimental tests and damage scenarios

328 We considered five different damage scenarios in addition to testing the intact structure. The first damage
329 scenario was introduced by loosening a bolt in a beam-column connection at the location shown in Figure
330 7(a). The second damage scenario was introduced by replacing a reduced cross section element with one of
331 the diagonal elements of the tower. Figures 7(b) and 7(d), respectively, show the location of this damage
332 scenario and the reduced cross section element. To build the reduced cross section element, both flanges of
333 a diagonal element were machined to reduce the element's cross section throughout its length. The third
334 damage scenario was made by taking out the reduced cross section element from the structure. Basically,
335 the location of this damage was the same as in the second scenario, but for the third case, the element was
336 taken out. The fourth and fifth damage scenarios were simulating the presence of multiple damages on a
337 structure. For these damage scenarios we removed the beams that are shown in Figure 7(c) from face C of
338 the structure, and combined this scenario with the second and the third scenarios. Table 1 summarizes the
339 tested damage scenarios on this structure.

340 To establish a baseline data set for the intact state, the structure was tested four times. In each test,
341 the structure was excited under a pink noise excitation with the spectrum that is shown in Figure 8 for two
342 minutes. Then, each test was segmented, by windowing with overlaps, to generate 14 pseudo tests. The
343 reason for this segmentation of the structural response was to model real world scenario, as it was explained
344 in Section 7.1.1. Each damage scenario was also tested two times with the same manner as testing the
345 intact structure. Therefore, we provided 56 and 28 pseudo tests, respectively, for the intact and each damage

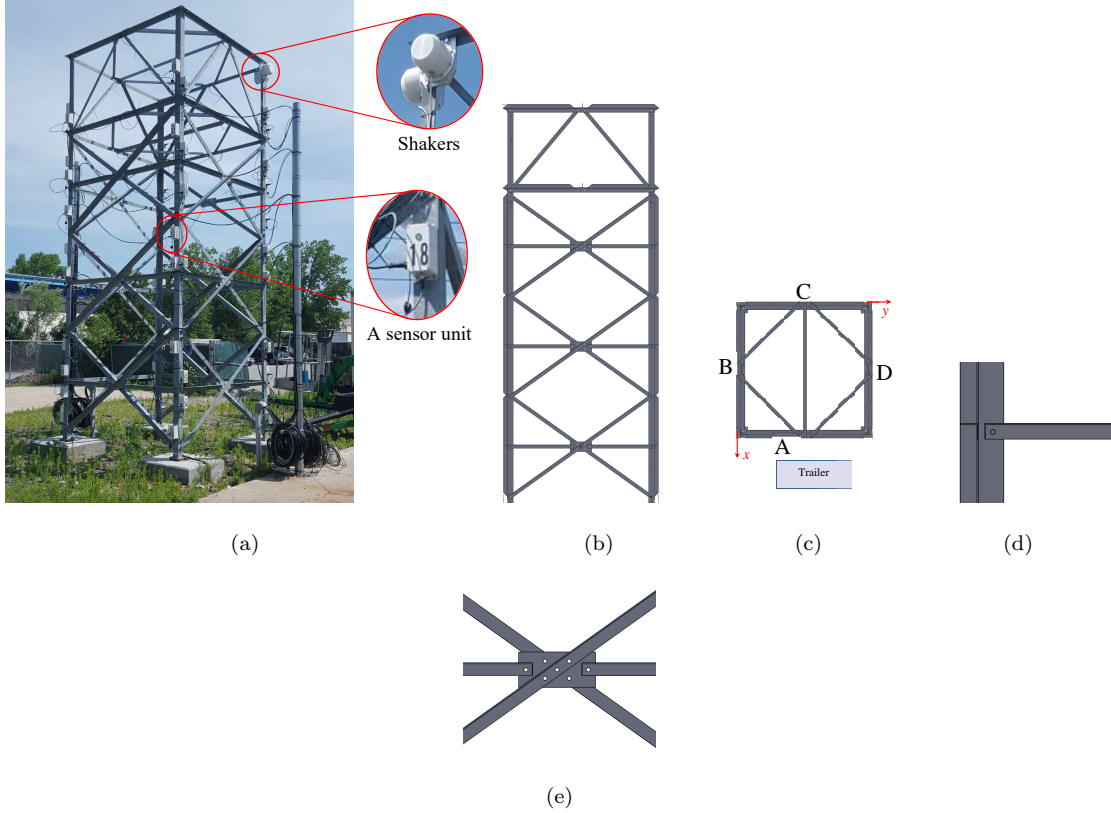


Figure 6: Large scale structure: (a) structure, sensor units, and shakers, (b) side view, (c) roof plan of the top floor, the location of the trailer with respect to the structure, and the name of each of its faces, (d) beam-column connection, (e) bracing connection

