

Article

The Quality of the Covariance Selection through Detection Problem and AUC Bounds

Navid Tafaghodi Khajavi * and Anthony Kuh

Department of Electrical Engineering,
University of Hawaii, Honolulu, HI 96822,
Email: {navidt, kuh}@hawaii.edu

* Correspondence: navidt@hawaii.edu

Version November 15, 2017 submitted to MDPI

Abstract: This paper considers the problem of quantifying the quality of a model selection problem for a graphical model. The model selection problem often uses a distance measure such as the Kulback-Leibler (KL) distance to quantify the quality of the approximation between the original distribution and the model distribution. We extend this work by formulating the problem as a detection problem between the original distribution and the model distribution. In particular, we focus on the covariance selection problem by Dempster, [1], and consider the cases where the distributions are Gaussian distributions. Previous work showed that if the approximation model is a tree, that the optimal tree that minimizes the KL divergence can be found by using the Chow-Liu algorithm [2]. While the algorithm minimizes the KL divergence it does not minimize other measures such as other divergences and the area under the curve (AUC). These measures all depend on the eigenvalues of the correlation approximation measure (CAM). We find expressions for KL divergence, log-likelihood ratio, and AUC as a function of the CAM. Easily computable upper and lower bounds are also found for the AUC. The paper concludes by computing these measures for real and synthetic simulation data.

Keywords: covariance selection; model approximation; detection problem; area under the curve; information divergences

1. Introduction

Graphical models are useful tools for describing the geometric structure of networks in numerous applications such as energy, social, sensor, biological, and transportation networks [3] that deal with high dimensional data. Learning from these high dimensional data requires large computation power which is not always available [4], [5]. The hardware limitation for different applications force us to compromise between the accuracy of the learning algorithm and its time complexity by using the best possible approximation algorithm given the constrained graph. In other words, the main concern is to compromise between model complexity and its accuracy by choosing a simpler, yet informative model. There are lots of approximation algorithms that are proposed for model selection to impose structure given data. For the Gaussian distribution, the covariance selection problem is presented and studied in [1] and [6].

The ultimate purpose of the covariance selection problem is to reduce the computational complexity in various applications. One of the special approximation models is the tree approximation model. Tree approximation algorithms are among the algorithms that reduce the number of computations to get quicker approximate solutions to a variety of problems. If a tree model is used, then distributed estimation algorithms such as message passing algorithm [7] and the belief propagation algorithm [8] can easily be applied and are guaranteed to converge to the maximum likelihood solution.

There are algorithms in the literature such as the Chow-Liu minimum spanning tree (MST) [2], the first order Markov chain approximation [9] and penalized likelihood methods such as LASSO [10] and

graphical LASSO [11] that can be used to approximate the correlation matrix and the inverse correlation matrix with a more sparse graph representation while retaining good accuracy. The Chow-Liu MST algorithm for Gaussian distribution is to find the optimal tree structure using a Kullback-Leibler (KL) divergence cost function [1]. The Chow-Liu MST algorithm constructs a weighted graph by computing pairwise mutual informations and then utilizes one of the MST algorithms such as the Kruskal algorithm [12] or the Prim algorithm [13]. The first order Markov chain approximation uses a regret cost function to output first order Markov chain structured graph [9] by utilizing a greedy type algorithm. Penalized likelihood methods uses an L1-norm penalty term in order to sparsify the graph representation and eliminate some edges. Recently, a tree approximation in a linear, underdetermined model is proposed in [14] where the solution is based on expectation, maximization (EM) algorithm combined with the Chow Liu algorithm.

Sparse modeling has many applications in distributed signal processing and machine learning over graphs. One of its applications is for the electric power grid at the distribution level. The *smart grid* is a promising solution that delivers reliable energy to consumers through the power grid when there are uncertainties such as distributed renewable energy generation sources. Smart grid technologies such as smart meters and communication links are added to the distribution grid in order to obtain the high dimensional, real-time data and information and overcome uncertainties and unforeseen faults. The future grid will incorporate distributed renewable energy generation such as solar photovoltaics (PV), with these energy sources being highly correlated. Thus, modeling is an essential part for signal processing and implementation of the smart grid.

We discuss the quality of the model selection problem, focusing on the Gaussian case, i.e. covariance selection problem. We ask the following important question: “*is the covariance approximation of the covariance matrix for the Gaussian model a good approximation?*” To answer this question, we need to pick a closeness criterion which has to be coherent and general enough to handle a wide variety of problems and also have asymptotic justification [15]. In many applications the Kullback-Leibler (KL) divergence has been proposed as a closeness criterion between the original distribution and its model approximation distribution [1] and [2]. Besides that, other closeness measures and divergences are used for the model selection problem. One example is the use of reverse KL divergence as the closeness criterion in variational methods to learn the desired approximation structure [16].

In this paper we bring a different perspective to the model approximation problem by formulating a general detection problem. The detection problem leads to calculation of the log-likelihood ratio test (LLRT) statistic, the receiver operating characteristic (ROC) curve, the KL divergence and the reverse KL divergence as well as the area under the curve (AUC) where the AUC is used as the accuracy measure for the detection problem. The detection problem formulation gives us a broader view as well as looking at different approaches of determining whether a particular model is a good approximation or not. More specifically, the AUC gives us additional incite about any approximation since it is a way to formalize the model approximation problem. For Gaussian data, the LLRT statistic simplifies to an indefinite quadratic form. A key quantity that we define is the correlation approximation matrix (CAM) as the product of the original correlation matrix and the inverse of the model approximation correlation matrix. For Gaussian data this matrix contains all the information needed to compute the information divergences, the ROC curve and the area under this curve, i.e. the AUC. We also show the relationship between the CAM, the AUC and the Jeffreys divergence [17], the KL divergence and the reverse KL divergence. We present an analytical expression to compute the AUC for a given CAM that can be efficiently evaluated numerically. We show the relation between the AUC, the KL divergence, the LLRT statistics and the ROC curve. We also present analytical upper and lower bounds for the AUC which are only depend on eigenvalues of the CAM. Throughout the discussion section, we pick the tree approximation model as a well-known subset of all graphical models. The tree approximations is considered since they are widely used in literature and it is much simpler performing inference and estimation on trees rather than graphs that have cycles or loops. We perform simulations over synthetic and real data for several examples to explore and discuss our results. Simulation results

87 indicate that $1 - \text{AUC}$ is decreasing exponentially as the number of nodes of the graph increases which
 88 is consistent with analytical results obtained from the AUC upper and lower bounds.

89 The rest of this paper is organized as follows. In section 2 we give a general framework for the
 90 detection problem and the corresponding sufficient test statistic, the log-likelihood ratio test. Moreover,
 91 the sufficient test statistic for Gaussian data as well as its distribution under both hypotheses are also
 92 presented in this section. The ROC curve and the AUC definition as well as analytical expression for
 93 the AUC are given in section 3. Section 4 provides analytical lower and upper bounds for the AUC.
 94 The lower bound for the AUC uses the Chernoff bound and is a function of the CAM eigenvalues.
 95 The upper bound is obtained by finding a parametric relationship between the AUC and the KL and
 96 reverse KL divergences. Moreover, Section 5 presents the tree approximation model and provides
 97 some simulations over synthetic examples as well as real solar data examples and investigates quality
 98 of the tree approximation based on the numerically evaluated AUC and also its analytical upper and
 99 lower bounds. Finally, Section 6 summarizes results of this paper.

100 2. Detection Problem Framework

101 In this section, we present a framework to quantify the quality of a model selection problem.
 102 More specifically, we formulate a detection problem to distinguish between the covariance matrix of
 103 a multivariate normal distribution and an approximation of the aforementioned covariance matrix
 104 based on the given model.

105 2.1. Model selection problem

106 We want to approximate a multivariate distribution by the product of lower order component
 107 distributions [18]. Let random vector $\underline{X} \in \mathbb{R}^n$, have a distribution with parameter Θ , i.e. $\underline{X} \sim f_{\underline{X}}(\underline{x})$.
 108 We want to approximate the random vector \underline{X} , with another random vector associated with the desired
 109 model¹. Let the model random vector $\underline{X}_{\mathcal{M}} \in \mathbb{R}^n$ have a distribution with parameter $\Theta_{\mathcal{M}}$, associated
 110 with the desired model, i.e. $\underline{X}_{\mathcal{M}} \sim f_{\underline{X}_{\mathcal{M}}}(\underline{x})$. Also, let $\mathcal{G} = (\mathcal{V}, \mathcal{E}_{\mathcal{M}})$ be the graph representation of the
 111 model random vector $\underline{X}_{\mathcal{M}}$ where sets \mathcal{V} and $\mathcal{E}_{\mathcal{M}}$ are the set of all vertices and the set of all edges of the
 112 graph representing of $\underline{X}_{\mathcal{M}}$, respectively. Moreover, $\mathcal{E}_{\mathcal{M}} \subseteq \psi$ where ψ is the set of all edges of complete
 113 graph with vertex set \mathcal{V} .

114 **Remark:** Covariance selection is presented in [1]. Moreover, tree model as a special case for the model
 115 selection problem is discussed in subsection 5.1.

116 2.2. General detection framework

117 The model selection problem is extensively studied in the literature [1]. In many state of the art
 118 works, minimizing the KL divergence between two distributions or the maximum likelihood criterion
 119 are proposed in order to quantify the quality of the model approximation. A different way to look
 120 at the problem of quantifying the quality of the model approximation algorithm is to formulate the
 121 problem as a detection problem [19]. Given the set of data in the detection problem, the goal is to
 122 distinguish between two hypotheses, *the null hypothesis* and *the alternative hypothesis*. To set up a
 123 detection problem, we need to define these two hypotheses for the model selection problem as follow

- 124 - The null hypothesis, \mathcal{H}_0 : The hypothesis that data is generated using the known distribution,
- 125 - The alternative hypothesis, \mathcal{H}_1 : The hypothesis that data is generated using the model
- 126 approximation distribution.

¹ Examples of possible models: tree structure, sparse structure and Markov chain.

Given the set up for the null hypothesis and the alternative hypothesis, we need to define a test statistic to quantify the detection problem. The likelihood ratio test (the Neyman-Pearson (NP) Lemma [20]) is the most powerful test statistic where we first define the log-likelihood ratio test (LLRT) as

$$l(\underline{x}) = \log \frac{f_{\underline{X}}(\underline{x}|\mathcal{H}_1)}{f_{\underline{X}}(\underline{x}|\mathcal{H}_0)} = \log \frac{f_{\underline{X},\mathcal{M}}(\underline{x})}{f_{\underline{X}}(\underline{x})}$$

127 where $f_{\underline{X}}(\underline{x}|\mathcal{H}_0)$ is the distribution of random vector \underline{X} under the null hypothesis while $f_{\underline{X}}(\underline{x}|\mathcal{H}_1)$ is
128 the distribution of random vector \underline{X} under the alternative hypothesis.

129 We then define *the false-alarm probability* and *the detection probability* by comparing the LLRT statistic
130 under each hypothesis with a given threshold, τ , and computing the following probabilities

- 131 - The false-alarm probability, $P_0(\tau)$, under the null hypothesis, \mathcal{H}_0 : $P_0(\tau) = \Pr(L(\underline{X}) \geq \tau|\mathcal{H}_0)$,
- 132 - The detection probability, $P_1(\tau)$, under the alternative hypothesis, \mathcal{H}_1 : $P_1(\tau) = \Pr(L(\underline{X}) \geq \tau|\mathcal{H}_1)$,
- 133 where random variable $L(\underline{X})$ is the LLRT statistic random variable.