346 scenarios. For damage sensitive features, we considered the autoregressive (AR) coefficients. To determine
 347 the order of the AR model, we followed the AR model selection approach using the Akaike information
 348 criterion (AIC) that is suggested in [43]. Figure 9 shows the variation of the AIC for different AR model
 349 orders. Each line in this plot corresponds to the AR model selection result of a specific sensor. It follows
 350 from these results that the variations of the AIC are negligible for model orders that are larger than 16 for all
 351 sensors' measurements; thus, we choose this number as the AR model order in our study. The weights $w_{ii'}$
 352 were chosen similar to the previous experiment due to the lack of prior information about the dependencies
 353 of sensor measurements, i.e. $w_{ii'} = 1/|N_i|$ for $(i, i') \in \mathcal{E}$, $\forall i \in \mathcal{V}$ and $\forall i' \in N_i$. For comparing the performance
 354 of the KDND with other methods on this structure we use the evaluation criteria that were defined in section
 355 7.1.2.

356 7.2.3. Detection results and comparison

357 For learning the KDND, we followed the same procedure that was explained in section 7.1.3 for the
 358 plate structure. For evaluating the performance of the proposed method on this structure, we compared the
 359 KDND with the one-class gradient boosting, one-class SVM, and one-class clustering via GMM. We followed
 360 [44, 3] for learning the SVM models and [5] for learning the clustering model. The detection results of these

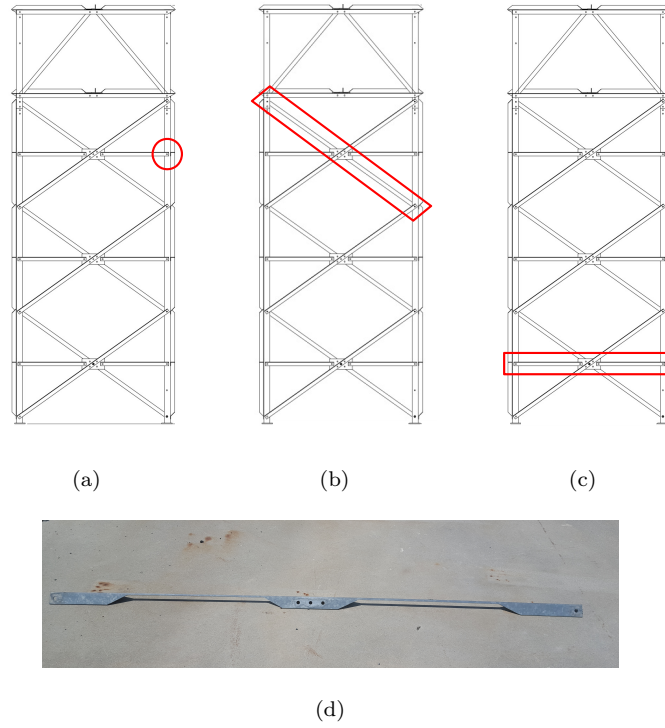


Figure 7: Location of damages on the structure. (a) location of the loosened bolt at the beam-column connection on face A, (b) location of the reduced cross section element on face A, (c) location of the removed beam on face C of the structure for damage scenarios 4 and 5, (d) the reduced cross section element

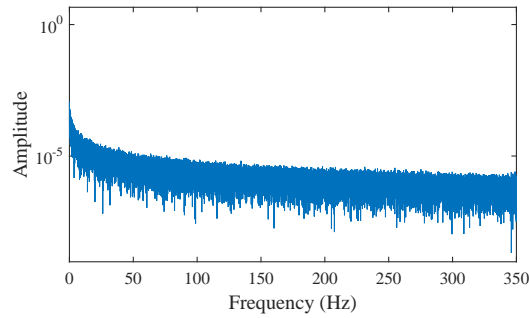


Figure 8: The spectrum of the pink noise excitation

361 algorithms for the first damage scenario are illustrated in Figure 10. The sensor locations are shown with
 362 circles in the plots of this figure. In Figure 10(a) the actual location of damage is marked by a red circle,
 363 and the black lines and circles, respectively, show the intact elements and sensor locations. Figures 10(b)
 364 to 10(e) show the detection results of the four above-mentioned methods. The red and blue circles in these
 365 plots are the sensor locations which are predicted as damaged and intact, respectively. It follows from these
 366 plots that all techniques can detect the damage; however, the localization results are significantly different.
 367 Basically, the localization with the gradient boosting is inconclusive. The results of SVM and clustering are

Table 1: Summary of the tested damage scenarios

Scenario No.	Description
1	Bolt loosening at a beam-column connection at the location shown in Figure 7(a) on face A
2	Reduced cross section element at the location shown in Figure 7(b) on face A
3	Element removal at the location shown in Figure 7(b) on face A
4	Multiple damage scenario by combining scenario #2 and element removal at the location shown in Figure 7(c) on face C
5	Multiple damage scenario by combining scenario #3 and element removal at the location shown in Figure 7(c) on face C

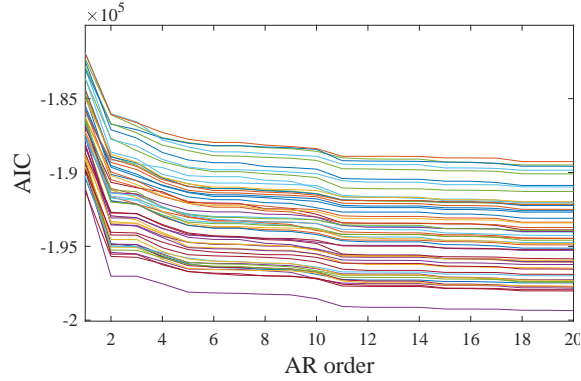


Figure 9: Variation of AIC for different AR model orders for the vibration response at all sensor locations. There are 56 lines in this plot and each line is associated with one sensor location.

368 more concentrated around the true location of damage, but these techniques also suffer from high false rates.
 369 The KDND, on the other hand, provides a more accurate localization result by detecting the true location
 370 of damage and one sensor location that is directly connected to the damaged location. This method has
 371 only one false detection that is two elements away from the actual location of damage. The predictions of
 372 these algorithms for other damage scenarios follow a similar pattern; thus, we skip showing the plots for
 373 brevity, and instead, summarize the results in Table 2. It should be noted that all algorithms can correctly
 374 detect the damage in all scenarios; therefore, we only report the false positive reduction rates in this table.
 375 The information in this table implies that considering the dependencies of RV via the KDND method can
 376 effectively reduce the false positive rates between 14% to 33% without affecting the true positive rates.

377 8. Conclusion

378 In this paper we have proposed a novelty detection method that uses kernel dependence analysis for
 379 considering the statistical dependencies of the problem's RVs to make predictions. The method considers a
 380 pairwise graphical model over the RVs and aims to detect statistically significant variations in the parameters