134 The most powerful test is defined by setting the false-alarm rate $P_0(\tau) = \bar{P}_0$ and then computing the
135 threshold value τ_0 such that $\Pr(L(\underline{X}) \geq \tau_0|\mathcal{H}_0) = \bar{P}_0$.

Definition 1. *The KL divergence between two multivariate continuous distributions $p(\underline{X})$ and $q(\underline{X})$ is defined as*

$$\mathcal{D}(p_{\underline{X}}(\underline{x})||q_{\underline{X}}(\underline{x})) = \int_{\mathcal{X}} p_{\underline{X}}(\underline{x}) \log \frac{p_{\underline{X}}(\underline{x})}{q_{\underline{X}}(\underline{x})} d\underline{x}$$

136 where \mathcal{X} is the feasible set. ■

137 Throughout this paper we may use other notations such as the KL divergence between two
138 random variable or the KL divergence between two covariance matrices for zero-mean Gaussian
139 distribution case in order to present the KL divergence between two distributions.

140 **Proposition 1.** *Expectation of the LLRT statistic under each hypothesis is*

- 141 - $E(L(\underline{X})|\mathcal{H}_0) = -\mathcal{D}(f_{\underline{X}}(\underline{x}|\mathcal{H}_0)||f_{\underline{X}}(\underline{x}|\mathcal{H}_1)) = -\mathcal{D}(f_{\underline{X}}(\underline{x})||f_{\underline{X},\mathcal{M}}(\underline{x}))$,
- 142 - $E(L(\underline{X})|\mathcal{H}_1) = \mathcal{D}(f_{\underline{X}}(\underline{x}|\mathcal{H}_1)||f_{\underline{X}}(\underline{x}|\mathcal{H}_0)) = \mathcal{D}(f_{\underline{X},\mathcal{M}}(\underline{x})||f_{\underline{X}}(\underline{x}))$.

143 **Proof.** Proof is based on the KL divergence definition. ■

144 **Remark:** Relationship between the NP lemma and the KL divergence is previously stated in [21]
145 with the similar straightforward calculation, where the LLRT statistic loses power when the wrong
146 distribution is used instead of the true distribution for one of these hypotheses.

147 In a regular detection problem framework, the NP decision rule is to accept the hypothesis \mathcal{H}_1 if
148 the LLRT statistic, $l(\underline{x})$, exceeds a critical value, and reject it otherwise. Moreover, the critical value is
149 set based on the rejection probability of the hypothesis \mathcal{H}_0 , i.e. false-alarm probability. Note that, we
150 pursue a different goal in the approximation problem scenario. We approximate a model distribution,
151 $f_{\underline{X},\mathcal{M}}(\underline{x})$, as close as possible to the given distribution, $f_{\underline{X}}(\underline{x})$. The closeness criterion is based on the
152 modified detection problem framework where we compute the LLRT statistic and compare it with a
153 threshold. In an ideal case where there is no approximation error, the detection probability must be
154 equal to the false-alarm probability for the optimal detector at all possible thresholds, i.e. the receiver
155 operating characteristic (ROC) curve [22] that represents best detectors for all threshold values should
156 be a line of slope 1 passing through the origin.

157 In the next, we assume that the random vector \underline{X} has zero-mean Gaussian distribution. Thus, the
158 covariance matrix of the random vector \underline{X} is the parameter of interest in the model selection problem,
159 i.e. covariance selection.

160 2.3. Multivariate Gaussian distribution

Let random vector $\underline{X} \in \mathbb{R}^n$, have a zero-mean jointly Gaussian distribution with covariance matrix $\Sigma_{\underline{X}}$, i.e. $\underline{X} \sim \mathcal{N}(\underline{0}, \Sigma_{\underline{X}})$ where the covariance matrix $\Sigma_{\underline{X}}$ is positive-definite, $\Sigma_{\underline{X}} \succ 0$. In this paper, the null hypothesis, \mathcal{H}_0 , is the hypothesis that the parameter of interest is known and is equal to $\Sigma_{\underline{X}}$ while the alternative hypothesis, \mathcal{H}_1 , is the hypothesis that the random vector \underline{X} is replaced by the model random vector $\underline{X}_{\mathcal{M}}$. In this scenario, the model random vector $\underline{X}_{\mathcal{M}}$ has a zero-mean jointly Gaussian distribution (the model approximation distribution) with covariance matrix $\Sigma_{\underline{X}_{\mathcal{M}}}$, i.e. $\underline{X}_{\mathcal{M}} \sim \mathcal{N}(\underline{0}, \Sigma_{\underline{X}_{\mathcal{M}}})$ where the covariance matrix $\Sigma_{\underline{X}_{\mathcal{M}}}$ is also positive-definite, $\Sigma_{\underline{X}_{\mathcal{M}}} \succ 0$. Thus, the LLRT statistic for the jointly Gaussian random vectors, \underline{X} and $\underline{X}_{\mathcal{M}}$, is simplified as

$$l(\underline{x}) = \log \frac{\mathcal{N}(\underline{0}, \Sigma_{\underline{X}_{\mathcal{M}}})}{\mathcal{N}(\underline{0}, \Sigma_{\underline{X}})} = -c + k(\underline{x}) \quad (1)$$

161 where $c = -\frac{1}{2} \log(|\Sigma_{\underline{X}} \Sigma_{\underline{X}_{\mathcal{M}}}^{-1}|)$ is a constant and $k(\underline{x}) = \underline{x}^T \mathbf{K} \underline{x}$ where $\mathbf{K} = \frac{1}{2}(\Sigma_{\underline{X}}^{-1} - \Sigma_{\underline{X}_{\mathcal{M}}}^{-1})$ is an indefinite
162 matrix with both positive and negative eigenvalues.

163 We define the correlation approximation matrix (CAM) associated with the covariance selection
164 problem and dissimilarity parameters of the CAM as follows.

165 **Definition 2. Correlation approximation matrix.** The CAM for the covariance selection problem is defined
166 as $\Delta \triangleq \Sigma_{\underline{X}} \Sigma_{\underline{X}_{\mathcal{M}}}^{-1}$ where $\Sigma_{\underline{X}_{\mathcal{M}}}$ is the model covariance matrix. ■

167 **Definition 3. Dissimilarity parameters for covariance selection problem.** Let $\alpha_i \triangleq \lambda_i + \lambda_i^{-1} - 2$ for
168 $i \in \{1, \dots, n\}$ be dissimilarity parameters of the CAM correspond to the covariance selection problem where
169 $\lambda_i > 0$ for $i \in \{1, \dots, n\}$ are eigenvalues of the CAM. ■

170 **Remark:** The CAM is a positive definite matrix. Moreover, eigenvalues of the CAM contains all
171 information necessary to compute cost functions associated with the model selection problem.

172 **Theorem 1. Covariance Selection [1].** Given a multivariate Gaussian distribution with covariance matrix
173 $\Sigma_{\underline{X}} \succ 0$, $f_{\underline{X}}(\underline{x})$, and a model \mathcal{M} , there exists a unique approximated multivariate Gaussian distribution with
174 covariance matrix $\Sigma_{\underline{X}_{\mathcal{M}}} \succ 0$, $f_{\underline{X}_{\mathcal{M}}}(\underline{x})$, that minimize the KL divergence, $\mathcal{D}(f_{\underline{X}}(\underline{x}) || f_{\underline{X}_{\mathcal{M}}}(\underline{x}))$ and satisfies the
175 covariance selection rules, i.e. the model covariance matrix satisfies the following covariance selection rules

$$\begin{aligned} 176 & - \Sigma_{\underline{X}_{\mathcal{M}}}(i, i) = \Sigma_{\underline{X}}(i, i), & \forall i \in \mathcal{V} \\ 177 & - \Sigma_{\underline{X}_{\mathcal{M}}}(i, j) = \Sigma_{\underline{X}}(i, j), & \forall (i, j) \in \mathcal{E}_{\mathcal{M}} \\ 178 & - \Sigma_{\underline{X}_{\mathcal{M}}}^{-1}(i, j) = 0, & \forall (i, j) \in \mathcal{E}_{\mathcal{M}}^c \end{aligned}$$

179 where the set $\mathcal{E}_{\mathcal{M}}^c = \psi - \mathcal{E}_{\mathcal{M}}$ represents the complement of the set $\mathcal{E}_{\mathcal{M}}$.

180 **Proof.** Proof for Gaussian distributions is given in Dempster 1972 paper [1]. ■

181 **Remark:** Using the CAM definition, the constant c can be written as $c = -\frac{1}{2} \log(|\Delta|)$. Moreover, given
182 any covariance matrix and its model covariance matrix satisfying theorem 1, we have $tr(\Delta) = n$. Thus,
183 from the result in theorem 1 and the definition of KL divergence for jointly Gaussian distributions, we
184 conclude $c = \mathcal{D}(f_{\underline{X}}(\underline{x}) || f_{\underline{X}_{\mathcal{M}}}(\underline{x}))$.

185 2.4. Covariance selection example

Here we choose tree approximation model as an example. Figure 1 indicates two graphs: (a) the complete graph and (b) its tree approximation model where edges in the graph represent non-zero coefficients in the inverse of the covariance matrix [1]. The correlation coefficient between each pair of adjacent nodes has been written on each edge. The correlation coefficient between each pair of

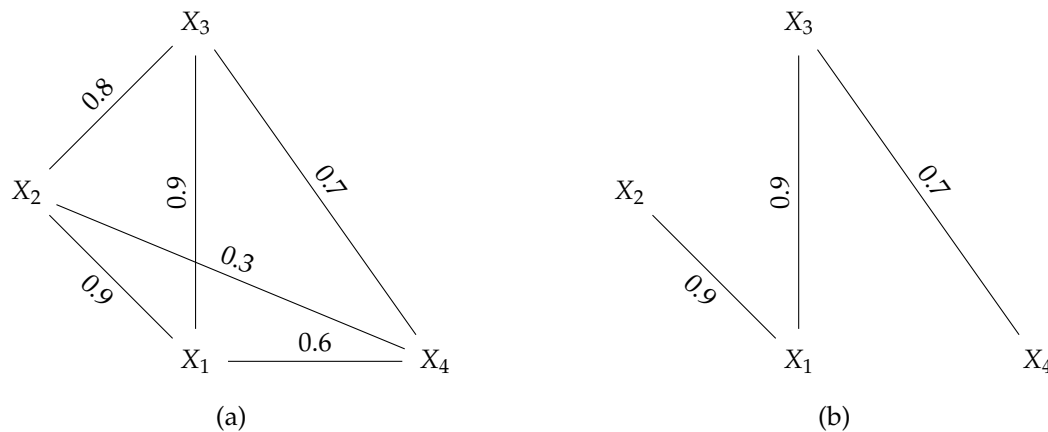


Figure 1. (a) The complete graph; (b) The tree approximation of the complete graph.

nonadjacent nodes is the multiplication of all correlations on the unique path that connects those nodes. The correlation matrix for each graph is

$$\Sigma_{\underline{X}} = \begin{bmatrix} 1 & 0.9 & 0.9 & 0.6 \\ 0.9 & 1 & 0.8 & 0.3 \\ 0.9 & 0.8 & 1 & 0.7 \\ 0.6 & 0.3 & 0.7 & 1 \end{bmatrix}$$

and

$$\Sigma_{\underline{X}_T} = \begin{bmatrix} 1 & 0.9 & 0.9 & 0.63 \\ 0.9 & 1 & 0.81 & 0.567 \\ 0.9 & 0.81 & 1 & 0.7 \\ 0.63 & 0.567 & 0.7 & 1 \end{bmatrix}.$$

The CAM for the above example is

$$\Delta = \begin{bmatrix} 1 & 0 & 0.0412 & -0.0588 \\ 0.0474 & 1 & 0.3042 & -0.5098 \\ 0.0474 & -0.0526 & 1 & 0 \\ 0.9789 & -1.2632 & 0.1421 & 1 \end{bmatrix}.$$

186 The CAM contains all information about the tree approximation². Here we assume cases that Gaussian
187 random variables have finite, nonzero variances.