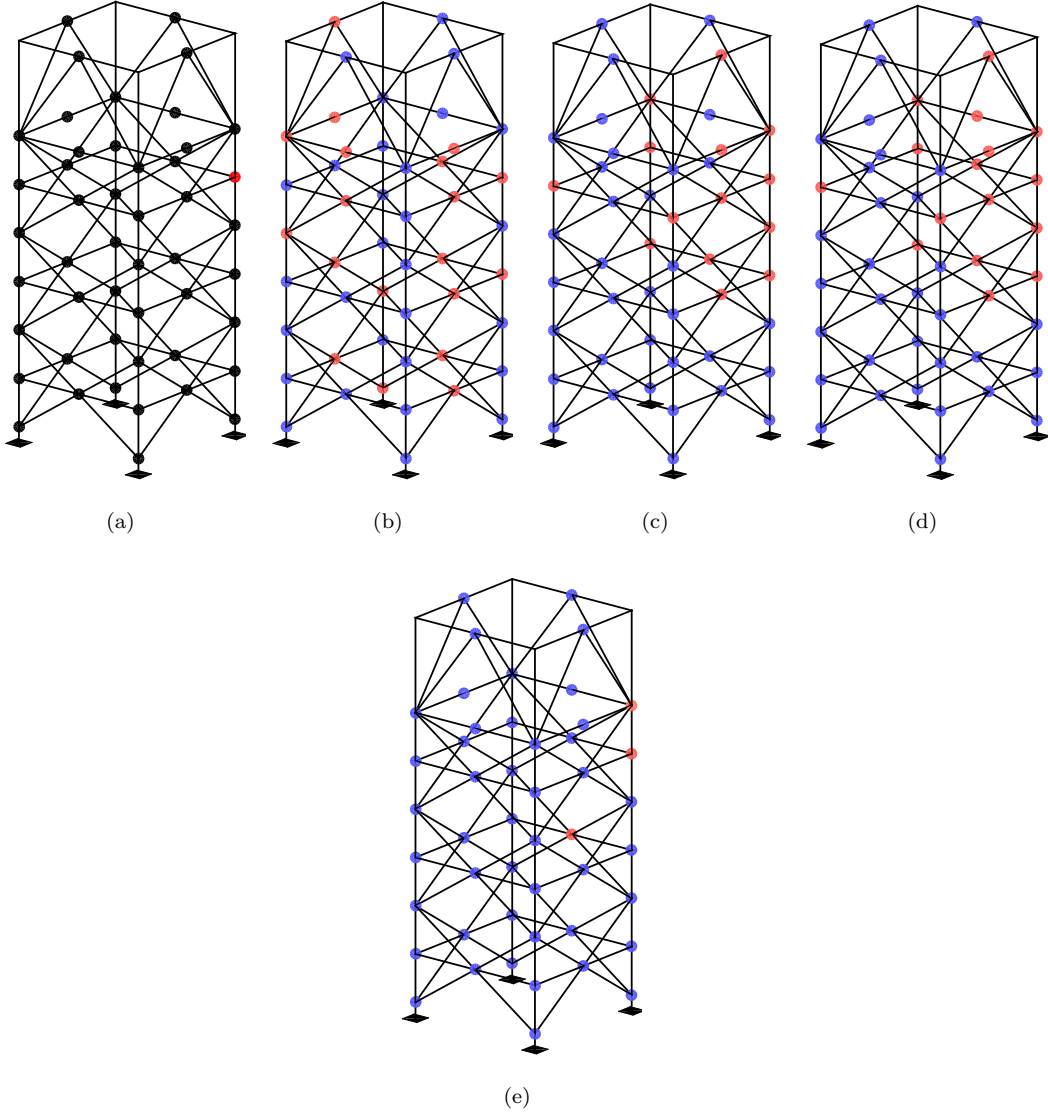


Figure 10: Damage detection/localization results for the bolt loosening damage scenario. The red color shows the damaged locations and the blue is used for the intact ones: (a) actual damage location, (b) predictions of one-class gradient boosting algorithm, (c) prediction of one-class SVM, (d) predictions of clustering technique using GMM, (e) predictions of our proposed KDND algorithm.

381 of this graph as a result of changes in the characteristics of those RVs. The main advantage of graph structure
 382 morphing using the kernel dependence technique is its robustness with respect to the dimensionality of the
 383 data sets. Therefore, the proposed KDND method is applicable to arbitrarily high dimensional data sets.

384 The experimental results of applying the proposed algorithm to realistic SHM application problems shows
 385 that considering the dependencies of the relevant RVs by tracking their dependence structures can potentially
 386 yield more accurate classification results compared to traditional classification based on tracking the changes
 387 in the marginal distributions of the RVs. Followed by the results, the KDND method reduced the false
 388 positive rates between 14% and 33% compared to peer techniques such as gradient boosting method, which

Table 2: Comparing the performance of KDND with one-class gradient boosting algorithm, one-class SVM, and clustering method using the false positive reduction criteria that was defined in section 7.1.3

FP reduction compared to (%)	Damage scenarios				
	1	2	3	4	5
one-class gradient boosting	32.5	33.3	36.1	29.6	29.6
one-class SVM	20.9	16.6	19.4	14.8	14.8
one-class clustering	20.9	22.2	14.8	22.2	18.5

389 is another robust method with respect to the dimensionality of data sets, one-class SVM, and clustering via
 390 GMM.