188 **Remark:** Without loss of generality, throughout this paper we are working with normalized correlation
189 matrices, i.e. the diagonal elements of the correlation matrices are normalized to be equal to one.

190 2.5. Distribution of the LLRT statistic

The random vector \underline{X} has Gaussian distribution under both hypotheses \mathcal{H}_0 and \mathcal{H}_1 . Thus under both hypotheses, the real random variable, $K(\underline{X}) = \underline{X}^T \mathbf{K} \underline{X}$ has a generalized chi-squared distribution, i.e. the random variable, $K(\underline{X})$, is equal to a weighted sum of chi-squared random variables with both positive and negative weights under both hypotheses. Let us define $\underline{W} = \Sigma_{\underline{X}}^{-\frac{1}{2}} \underline{X}$ under \mathcal{H}_0 and $\underline{Z} = \Sigma_{\underline{X}_M}^{-\frac{1}{2}} \underline{X}$ under \mathcal{H}_1 , where $\Sigma_{\underline{X}}^{\frac{1}{2}}$ and $\Sigma_{\underline{X}_M}^{\frac{1}{2}}$ are the square root of covariance matrices $\Sigma_{\underline{X}}$ and

² Dissimilarity parameters α_i 's and eigenvalues of CAM contains all information about the tree approximation.

$\Sigma_{\underline{X}_M}$, respectively. Then let random vectors $\underline{W} \sim \mathcal{N}(\underline{0}, \mathbf{I})$ and $\underline{Z} \sim \mathcal{N}(\underline{0}, \mathbf{I})$ have zero-mean Gaussian distributions with the same covariance matrices, \mathbf{I} , where \mathbf{I} is the identity matrix of dimension n . Note that, the CAM is a positive definite matrix with $\lambda_i > 0$ where $1 \leq i \leq n$. Thus, the random variable $K(\underline{X})$, under both hypotheses \mathcal{H}_0 and \mathcal{H}_1 can be written as:

$$K_0(\underline{X}) \triangleq K(\underline{X})|\mathcal{H}_0 = \frac{1}{2} \sum_{i=1}^n (1 - \lambda_i) W_i^2$$

and

$$K_1(\underline{X}) \triangleq K(\underline{X})|\mathcal{H}_1 = \frac{1}{2} \sum_{i=1}^n (\lambda_i^{-1} - 1) Z_i^2$$

respectively, where random variables W_i and Z_i , are the i -th element of random vectors \underline{W} and \underline{Z} , respectively. Moreover, random variables W_i^2 and Z_i^2 , follow the first order central chi-squared distribution. Note that, similarly random variable $L(\underline{X}) \triangleq -c + K(\underline{X})$ is defined under each hypothesis as

$$L_0 \triangleq L(\underline{X})|\mathcal{H}_0 = -c + K_0(\underline{X})$$

and

$$L_1 \triangleq L(\underline{X})|\mathcal{H}_1 = -c + K_1(\underline{X}).$$

191 **Remark:** As a simple consequence of the covariance selection theorem, the summation of weights for
 192 the generalized chi-squared random variable, the expectation of $K(\underline{X})$, is zero under the hypothesis
 193 \mathcal{H}_0 , i.e. $E(K_0(\underline{X})) = \frac{1}{2} \sum_{i=1}^n (1 - \lambda_i) = 0$ [1], and this summation is positive under the hypothesis \mathcal{H}_1 ,
 194 i.e. $E(K_1(\underline{X})) = \frac{1}{2} \sum_{i=1}^n (\lambda_i^{-1} - 1) \geq 0$.

195 3. The ROC Curve and the AUC Computation

196 3.1. The receiver operating characteristic curve

197 The receiver operating characteristic (ROC) curve is the parametric curve where the detection
 198 probability is plotted versus the false-alarm probability for all thresholds, i.e. each point on the ROC
 199 curve represents a pair of $(P_0(\tau), P_1(\tau))$ for a given threshold τ . Set $z = P_0(\tau)$ and $\eta = P_1(\tau)$, the
 200 ROC curve is $\eta = h(z)$. If $P_0(\tau)$ has an inverse function, then the ROC curve is $h(z) = P_1(P_0^{-1}(z))$. In
 201 general, the ROC curve, $h(z)$, has the following properties [22]

- 202 - $h(z)$ is concave and increasing,
- 203 - $h'(z)$ is positive and decreasing,
- 204 - $\int_0^1 h'(z) dz \leq 1$.

205 Note that, for the ROC curve, the slope of the tangent line at a given threshold, $h'(z)$, gives the
 206 likelihood ratio for the value of the test.

Remark: For the ROC curve for our Gaussian random vectors we have $h'(z)$ is positive, continuous and decreasing in interval $[0, 1]$ with right continuity at 0 and left continuity at 1. Moreover,

$$\int_0^1 h'(z) dz = 1$$

207 since $h(0) = 0$ and $h(1) = 1$.

208 **Definition 4.** Let $f_{L_0}(l)$ and $f_{L_1}(l)$ be the probability density function (PDF) of the random variables L_0 and
 209 L_1 , respectively. ■

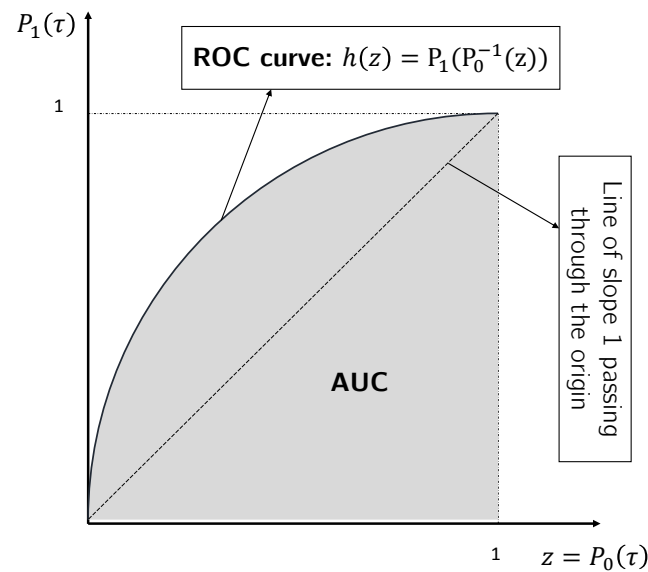


Figure 2. The ROC curve and the area under the ROC curve. Each point on the ROC curve indicates a detector with given detection and false-alarm probabilities.

Lemma 1. Given the ROC curve, $h(z)$, we can compute following KL divergences

$$\mathcal{D}(f_{L_1}(l) || f_{L_0}(l)) = - \int_0^1 \log(h'(z)) dz.$$

and

$$\begin{aligned} \mathcal{D}(f_{L_0}(l) || f_{L_1}(l)) &= - \int_0^1 h'(z) \log(h'(z)) dz \\ &\stackrel{(*)}{=} - \int_0^1 \log\left(\frac{dh^{-1}(\eta)}{d\eta}\right) d\eta \end{aligned}$$

210 where (*) holds if the ROC curve, $\eta = h(z)$, has an inverse function.

211 **Proof.** These results are consequence of the Radon-Nikodým theorem [23]. Simple, alternative calculus
212 based proofs are given in appendix A. ■

213 3.2. Area under the curve

214 As discussed previously we examine the ROC with a goal that the model approximation results
215 with the ROC being a line of slope 1 passing through the origin. This is in contrast to the conventional
216 detection problem where we want to distinguish between the two hypotheses and ideally have an
217 ROC that is a unit step function. Area under the curve (AUC) is defined as the integral of the ROC
218 curve (figure 2) and is a measure of accuracy in decision problems.

Definition 5. The area under the ROC curve (AUC) is defined as

$$AUC = \int_0^1 h(z) dz = \int_0^1 P_1(\tau) dP_0(\tau) \quad (2)$$

219 where τ is the detection problem threshold. ■

220 **Remark:** The AUC is a measure of accuracy for the detection problem and $1/2 \leq \text{AUC} \leq 1$. Note
 221 that, in conventional decision problems, the AUC is desired to be as close as possible to 1 while in
 222 approximation problem presented here we want the AUC to be close to $1/2$.

Theorem 2. Statistical property of AUC [24]. *The AUC for the LLRT statistic is*

$$\text{AUC} = \Pr(L_1 > L_0).$$

223

Corollary 1. *From theorem 2, when PDFs for the LLRT statistic under both hypotheses exist, we can compute the AUC as*

$$\text{AUC} = \int_0^\infty f_{L_1}(l) \star f_{L_0}(l) dl \quad (3)$$

224 where $f_{L_1}(l) \star f_{L_0}(l) \triangleq \int_{-\infty}^\infty f_{L_1}(\tau) f_{L_0}(l + \tau) d\tau$ is the cross-correlation between $f_{L_1}(l)$ and $f_{L_0}(l)$.

225 **Proof.** A proof based on the definition of the AUC (2), is given in [25]. ■

Let us define the difference LLRT statistic random variable as $L_\Delta = L_1 - L_0$. Then, we get

$$\begin{aligned} \text{AUC} &= \Pr(L_\Delta > 0) \\ &= 1 - F_{L_\Delta}(0) \end{aligned}$$

226 where $F_{L_\Delta}(l)$ is the cumulative distribution function (CDF) for random variable L_Δ .

The two conditional random variables L_0 and L_1 are independent³. Thus, the cross-correlation between the corresponding two distributions is the distribution of the difference LLRT statistic, L_Δ . We can write the random variable L_Δ as

$$\begin{aligned} L_\Delta &= -c + K_1(\underline{X}) - (-c + K_0(\underline{X})) \\ &= K_1(\underline{X}) - K_0(\underline{X}). \end{aligned}$$

Replacing the definition for $K_0(\underline{X})$ and $K_1(\underline{X})$, we have

$$L_\Delta = \frac{1}{2} \sum_{i=1}^n (\lambda_i^{-1} - 1) Z_i^2 - \frac{1}{2} \sum_{i=1}^n (1 - \lambda_i) W_i^2.$$

We can rewrite the difference LLRT statistic, L_Δ , in an indefinite quadratic form as

$$L_\Delta = \frac{1}{2} \underline{V}^T (\underline{\Lambda} - \mathbf{I}) \underline{V}$$

where

$$\underline{V} = \begin{bmatrix} \underline{W} \\ \underline{Z} \end{bmatrix}$$

³ By the definition in the detection problem.