391 The main trade-off of using the proposed technique over the alternative methods is its higher computa-
 392 tional demand as a result of computing kernel matrices and iterative operations on such matrices. Also, only
 393 a few dimensionality reduction methods, such as random subspace feature selection which are also usually
 394 demanding, can be used along with the KDND classifier. Moreover, the formulation of the KDND algorithm
 395 requires a fixed set of features to be used for all RVs. This can be viewed as an additional constraint when
 396 it comes to feature selection. Thus, developing specific feature selection techniques for the KDND classifier,
 397 and similar techniques that impose the same constraint, can be pursued in future studies.

398 9. Acknowledgment

399 The authors acknowledge the support provided by Royal Dutch Shell through the MIT Energy Initiative,
 400 and thank chief scientists Dr. Dirk Smit and Dr. Sergio Kapusta.

401 Appendix A. MMD and HSIC for comparing Gaussian RVs

402 Consider two Gaussian RV \mathbf{z}_1 and \mathbf{z}_2 , where $\mathbf{z}_1 \sim \mathcal{N}(0, 1)$ and $\mathbf{z}_2 \sim \mathcal{N}(\mu, \sigma)$, and $\mathcal{N}(\cdot, \cdot)$ denotes Gaussian
 403 RV with the mean and variance as its first and second arguments, respectively. The reason for considering
 404 Gaussian RVs is that their correlation coefficient is an exact measure of their dependency. To study the
 405 variations of HSIC and MMD for Gaussian RVs, we keep the parameters of \mathbf{z}_1 unchanged while changing
 406 the parameters of \mathbf{z}_2 . For each new set of parameters of \mathbf{z}_2 we draw two sets of samples from \mathbf{z}_1 and \mathbf{z}_2 and
 407 compute their associated HSIC and MMD.

408 We consider 41 linearly spaced values between 0.0 and 3.0 for changing μ . For each new value of μ , we
 409 use 20 different values, which are linearly spaced between 0.0 and 1.0, as the correlation coefficient between
 410 \mathbf{z}_1 and \mathbf{z}_2 . For a desired correlation coefficient, we adjust σ accordingly. Thus, the new parameter of \mathbf{z}_2 can
 411 sit on a grid of 41×20 . For each pair of (μ, σ) on this grid, we run 30 simulations by sampling from the
 412 distributions of \mathbf{z}_1 and \mathbf{z}_2 with its new parameters and computing the associated HSIC and MMD of the two
 413 RVs. The result of these simulations is shown in Figure A.11. It follows from this figure that the HSIC is

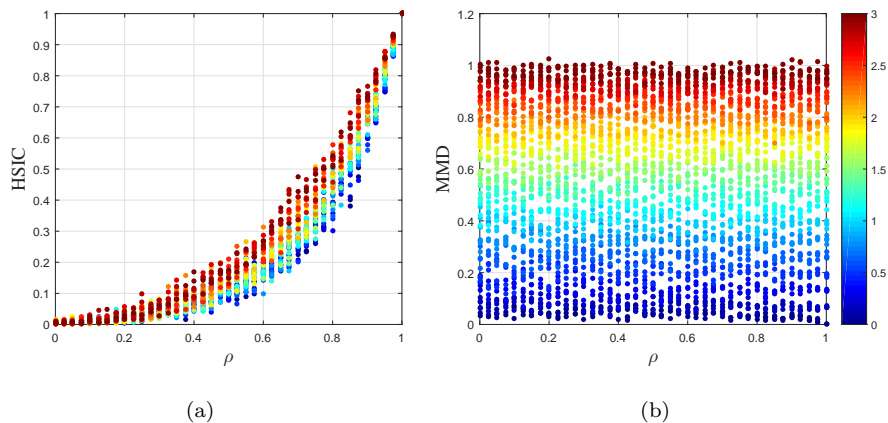


Figure A.11: Variations of (a) HSIC and (b) MMD as a result of changing the mean value and correlation coefficient of two Gaussian RV. The distance between the mean values of these RV is color coded such that dark blue is used for zero distance between the mean values, red is used for a shift of mean value of 3, and other colors are in between these extreme cases.