231 3.4. Analytical expression for AUC

To compute the CDF of random variable L_{Δ} , we need to evaluate a multi-dimensional integral of jointly Gaussian distributions [28] or we need to approximate this CDF [29]. More efficiently, as discussed in [30] for the real valued case, the CDF of the random variable L_{Δ} can be expressed as a single-dimensional integral of a complex function⁵ in the following form

$$F_{L_{\Delta}}(l) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{e^{\frac{1}{2}(j\omega + \beta)l}}{j\omega + \beta} \frac{1}{\sqrt{|\mathbf{I} + \frac{1}{2}(\mathbf{\Lambda} - \mathbf{I})(j\omega + \beta)|}} d\omega$$

232 where $\beta > 0$ is chosen such that matrix $\mathbf{I} + \frac{\beta}{2}(\mathbf{\Lambda} - \mathbf{I})$, is positive definite and simplifies the evaluation
233 of the multivariate Gaussian integral [30].

Special case: When $\mathbf{\Lambda} = \mathbf{I}$, i.e. the given covariance obeys the model structure, then

$$AUC = 1 - F_{L_{\Delta}}(0) = 1 - \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{1}{j\omega + \beta} = \frac{1}{2}$$

234 for $\beta > 0$ and is also independent of the value of the parameter β .

Picking an appropriate value for the parameter β ⁶, the AUC can be numerically computed by evaluating the following one dimension complex integral

$$AUC = 1 - \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{1}{j\omega + \beta} \frac{1}{\sqrt{|\mathbf{I} + \frac{1}{2}(\mathbf{\Lambda} - \mathbf{I})(j\omega + \beta)|}} d\omega.$$

Furthermore, since $\mathbf{\Lambda} \succ 0$, choosing $\beta = 2$ and changing variable as $\nu = \omega/2$, we have

$$AUC = 1 - \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{1}{j\nu + 1} \frac{1}{\sqrt{|\mathbf{\Lambda} + j\nu(\mathbf{\Lambda} - \mathbf{I})|}} d\nu. \quad (7)$$

235 Moreover, $|\mathbf{\Lambda} + j\nu(\mathbf{\Lambda} - \mathbf{I})| = \prod_{i=1}^p (1 + \alpha_i \nu^2 - j\alpha_i \nu)$. This equation shows that the AUC only depends
236 on α_i 's.

237 **Remark:** Since the AUC integral in (7) can not be evaluated in closed form, it can not be used directly
238 in obtaining model selection algorithms. Numerical evaluation of the AUC using the one dimensional
239 complex integral (7) is very efficient and fast comparing to numerical evaluation of a multi-dimensional
240 integral of jointly Gaussian CDF.

241 4. Analytical Bounds for the AUC

242 As in the previous section we present an analytical expression for the AUC, in this section, we
243 presents analytical lower and upper bounds for the AUC. These bounds will give us insight on how
244 the AUC behave.

245 4.1. Lower bound for the AUC (Chernoff bound application)

246 Given the MGF for the difference LLRT statistic distribution (6), we can apply the Chernoff bound
247 [31] to find a lower bound for the AUC or upper bound for the CDF of the difference LLRT statistic
248 random variable, L_{Δ} , evaluated at zero).

⁵ This is the transform to the frequency domain for an arbitrary β .

⁶ The parameter β is picked such that $\mathbf{I} + \frac{\beta}{2}(\mathbf{\Lambda} - \mathbf{I}) \succ 0$ and $\beta = 2$ always satisfies this condition since $\mathbf{\Lambda} \succ 0$.

Proposition 2. Lower bound for the AUC is

$$\Pr(L_{\Delta} > 0) \geq \max \left\{ \frac{1}{2}, 1 - e^{-\frac{1}{2} \sum_{i=1}^n \log(1 + \frac{\alpha_i}{4})} \right\} \quad (8)$$

Proof. One-half is a trivial lower bound for AUC. To achieve a non-trivial lower bound, we apply Chernoff bound [31] as follow

$$\Pr(L_{\Delta} < 0) \leq \inf_t M_{L_{\Delta}}(t).$$

To complete the proof we need to solve the right-hand-side (RHS) optimization problem.

Step 1: First derivatives of $M_{L_{\Delta}}(t)$ is

$$\begin{aligned} \frac{d}{dt} M_{L_{\Delta}}(t) &= M_{L_{\Delta}}(t) \\ &\left(\frac{1}{2} \sum_{i=1}^n \frac{\lambda_i - 1}{1 - (\lambda_i - 1)t} + \frac{\lambda_i^{-1} - 1}{1 - (\lambda_i^{-1} - 1)t} \right) \\ &= M_{L_{\Delta}}(t) (1 + 2t) \sum_{i=1}^n \frac{\alpha_i}{2(1 - \alpha_i t - \alpha_i t^2)}. \end{aligned}$$

Clearly, first derivative is zero for $t = -1/2$ which is in the feasible domain of the MGF for the difference LLRT statistic. Note that, the smallest feasible domain is $-1 < t < 0$.

Step 2: Second derivatives of $M_{L_{\Delta}}(t)$ is

$$\begin{aligned} \frac{d^2}{dt^2} M_{L_{\Delta}}(t) &= \\ M_{L_{\Delta}}(t) &\left(\frac{1}{4} \sum_{i=1}^n \frac{\lambda_i - 1}{1 - (\lambda_i - 1)t} + \frac{\lambda_i^{-1} - 1}{1 - (\lambda_i^{-1} - 1)t} \right)^2 + \\ M_{L_{\Delta}}(t) &\left(\frac{1}{4} \sum_{i=1}^n \frac{(\lambda_i - 1)^2}{(1 - (\lambda_i - 1)t)^2} + \frac{(\lambda_i^{-1} - 1)^2}{(1 - (\lambda_i^{-1} - 1)t)^2} \right). \end{aligned}$$

Therefore, we conclude that the second derivative is positive and thus the optimal solution to the RHS optimization problem is at $t = -\frac{1}{2}$. Replacing that in the definition of the moment generation function which results in the following bound

$$\Pr(L_{\Delta} \leq 0) < \prod_{i=1}^n \frac{2}{\sqrt{4 + \alpha_i}}$$

which can be written as

$$\Pr(L_{\Delta} > 0) \geq 1 - \prod_{i=1}^n \frac{2}{\sqrt{4 + \alpha_i}}$$

249 which completes the proof. ■

250 4.2. Upper Bound for the AUC

251 In this section, we present a parametric upper bound for the AUC, but first, we need to present
252 the following results.

Lemma 2. Invariance property of the KL divergence for the LLRT statistic. We have

$$\mathcal{D}(f_{L_1}(l) || f_{L_0}(l)) \leq \mathcal{D}(f_{\underline{X}}(\underline{x} | \mathcal{H}_1) || f_{\underline{X}}(\underline{x} | \mathcal{H}_0))$$

and

$$\mathcal{D}(f_{L_0}(l)||f_{L_1}(l)) \leq \mathcal{D}(f_{\underline{X}}(\underline{x}|\mathcal{H}_0)||f_{\underline{X}}(\underline{x}|\mathcal{H}_1)).$$

253 **Proof.** This lemma is an special case of the invariance property of the KL divergence [32]. By picking
 254 appropriate measurable mapping, here appropriate quadratic function for each equation of the above
 255 equations, we conclude the lemma. ■

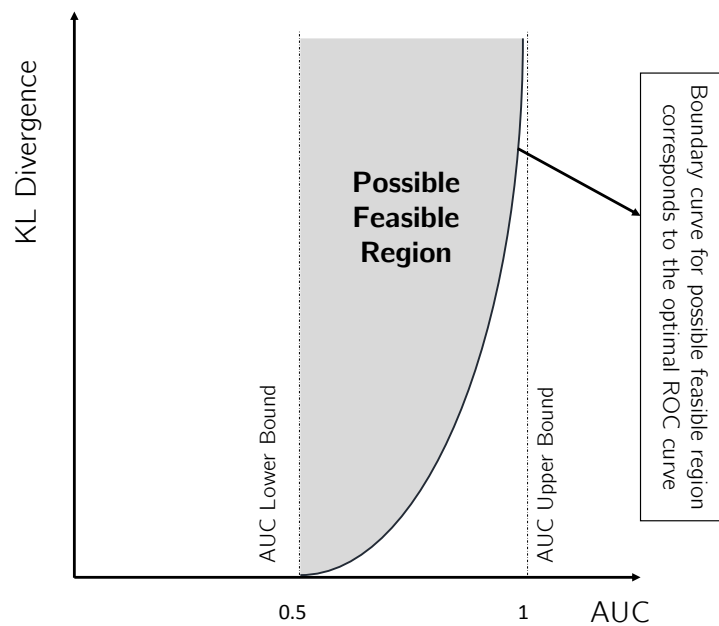


Figure 3. Possible feasible region for the AUC and the Kl divergence pair for all possible detectors or equivalently all possible ROC curves (the KL divergence is between the LLRT statistics under different hypotheses, i.e. $\mathcal{D}(f_{L_0}(l)||f_{L_1}(l))$ or $\mathcal{D}(f_{L_1}(l)||f_{L_0}(l))$.)

256 **Definition 6. Feasible Region.** The AUC and the KL divergence pair is lying in the feasible region (figure 3)
 257 for all possible detectors (ROC curves), i.e. no detector with the AUC and the KL divergence pair lie outside the
 258 feasible region⁷.

Theorem 3. Possible feasible region for the AUC and the KL divergence. Given the ROC curve, the parametric possible feasible region as shown in figure 3 can be expressed using the positive parameter $a > 0$ as

$$\Pr(L_{\Delta} > 0) = \frac{1}{1 - e^{-a}} - \frac{1}{a}$$

and

$$\mathcal{D}_l^* \geq \log(a) + \frac{a}{e^a - 1} - 1 - \log(1 - e^{-a})$$

where

$$\mathcal{D}_l^* = \min \{ \mathcal{D}(f_{L_1}(l)||f_{L_0}(l)), \mathcal{D}(f_{L_0}(l)||f_{L_1}(l)) \}.$$

259 **Proof.** Proof is given in the appendix B. ■

⁷ The definition of the feasible region here is inspired by the joint range of f- divergences [33].

260 Theorem 3 formulates the relationship between the AUC and the KL divergence. *The results of this*
 261 *theorem is generally true for any LLRT statistic.* Theorem 3 states that for any valid ROC corresponds to
 262 a detector, the pair of AUC and KL divergence *must* lie in the possible feasible region (figure 3), i.e.
 263 outside of this region is infeasible. This possible feasible region results in the general upper bound for
 264 AUC.

265 Since computing the distribution of the LLRT statistics is not straight forward in most cases,
 266 proposition 3, relaxes the Theorem 3 by bounding the KL divergence between the LLRT statistics using
 267 the the invariance property of KL divergence for the LLRT statistic (lemma 2).

Proposition 3. *The parametric upper bound for AUC is*

$$\Pr(L_{\Delta} > 0) = \frac{1}{1 - e^{-a}} - \frac{1}{a}$$

and

$$\mathcal{D}^* \geq \log(a) + \frac{a}{e^a - 1} - 1 - \log(1 - e^{-a})$$

where $a > 0$ is a positive parameter and

$$\mathcal{D}^* = \min\{ \mathcal{D}(f_{\underline{X}}(\underline{x}|\mathcal{H}_1)||f_{\underline{X}}(\underline{x}|\mathcal{H}_0)), \mathcal{D}(f_{\underline{X}}(\underline{x}|\mathcal{H}_0)||f_{\underline{X}}(\underline{x}|\mathcal{H}_1)) \}. \quad (9)$$

Proof. Proof is based on the lemma 2 and the possible feasible region presented in the theorem 3. From the lemma 2, we have

$$\mathcal{D}_l^* \leq \mathcal{D}^*.$$

268 Then, using the result in the theorem 3, we get the parametric upper bound. ■

269 4.3. Asymptotic behavior for AUC bounds

Proposition 4. Asymptotic behavior of the lower bound. *We have*

$$\Pr(L_{\Delta} > 0) \geq 1 - e^{-n(1 - \frac{1}{n} \sum_{i=1}^n (1 + \frac{\kappa_i}{8})^{-1})}.$$

Proof. Applying the inequality

$$\frac{2x}{2+x} < \log(1+x)$$

270 for $x > 0$, we achieve the result. ■

Proposition 5. Asymptotic behavior of the upper bound. *The parametric upper bound for AUC has the following asymptotic behavior*

$$\Pr(L_{\Delta} > 0) \leq 1 - e^{-\mathcal{D}^* - 1}$$

271 where \mathcal{D}^* is given in (9).

Proof. Proof is as follows.

$$\begin{aligned} -\log(1 - \Pr(L_{\Delta} > 0)) &= -\log\left(\frac{1}{e^a - 1} + \frac{1}{a}\right) \\ &\leq \log(a) \\ &\leq \mathcal{D}^* + 1. \end{aligned}$$

272 Applying the exponential function to both sides of the above inequality we conclude the upper bound.
 273 ■

274 Figure 4 shows the possible feasible region and the asymptotic behavior log-scale. As it is shown
 275 in this figure, the parametric upper bound can be approximated with a straight line especially for large
 276 values of the parameter a (the result in proposition 5). Also, figure 5 shows the possible feasible region
 277 and the asymptotic behavior in regular-scale.

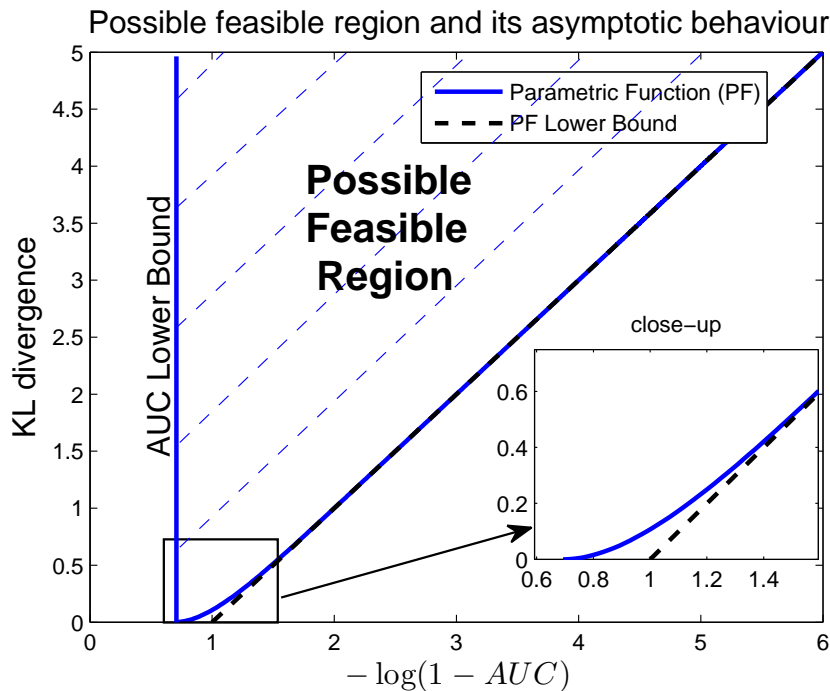


Figure 4. Log-scale of the possible feasible region and its asymptotic behavior (linear line) for the AUC and the KL divergence pair for all possible detectors or equivalently all possible ROC curves (the KL divergence is between the LLRT statistics under different hypotheses, i.e. $\mathcal{D}(f_{L_1}(l)||f_{L_0}(l))$ or $\mathcal{D}(f_{L_0}(l)||f_{L_1}(l))$.) Close-up part shows the non-linear behavior of the possible feasible region around one.

278 5. Examples and Simulation Results

279 In this section, we consider some examples of covariance matrices for Gaussian random vector \underline{X} .
 280 We pick the tree structure as the graphical model corresponds to the covariance selection problem. In
 281 our simulations, we compare the numerically evaluated AUC and its lower and upper bounds and
 282 discuss their asymptotic behavior as the dimension of the graphical model, n , increases.

283 5.1. Tree approximation model

The maximum order of the lower order distributions in tree approximation problem is two, i.e. no more than pairs of variables. Let $\underline{X}_{\mathcal{T}} \sim \mathcal{N}(0, \Sigma_{\underline{X}_{\mathcal{T}}})$ have the graph representation $\mathcal{G}_{\mathcal{T}} = (\mathcal{V}, \mathcal{E}_{\mathcal{T}})$ where $\mathcal{E}_{\mathcal{T}} \subseteq \psi$ is a set of edges that represents a tree structure. Let $\underline{X}_r \sim \mathcal{N}(0, \Sigma_{\underline{X}_r})$ have the graph representation $\mathcal{G}_r = (\mathcal{V}, \mathcal{E}_r)$ where $\mathcal{E}_r \subseteq \mathcal{E}_{\mathcal{T}}$ is the set of all edges in the graph of \underline{X}_r . The joint PDF for elements of random vector \underline{X}_r can be represented by joint PDFs of two variables and marginal PDFs in the following convenient form

$$f_{\underline{X}_r}(\underline{x}_r) = \prod_{(u,v) \in \mathcal{E}_r} \frac{f_{\underline{X}^u, \underline{X}^v}(\underline{x}^u, \underline{x}^v)}{f_{\underline{X}^u}(\underline{x}^u) f_{\underline{X}^v}(\underline{x}^v)} \prod_{u \in \mathcal{V}} f_{\underline{X}^u}(\underline{x}^u). \quad (10)$$

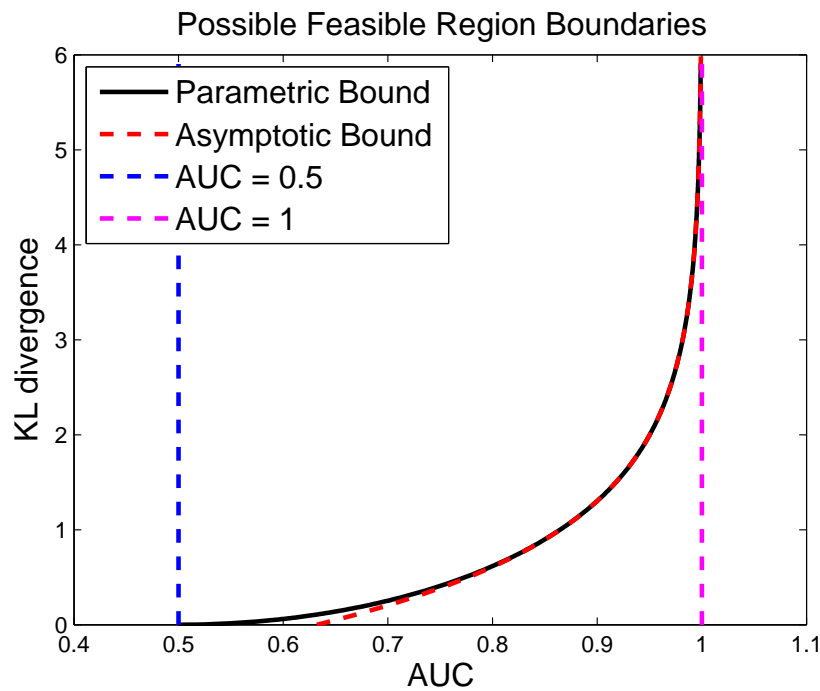


Figure 5. The possible feasible region boundaries and its asymptotic behavior for the AUC and the KL divergence pair for all possible detectors or equivalently all possible ROC curves (the KL divergence is between the LLRT statistics under different hypotheses, i.e. $\mathcal{D}(f_{L_0}(l)||f_{L_1}(l))$ or $\mathcal{D}(f_{L_1}(l)||f_{L_0}(l))$).

Using equation (10) we can then easily construct a tree using iterative algorithms (such as the Chow-Liu algorithm [2] combined with the Kruskal [12] algorithm or the Prim [13] algorithm) by adding edges one at a time [34]. Consider the sequence of random vectors \underline{X}_r with $0 \leq r \leq |\mathcal{E}_{\mathcal{T}}|$, where \underline{X}_r is recursively generated by augmenting a new edge, $(i, j) \in \mathcal{E}_r$, to the graph representation of \underline{X}_{r-1} . For the special case of Gaussian distributions, $\Sigma_{\underline{X}_r}$ has the following recursive formulation [34]

$$\Sigma_{\underline{X}_r}^{-1} = \Sigma_{\underline{X}_{r-1}}^{-1} + \Sigma_{i,j}^{\dagger} - \Sigma_i^{\dagger} - \Sigma_j^{\dagger}, \quad \forall 0 \leq r \leq |\mathcal{E}_{\mathcal{T}}|$$

284 where $\Sigma_{i,j}^{\dagger} = [e_i \ e_j] \Sigma_{i,j}^{-1} [e_i \ e_j]^T$ and $\Sigma_i^{\dagger} = e_i \Sigma_i^{-1} e_i^T$ where e_i is a unitary vector with 1 at the i -th place
 285 and $\Sigma_{i,j}$ and Σ_i are the 2-by-2 and 1-by-1 principle sub-matrices of $\Sigma_{\underline{X}_r}$, with initial step $\Sigma_{\underline{X}_0} = \text{diag}(\Sigma_{\underline{X}})$
 286 where $\text{diag}(\Sigma_{\underline{X}})$ represents a diagonal matrix with diagonal elements of $\Sigma_{\underline{X}}$.

287 **Remark:** For all $0 \leq r \leq |\mathcal{E}_{\mathcal{T}}|$, we have

- 288 1. $\text{tr}(\Sigma_{\underline{X}_r}) = \text{tr}(\Sigma_{\underline{X}})$
- 289 2. $\text{tr}(\Sigma_{\underline{X}} \Sigma_{\underline{X}_r}^{-1}) = n$.
- 290 3. $\mathcal{D}(f_{\underline{X}}(\underline{x})||f_{\underline{X}_r}(\underline{x})) = -\frac{1}{2} \log(|\Sigma_{\underline{X}} \Sigma_{\underline{X}_r}^{-1}|)$
- 291 4. $|\Sigma_{\underline{X}}| \leq \dots \leq |\Sigma_{\underline{X}_r}| \leq \dots \leq |\Sigma_{\underline{X}_0}| = |\text{diag}(\Sigma_{\underline{X}})|$
- 292 5. $H(\underline{X}) \leq \dots \leq H(\underline{X}_r) \leq \dots \leq H(\underline{X}_0)$.