414 monotonically changed from zero to unity as the correlation coefficient, ρ , varies from zero to one; however,
 415 the independence criterion is not as sensitive to the shift of mean values of the Gaussian RVs. In contrast to
 416 HSIC, the MMD is capable of capturing the translational discrepancy between two Gaussian clusters, while it
 417 is not sensitive to the change of their correlation. Therefore, by using both of these measures, as suggested in
 418 Section 6.1, we should be able to track the translational discrepancies as well as the change of dependencies
 419 between two clusters of data.

420 Another important characteristic of these measures is the relation between their magnitude and variance.
 421 Figure 12(a) shows the same data as Figure 11(b), but in logarithmic scale. As is shown, the variation of
 422 HSIC increases when its magnitude decreases. This can be quantified by the coefficient of variation (CoV) of
 423 the HSIC which is shown in Figure 12(b) for different correlation coefficients and shifts of mean values. μ_h
 424 and σ_h in this plot are, respectively, the mean value and the standard deviation of the HSIC for a given mean
 425 shift and correlation coefficient between \mathbf{z}_1 and \mathbf{z}_2 . Due to the high variations of HSIC for weak dependencies,
 426 direct use of the likelihood difference stated in (5c) in our classification problem results in the dominance of
 427 the weakly dependent variables in the final decision making. This is in contrast with our main objective; thus,
 428 we proposed the voting strategy to avoid the direct use of likelihood ratios in our decision making process.

429 References

430 References

- 431 [1] M. A. Pimentel, D. A. Clifton, L. Clifton, L. Tarassenko, A review of novelty detection, *Signal Processing*
 432 99 (2014) 215–249.
- 433 [2] F. Angiulli, C. Pizzuti, Support vector method for novelty detection, in: *Proceedings of the 6th European*
 434 *Conference on Principles of Data Mining and Knowledge Discovery, PKDD '02, 2002*, pp. 15–26.

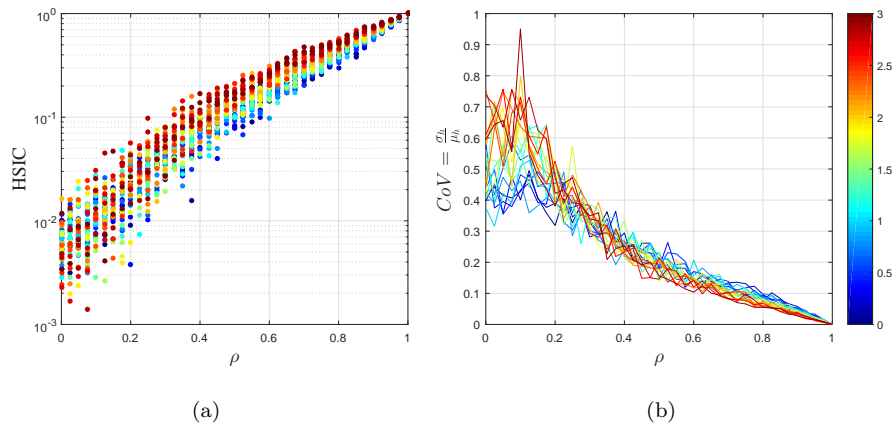


Figure A.12: The relationship between the magnitude and variance of the HSIC for the Gaussian RVs. (a) variation of HSIC as a function of correlation coefficient and shift of mean values of two Gaussian clusters in log-scale, (b) coefficient of variation for HSIC. The color coded lines and points follow the same rule as in Figure A.11