293 Tree approximation models are interesting to study since there are algorithms such as Chow-Liu
 294 [2] combined by the Kruskal [12] or the Prim's [13] that efficiently compute the model covariance
 295 matrix from the graph covariance matrix.

296 5.2. Toeplitz example

Here, we assume that the covariance matrix $\Sigma_{\underline{X}}$ has a Toeplitz structure with ones on the diagonal elements and the correlation coefficient $\rho > -\frac{1}{(n-1)}$ as off diagonal elements

$$\Sigma_{\underline{X}} = \begin{bmatrix} 1 & \rho & \dots & \rho \\ \rho & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \rho \\ \rho & \dots & \rho & 1 \end{bmatrix}.$$

297 For the tree structure model, all possible tree structured distributions satisfying (10) have the same
 298 KL divergence to the original graph, i.e. $\mathcal{D}(f_{\underline{X}}(x)||f_{\underline{X}_{\mathcal{T}}}(x))$ is constant for all possible connected tree
 299 approximation model for this example. The reason is that all the weights computed by the Chow-Liu
 300 algorithm to construct the weighted graph associated with this problem are the same and are equal to
 301 $-\frac{1}{2}\log(1+\rho^2)$, which only depends on the correlation coefficient ρ . In the sequel, we test our results
 302 for two tree structured networks: a star network and a chain network.

303 5.2.1. Star approximation

The star covariance matrix is as follows (all the nodes are connected to the first node)⁸

$$\Sigma_{\underline{X}_{\mathcal{T}}}^{star} = \begin{bmatrix} 1 & \rho & \dots & \dots & \rho \\ \rho & \ddots & \rho^2 & \dots & \rho^2 \\ \vdots & \rho^2 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \rho^2 \\ \rho & \rho^2 & \dots & \rho^2 & 1 \end{bmatrix}.$$

For this example, the KL divergence and the Jeffreys divergence can be computed in closed form as

$$\mathcal{D}(\underline{X}||\underline{X}_{star}) = \frac{1}{2}(n-1)\log(1+\rho) - \frac{1}{2}\log(1+(n-1)\rho)$$

and

$$\mathcal{D}_{\mathcal{J}}(\underline{X}, \underline{X}_{star}) = \frac{(n-1)(n-2)\rho^2}{2(1+(n-1)\rho)}$$

respectively, where

$$\mathcal{D}_{\mathcal{J}}(\underline{X}, \underline{X}_{star}) = \mathcal{D}(\underline{X}||\underline{X}_{star}) + \mathcal{D}(\underline{X}_{star}||\underline{X})$$

is the Jeffreys divergence [17]. Moreover, for large values of n we have that

$$\mathcal{D}(\underline{X}||\underline{X}_{star}) \approx \frac{n}{2}\log(1+\rho)$$

and

$$\mathcal{D}_{\mathcal{J}}(\underline{X}, \underline{X}_{star}) \approx \frac{n}{2}\rho.$$

304 Figure 6 plots the 1-AUC v.s. the dimension of the graph, n for different correlation coefficients,
 305 $\rho = 0.1$ and $\rho = 0.9$. This figure also indicates the upper bound and the lower bound for the 1-AUC.

⁸ All n possible star networks have the same performance.

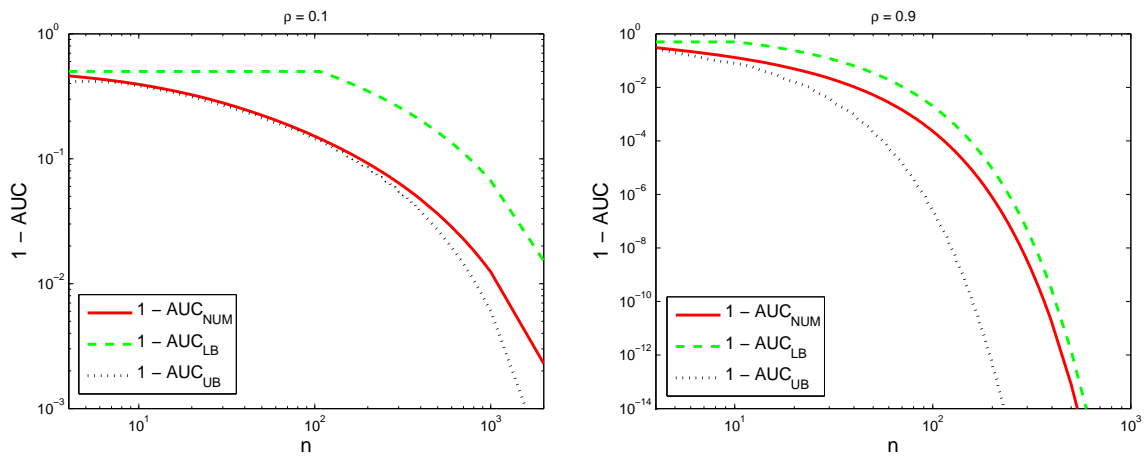


Figure 6. $1 - \text{AUC}$ v.s. the dimension of the graph, n for Star approximation of the Toeplitz example with $\rho = 0.1$ (left) and $\rho = 0.9$ (right). In both figures, the numerically evaluated AUC is compared with its bounds.

306 5.2.2. Chain approximation

307 The chain covariance matrix is as follows (nodes are connected like a first order Markov chain, 1
308 to n)

$$\Sigma_{\underline{X}_{\mathcal{T}}}^{\text{chain}} = \begin{bmatrix} 1 & \rho & \rho^2 & \dots & \rho^{n-1} \\ \rho & \ddots & \ddots & \ddots & \vdots \\ \rho^2 & \ddots & \ddots & \ddots & \rho^2 \\ \vdots & \ddots & \ddots & \ddots & \rho \\ \rho^{n-1} & \dots & \rho^2 & \rho & 1 \end{bmatrix}.$$

For this example, the KL divergence and the Jeffreys divergence can be computed in closed form as

$$\mathcal{D}(\underline{X} || \underline{X}_{\text{chain}}) = \mathcal{D}(\underline{X} || \underline{X}_{\text{star}})$$

and

$$\mathcal{D}_{\mathcal{J}}(\underline{X}, \underline{X}_{\text{chain}}) = \frac{\rho^2}{(1 + (n-1)\rho)(1-\rho)} \times \left(\frac{n(n-1)}{2} - \frac{n(1-\rho^n)}{1-\rho} + \frac{1 - (n+1)\rho^n + n\rho^{n+1}}{(1-\rho)^2} \right)$$

respectively. Moreover, for large values of n we have the following approximation

$$\mathcal{D}_{\mathcal{J}}(\underline{X}, \underline{X}_{\text{chain}}) \approx \frac{n}{2} \frac{\rho}{1-\rho}.$$

309 Figure 7 plots the $1 - \text{AUC}$ v.s. the dimension of the graph, n for different correlation coefficients,
310 $\rho = 0.1$ and $\rho = 0.9$ as well as its upper and lower bounds.

311 In both figure 6 and figure 7, $(1 - \text{AUC})$ and its bounds rapidly goes to 0 which means that AUC
312 goes to one as we increase the number of nodes, n , in the graph. More precisely, bounds for $1 - \text{AUC}$
313 are decaying exponentially as the dimension of the graph, n , increases which is consistent with the
314 theory obtained for analytical bounds. Furthermore, we can conclude from these figures that a smaller
315 ρ results in a better tree approximation, i.e. covariance matrices with smaller correlation coefficients

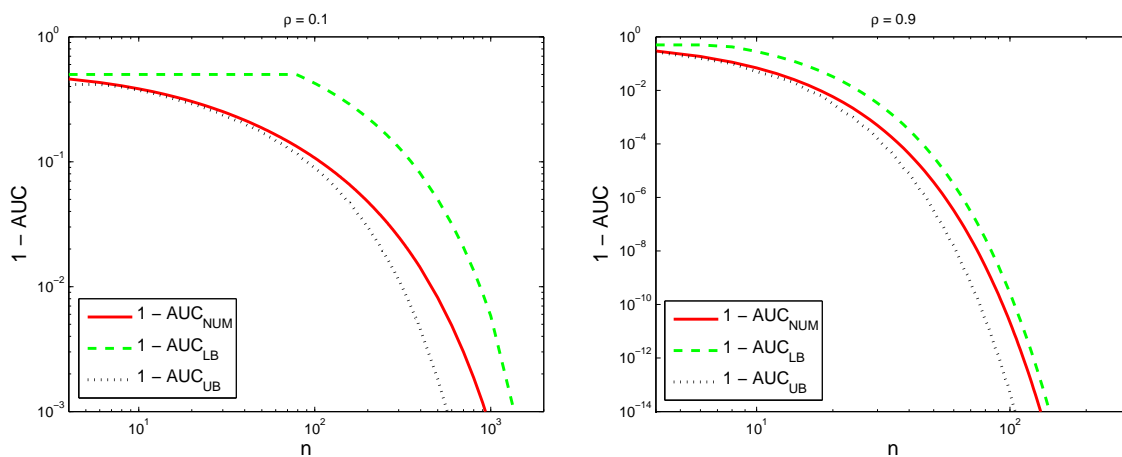


Figure 7. $1 - \text{AUC}$ v.s. the dimension of the graph, n for Chain approximation of the Toeplitz example with $\rho = 0.1$ (left) and $\rho = 0.9$ (right). In both figures, the numerically evaluated AUC is compared with its bounds.

316 are more like tree structure model. Moreover, comparing the AUC for the star network approximation
 317 with the AUC for the chain network approximation we conclude that the star network is a much better
 318 approximation than the chain network even though that both approximation networks have the same
 319 KL divergences. We can also interpret this fact through the analytical bounds obtained in this paper.
 320 The star network is a better approximation than the chain network since the decay rate of $1 - \text{AUC}$ for
 321 the star network is less than its decay rate for the chain network.

322 **Remark:** The star approximation in the above example has lower AUC than the chain approximation.
 323 Practically, it means the correlation between nodes that are not connected in the approximated graphical
 324 structure is more realistic in star network than the chain network.

325 5.3. Solar data

326 In this Example, covariance matrix is calculated based on datasets presented in [35]. Two datasets
 327 which are obtained from the National Renewable Energy Laboratory (NREL) website [36]. The first
 328 data set is the Oahu solar measurement grid which consists of 19 sensors (17 horizontal sensors and
 329 two tilted sensors) and the second one is the NREL solar data for 6 sites near Denver, Colorado. These
 330 two data sets are normalized using standard normalization method and the zenith angle normalization
 331 method [35] and then the unbiased estimate of the correlation matrix is computed⁹.

332 5.3.1. The Oahu solar measurement grid dataset

333 From data obtained from 19 solar sensors at the island of Oahu, we computed the spatial
 334 covariance matrix during the summer season at 12:00 PM averaged over a window of 5 minutes.
 335 Then, the AUC and the KL divergence are computed for those tree structures that are generated using
 336 Markov Chain Monte-Carlo (MCMC) method. Figure 8 shows the distribution of those tree structures
 337 generated using MCMC method versus the KL divergence (left) and v.s. $\log_{10}(1 - \text{AUC})$ (right)¹⁰.