- 435 [3] B. Scholkopf, R. Williamson, A. Smola, J. Shawe-Taylor, J. Platt, Support vector method for novelty
436 detection, *Advances in Neural Information Processing Systems*, MIT Press 3 (2000) 582–588.
- 437 [4] M. Markou, S. Singh, Novelty detection: a review - part 2: neural network based approaches, *Signal*
438 *Processing*, Elsevier 83 (2003) 2499–2521.
- 439 [5] P. Paalanen, J.-K. Kamarainen, J. Ilonen, H. Kalviainen, Feature representation and discrimination
440 based on gaussian mixture model probability densities—practices and algorithms, *Pattern Recognition*,
441 Elsevier 39 (2006) 1346–1358.
- 442 [6] C. Desir, S. B. an Caroline Petitjean, H. Laurent, One class random forests, *Pattern Recognition*, Elsevier
443 46 (2013) 3490–3506.
- 444 [7] Y. Sheikh, M. Shah, Bayesian object detection in dynamic scenes, in: *Proceedings of the 2005 IEEE*
445 *Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, IEEE, 2005.
- 446 [8] R. Mohammadi-Ghazi, J. Chen, O. Buyukozturk, Pairwise graphical models for structural health moni-
447 toring with dense sensor arrays, *Journal of Mechanical Systems and Signal Processing* 93 (2017) 578–592.
- 448 [9] R. Mohammadi-Ghazi, O. Buyukozturk, Non-planar ising graphical model for efficient inference in struc-
449 tural health monitoring, in: *10th International Workshop on Structural Health Monitoring (IWSHM*
450 *2015)*, Stanford, CA, USA, 2015.
- 451 [10] R. Mohammadi-Ghazi, Y. M. Marzouk, O. Buyukozturk, Conditional classifiers and boosted conditional
452 gaussian mixture model for novelty detection, *Pattern Recognition*, Elsevier 81 601–614.
- 453 [11] S. Mahamud, Comparing belief propagation and graph cuts for novelty detection, in: *Proceedings of*

- 454 the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06),
455 IEEE, 2006.
- 456 [12] T. Hastie, R. Tibshirani, J. Friedman, The Elements of Statistical Learning: Data Mining, Inference,
457 and Prediction, 2nd Edition, Springer Series in Statistics, 2009.
- 458 [13] A. Smola, A. Gretton, L. Song, B. Scholkopf, A hilbert space embedding for distributions, 18th Inter-
459 national Conference on Algorithmic Learning Theory, ALT 2007 (2007) 13–31.
- 460 [14] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Scholkopf, A. Smola, A kernel method for the two-sample
461 problem, Journal of Machine Learning Research 1 (2008) 1–43.
- 462 [15] A. Gretton, K. Fukumizu, Z. Harchaoui, B. K. Sriperumbudur, A fast, consistent kernel two-sample test,
463 Advances in Neural Information Processing Systems 22 (NIPS 2009), Red Hook, NY.
- 464 [16] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Scholkopf, A. Smola, A kernel two-sample test, Journal
465 of Machine Learning Research 13 (2012) 723–773.
- 466 [17] K. Chwialkowski, A. Gretton, A kernel independence test for random processes, 31th International
467 Conference on Machine Learning, Beijing, China 32.
- 468 [18] M. D. Lozzo, A. Marrel, New improvements in the use of dependence measures for sensitivity analysis
469 and screening, Journal of Statistical Computation and Simulation (2014) 1–21 [doi:arXiv:1412.1414v1](https://doi.org/10.1080/15227011.2014.914141).
- 470 [19] S. D. Veiga, Global sensitivity analysis with dependence measures, Journal of Statistical Computation
471 and Simulation 85 (2014) 1283–1305.
- 472 [20] V. J. Hodge, J. Austin, A survey of outlier detection methodologies, Artificial Intelligence Review 22
473 (2004) 85–126.
- 474 [21] i. Steinwart, The influence of the kernel on the consistency of support vector machines, Journal of
475 Machine Learning Research 2.
- 476 [22] T. Hofmann, B. Scholkopf, A. J. Smola, Kernel methods in machine learning, The Annals of Statistics
477 36 (2008) 1171–1220.
- 478 [23] A. Smola, A. Gretton, L. Song, B. Schölkopf, A hilbert space embedding for distributions, in: Algorithmic
479 Learning Theory: 18th International Conference, 2007.
- 480 [24] L. Song, Learning via hilbert space embedding of distributions, Ph.D. thesis, University of Sydney
481 (2008).
- 482 [25] A. Shamsheyeva, A. Sowmya, The anisotropic gaussian kernel for svm classification of hrct images of
483 the lung, Intelligent Sensors, Sensor Networks and Information Processing Conference, IEEE.