338 Looking back at figure 4, for very small value of $1 - \text{AUC}$ the relationship between the KL
 339 divergence and the boundary of the possible feasible region for $-\log(1 - \text{AUC})$ is linear. This
 340 means that if the upper bound is tight then the relationship between the KL divergence and the

⁹ See [35] for fields definition and other details about the normalization methods for the solar irradiation covariance matrix.

¹⁰ In this example, since the AUC for all generated tree structures is close to one, we plots the distribution of generated trees v.s. $\log_{10}(1 - \text{AUC})$.

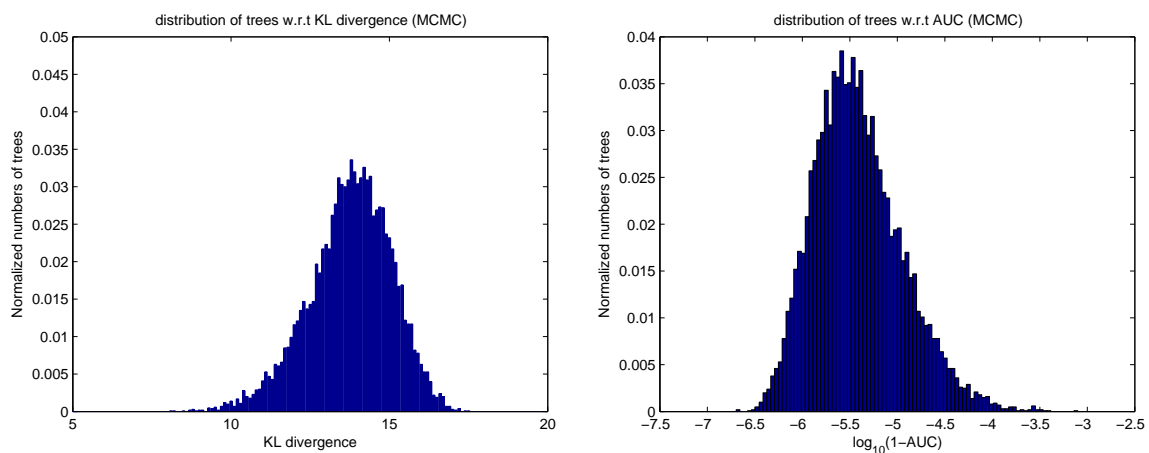


Figure 8. Left: distribution of the generated trees (Normalized histogram) using MCMC v.s. the KL divergence and **Right:** distribution of the generated trees (Normalized histogram) using MCMC v.s. $\log_{10}(1 - \text{AUC})$ for the Oahu solar measurement grid dataset in summer season at 12:00 PM.

341 $-\log(1 - \text{AUC})$ is almost linear. In figure 8, the maximum value of $1 - \text{AUC}$ for this model is less than
 342 10^{-3} which justifies why two distributions in figure 8 are scaled/mirrored of each other. Moreover,
 343 just by looking at the distribution of tree models in this example, it is obvious that most tree models
 344 have similar performance. Only a small portion of the tree models have better performance than the
 345 most trees, but the difference is not that significant.

346 5.3.2. The Colorado dataset

347 From the solar data obtained from 6 sensors near Denver, Colorado, we computed the spatial
 348 covariance matrix during the summer season at 12:00 PM averaged over a window of 5 minute. Then,
 349 the AUC and the KL divergence are computed for all possible tree structures. Figure 9 shows the
 350 distribution of all possible tree structures v.s the KL divergence (**left**) and v.s the AUC (**right**).

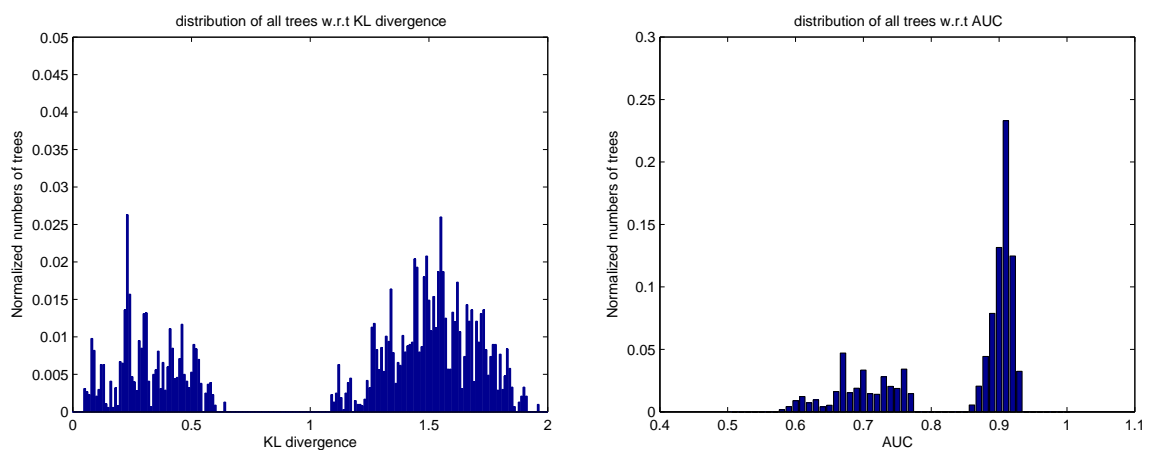


Figure 9. Left: distribution of all trees (Normalized histogram) v.s. the KL divergence and **Right:** distribution of all trees (Normalized histogram) v.s. the AUC for the Colorado dataset in summer season at 12:00 PM.

351 In the Colorado dataset, there are two sensors that are very close to each other compared to the
 352 distance between all other pairs of sensors. As a result, if the particular edge between these two sensors
 353 is in the approximated tree structure we get a smaller AUC and KL divergence compared to when that

354 particular edge is not in the tree structure. This explains why the distributions of all trees in this case
 355 looks like a mixture of two distributions. This result also gives us valuable incite on how to answer the
 356 following question, "How to construct informative approximation algorithms for model selection in
 357 general." One catch as an example is that for the Colorado dataset, almost all trees that contain the
 358 particular edge between the two aforementioned sensors are good approximations while the rest of
 359 tree models' performances are not desirable.

360 5.3.3. Two-dimensional sensor network

In this example, we create a 2-dimensional (2D) sensor network using Gaussian kernel [37] as follows

$$\Sigma_{\underline{X}}(i, j) = \left[e^{-\frac{d(i, j)^2}{2\sigma^2}} \right]$$

361 where $d(i, j)$ is the Euclidean distance between the i -th sensor and the j -th sensor in the 2D space. All
 362 sensors are located randomly in 2D space¹¹. We set $\sigma = 1$ and generate a 2D sensor network with 20
 363 sensors. For the 2D sensor network example, figure 10 shows the distribution of the generated tree
 364 structures using MCMC method v.s KL divergence (**left**) and v.s $\log_{10}(1 - \text{AUC})$ (**right**). Again we see
 365 the mirroring effect in Fig. 10 as we have an almost linear relationship between the KL divergence and
 366 $-\log(1 - \text{AUC})$. Note that, the covariance matrix generated has one dominant eigenvalue in most
 367 cases. Furthermore, figure 11 plots $1 - \text{AUC}$ as well as its analytical upper bound and lower bound
 368 v.s. the dimension of the graph, n for $\sigma = 1.3$ (**left**) and $\sigma = 1.8$ (**right**). To generate this figure, we
 369 randomly generated 1000 sensor networks and then plot the averaged AUC. As we can see in this
 370 figure, the $1 - \text{AUC}$ and its bounds decay exponentially which is consistent with the theoretical results
 of this paper.

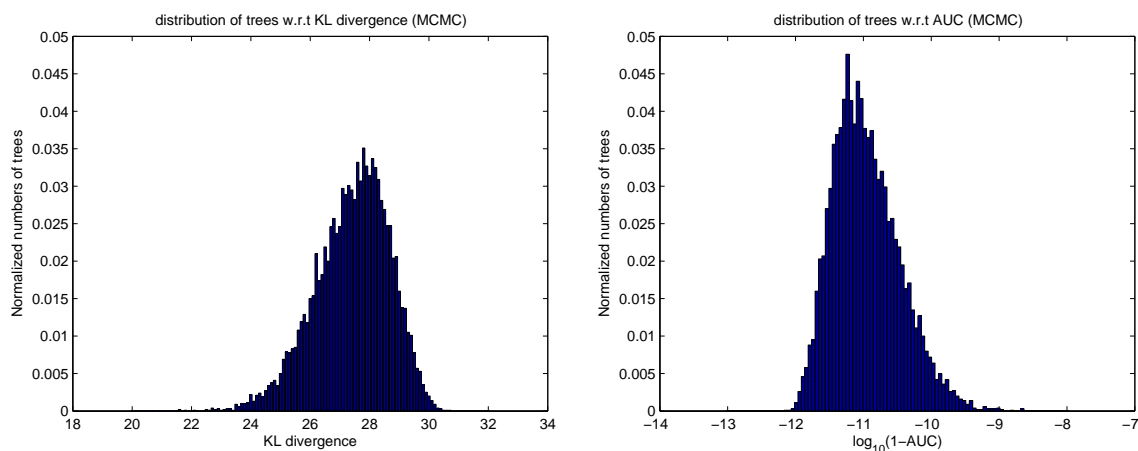


Figure 10. Left: distribution of the generated trees (Normalized histogram) using MCMC v.s. the KL divergence and **Right:** distribution of the generated trees (Normalized histogram) using MCMC v.s. $\log_{10}(1 - \text{AUC})$ for the 2D sensor network example with 20 sensors and $\sigma = 1$.

371

372 6. conclusion

373 In this paper, we formulate a detection problem and investigate the quality of model selection
 374 problem. More specifically, we consider Gaussian distributions and discuss the covariance selection
 375 quality of a given model. We present the correlation approximation matrix (CAM), and show its
 376 relationship with information theory divergences such as the KL divergence, the reverse KL divergence

¹¹ Sensors location in each dimension are drawn randomly from a Normal distribution.

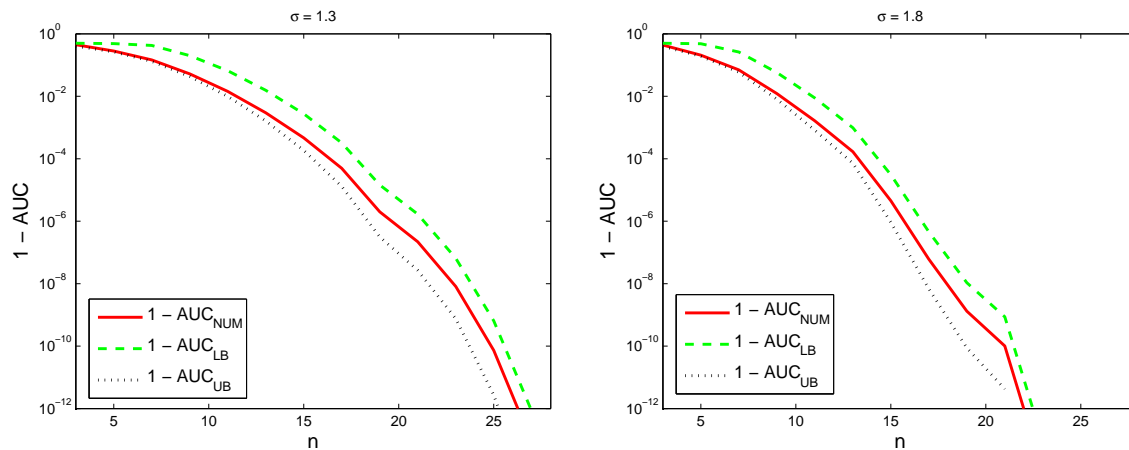


Figure 11. $1 - \text{AUC}$ and its bounds v.s. the dimension of the graph, n for $\sigma = 1.3$ (left) and $\sigma = 1.8$ (right), averaged over 1000 runs of sensor networks generated randomly.