- 484 [26] D. Brodic, Optimization of the anisotropic gaussian kernel for text segmentation and parameter extrac-
485 tion, *Theoretical Computer Science (TCS)* 323 (2010) 140–152.
- 486 [27] F. Aioli, M. Donini, Learning anisotropic rbf kernels, *Artificial Neural Networks and Machine Learning*,
487 ICANN.
- 488 [28] J. Rice, *Mathematical Statistics and Data Analysis*, 3rd Edition, Duxbury Press: Belmont, CA, 2007.
- 489 [29] R. J. Hyndman, Computing and graphing highest density regions, *The American Statistician*, Taylor
490 and Francis 50 (1996) 120–126.
- 491 [30] N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, SMOTE: Synthetic minority over-sampling
492 technique, *Journal of Artificial Intelligence Research* 16 (2002) 321–357.
- 493 [31] E. Figueiredo, G. Park, C. R. Farrar, K. Worden, J. Figueiras, Machine learning algorithms for damage
494 detection under operational and environmental variability, *Structural Health Monitoring* 10.
- 495 [32] C. R. Farrar, K. Worden, *Structural Health Monitoring: A Machine Learning Perspective*, John Wiley
496 Sons, 2012.
- 497 [33] K. Worden, G. Manson, The application of machine learning to structural health monitoring, *Philosophical*
498 *Transactions of the Royal Society A* (2007) 515–537.
- 499 [34] K. Worden, G. Manson, The application of machine learning to structural health monitoring, *Philosophical*
500 *Transactions of Royal Society A*.
- 501 [35] Y. Ying, J. H. G. Jr., I. J. Oppenheim, L. Soibelman, J. B. Harley, J. Shi, Y. Jin, Toward data-driven
502 structural health monitoring: Application of machine learning and signal processing to damage detection,
503 *Journal of Computing in Civil Engineering* 27.
- 504 [36] L. Bornn, C. R. Farrar, G. Park, Damage detection in initially nonlinear systems, *International Journal*
505 *of Engineering Science* 48 (2010) 909–920.
- 506 [37] S. da Silva, M. D. Junior, V. L. Junior, M. J. Brennan, Structural damage detection by fuzzy clustering,
507 *Mechanical Systems and Signal Processing* 22 (2008) 1636–1649.
- 508 [38] M. H. Rafiei, H. Adelib, A novel unsupervised deep learning model for global and local health condition
509 assessment of structures, *Engineering Structures* 156 (2018) 598–607.
- 510 [39] D. J. Fleet, A. D. Jepson, Computation of component image velocity from local phase information, *Int.*
511 *J. Comput. Vision* 5 (1) (1990) 77–104. doi:10.1007/BF00056772.
512 URL <http://dx.doi.org/10.1007/BF00056772>

- 513 [40] T. Gautama, M. Van Hulle, A phase-based approach to the estimation of the optical flow field using
514 spatial filtering, *Neural Networks, IEEE Transactions on* 13 (5) (2002) 1127 – 1136. doi:10.1109/TNN.
515 2002.1031944.
- 516 [41] E. P. Simoncelli, W. T. Freeman, The steerable pyramid: A flexible architecture for multi-scale derivative
517 computation, in: *icip, IEEE*, 1995, p. 3444.
- 518 [42] N. Wadhwa, M. Rubinstein, F. Durand, W. T. Freeman, Phase-based video motion processing, *ACM*
519 *Trans. Graph. (Proceedings SIGGRAPH 2013)* 32 (4).
- 520 [43] E. Figueiredo, J. Figueiras, G. Park, C. R. Farrar, K. Worden, Influence of the autoregressive model
521 order on damage detection, *Computer-Aided Civil and Infrastructure Engineering* 26 (2011) 225–238.
- 522 [44] S. Khazai, S. Homayouni, A. Safari, B. Mojaradi, Anomaly detection in hyperspectral images based
523 on an adaptive support vector method, *IEEE GEOSCIENCE AND REMOTE SENSING LETTERS* 8
524 (2011) 646–650.