377 and the Jeffreys divergence as well as the ROC curve and the area under it, i.e. the AUC, as a
 378 measure of accuracy in the detection problem framework. Moreover, this paper presents an analytical
 379 expression for the AUC that can efficiently be evaluated numerically. Also, the AUC analytical lower
 380 and upper bounds are provided in this paper. We show that the AUC and the lower bound for the
 381 AUC depend on the eigenvalues of the CAM. Upper bounds for the AUC are obtained from finding
 382 a parametric relationship between the AUC and the KL/reverse KL divergences. We pick the tree
 383 structure as an example of an approximation model and use the Chow-Liu MST algorithm to compute
 384 the maximum likelihood tree structure approximation. Then, the quality of the Chow-Liu MST tree
 385 algorithm is investigated using the formulated detection problem. Through some examples, we show
 386 that in general, the tree approximation is not a good model as the number of nodes in the graphical
 387 model increases which is the case in high dimensional problems such as modeling the electrical
 388 distribution grid using smart grid sensor measurements and distributed renewable energy sources.
 389 The aforementioned result is also consistent with the analytical results provided in this paper that is
 390 $1 - \text{AUC}$ decays exponentially as the dimension of graph increases.

391 The detection framework presented in this paper, can be generalized for non-Gaussian models.
 392 Moreover, the AUC analytical bounds obtained in this paper can also be used in other applications
 393 that are using AUC as a relevant criterion. One example is in medicine when the AUC is used for
 394 diagnostic tests between positive instance and negative instance [38] where instead of changing the
 395 coordinates we can look at the exponent of the AUC bounds. In ongoing work we are looking at more
 396 accurate graphical approximations that involve non-tree graphs. These approximations use a variation
 397 of the CAM which we call the symmetric CAM and simple linear transformations.

398 **Acknowledgments:** This paper was presented for the special case of tree approximation in part at 2016 Information
 399 Theory and Application Workshop [25]. Authors would like to thank Prof. Peter Harremoës for his helpful
 400 discussions on information divergences and assistance with Theorem 3. This work was supported in part by NSF
 401 grant ECCS-1310634, the Center for Science of Information (CSoI), an NSF Science and Technology Center, under
 402 grant agreement CCF-0939370, and the University of Hawaii REIS project.

403 **Conflicts of Interest:** The authors declare no conflict of interest. The founding sponsors had no role in the design
 404 of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, and in the
 405 decision to publish the results.

406 Appendix A Proof of Lemma 1

The calculus based proof for the special case of continuous PDFs is as follow. We can apply the Leibniz integral rule [39] and compute the derivative of CDFs $P_0(l)$ and $P_1(l)$ as

$$f_{L_0}(l) = -\frac{dP_0(l)}{dl}$$

and

$$f_{L_1}(l) = -\frac{dP_1(l)}{dl}$$

since $f_{L_0}(l)$ and $f_{L_1}(l)$ are continuous functions.¹² We have

$$\begin{aligned} \mathcal{D}(f_{L_0}(l)||f_{L_1}(l)) &= \int_{-\infty}^{+\infty} \log \frac{f_{L_0}(l)}{f_{L_1}(l)} f_{L_0}(l) dl \\ &\stackrel{(a)}{=} - \int_0^1 \log \frac{dP_1}{dP_0} dP_0 \\ &\stackrel{(b)}{=} - \int_0^1 \log h'(z) dz \end{aligned}$$

407 where equality (a) is true since we can replace PDFs $f_{L_0}(l)$ and $f_{L_1}(l)$ using the derivative of their
408 CDFs. Equality (b) is just a change of variable, $z = P_0(l)$, in order to write the integral in terms of the
409 derivative of the ROC curve. Proof for the second part of this lemma is similar to the proof of the first
410 part. ■

411 Appendix B Proof of Theorem 3

412 Looking back at properties of the ROC curve, $h(z)$, where $z \in [0, 1]$, the ROC curve have to satisfy
413 the following conditions

- 414 • **C1:** $\int_0^1 h'(z) dz = 1$
- 415 • **C2:** $h'(z) \geq 0$
- 416 • **C3:** $h'(z)$ is decreasing

where $h'(z)$ is the derivative of the ROC curve, $h(z)$. Also for a given ROC curve, $h(z)$, we can compute the AUC as

$$\Pr(L_\Delta > 0) = \int_0^1 h(z) dz.$$

Then, using integration by parts, we can show that

$$1 - \Pr(L_\Delta > 0) = \int_0^1 z h'(z) dz.$$

417 To compute the possible feasible region stated in the theorem 3, we need to optimize both of
418 following KL divergences, $\mathcal{D}(f_{L_1}(l)||f_{L_0}(l))$ and $\mathcal{D}(f_{L_0}(l)||f_{L_1}(l))$, with respect to the derivative of
419 the ROC curve given a fixed AUC, $\Pr(L_\Delta > 0)$, while conditions, C1, C2 and C3 hold. To solve this
420 optimization, we can use the method of Lagrange multiplier.

¹² Both $f_{L_0}(l)$ and $f_{L_1}(l)$ are PDFs in generalized Chi-squared distributions class. This means that each of these PDFs are convolution of weighted Chi-squared distributions. Weighted Chi-squared distribution is continuous in its domain thus, convolution of these distributions is continuous in its domain.

First step: Here we minimize $\mathcal{D}(f_{L_1}(l)||f_{L_0}(l))$ with respect to the derivative of the ROC curve given the constraints. Optimization problem is as follow

$$\begin{aligned} \arg \min_{h'(z)} \quad & - \int_0^1 \log h'(z) dz & (A1) \\ \text{s. t.} \quad & \int_0^1 z h'(z) dz = 1 - \Pr(L_\Delta > 0) \\ & C1, C2 \ \& \ C3. \end{aligned}$$

To solve this optimization problem, we first write the Lagrangian. We need two coefficients a and b corresponding to conditions in optimization problem (A1). Then, we can write the Lagrange multiplier as a function of the derivative of the ROC curve, z , a and b as follow

$$\begin{aligned} L(h'(z), z, a, b) = & - \int_0^1 \log h'(z) dz \\ & + a \left(\int_0^1 z h'(z) dz - (1 - \Pr(L_\Delta > 0)) \right) \\ & + b \left(\int_0^1 h'(z) dz - 1 \right). \end{aligned}$$

Note that, the Lagrangian, $L(h'(z), z, a, b)$ is a convex function of $h'(z)$. Thus, we can compute its minimum by taking its derivative with respect to $h'(z)$. Doing so, we get

$$\frac{\partial L(h'(z), z, a, b)}{\partial h'(z)} = \int_0^1 \left(az + b - \frac{1}{h'(z)} \right) dz.$$

Set $\frac{\partial L(h'(z), z, a, b)}{\partial h'(z)} = 0$ we get

$$h'(z) = \frac{1}{az + b}$$

421 for all $z \in [0, 1]$. From C3, since $h'(z)$ is decreasing, we can conclude that $a > 0$. Moreover, from C1, at
422 optimum we have $\int_0^1 h'(z) dz = 1$ and thus, we can compute one of the coefficients as $b = \frac{a}{e^a - 1}$.

Computing the AUC integral and the KL divergence using the ROC curve we get the following parametric boundary for the possible feasible region

$$\Pr(L_\Delta > 0) = \frac{1}{1 - e^{-a}} - \frac{1}{a} \quad (A2)$$

and

$$\mathcal{D} = \log(a) + \frac{a}{e^a - 1} - 1 - \log(1 - e^{-a}) \quad (A3)$$

423 where $\mathcal{D} = \mathcal{D}(f_{L_1}(l)||f_{L_0}(l))$.

Second step: Here we minimize $\mathcal{D}(f_{L_0}(l)||f_{L_1}(l))$. The Lagrange multiplier for this step is similar to the first step but it is more straight forward if we define $g(\eta) = h^{-1}(\eta)$. Note that using integration by parts, we can show that AUC is

$$\Pr(L_\Delta > 0) = \int_0^1 \eta g'(\eta) d\eta.$$

Now, we can write the Lagrangian for the optimization problem with respect to $g'(\eta)$. The Lagrangian is convex with respect to $g'(\eta)$, thus taking the derivative and set it equal to zero as follow

$$\frac{\partial L(g'(\eta), \eta, a, b)}{\partial g'(\eta)} = 0$$

we can compute the parametric boundary for the possible feasible region. The parametric boundary in this case is the same as solution in (A2) and (A3) with $\mathcal{D} = \mathcal{D}(f_{L_0}(l)||f_{L_1}(l))$. Thus, combining these two steps, for the optimal boundary we have

$$\mathcal{D}_l^* = \min\{\mathcal{D}(f_{L_1}(l)||f_{L_0}(l)), \mathcal{D}(f_{L_0}(l)||f_{L_1}(l))\}.$$

424

425 **References**

- 426 1. Dempster, A.P. Covariance selection. *Biometrics* **1972**, *28*, 157–175.
- 427 2. Chow, C.K.; Liu, C.N. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory* **1968**, pp. 462–467.
- 428 3. Shuman, D.I.; Narang, S.K.; Frossard, P.; Ortega, A.; Vandergheynst, P. The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *Signal Processing Magazine, IEEE* **2013**, *30*, 83–98.
- 429 4. Koller, D.; Friedman, N. *Probabilistic graphical models: principles and techniques*; MIT press, 2009.
- 430 5. Jordan, M.I. *Learning in graphical models*; Vol. 89, Springer Science & Business Media, 1998.
- 431 6. Lauritzen, S.L. *Graphical models*; Clarendon Press, 1996.
- 432 7. Kschischang, F.R.; Frey, B.J.; Loeliger, H.A. Factor graphs and the sum-product algorithm. *Information Theory, IEEE Transactions on* **2001**, *47*, 498–519.
- 433 8. Loeliger, H.A.; Dauwels, J.; Hu, J.; Korl, S.; Ping, L.; Kschischang, F.R. The Factor Graph Approach to Model-Based Signal Processing. *Proceedings of the IEEE* **2007**, *95*, 1295–1322.
- 434 9. Khajavi, N.T.; Kuh, A. First order Markov chain approximation of microgrid renewable generators covariance matrix. Proc. of IEEE International Symposium on Information Theory, Istanbul, Turkey (ISIT'13), 2013, pp. 1207–1211.
- 435 10. Meinshausen, N.; Buhlmann, P. Model selection through sparse maximum likelihood estimation. *Annals of Statistics* **2006**, pp. 1436–1464.
- 436 11. Friedman, J.; Hastie, T.; Tibshirani, R. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **2008**, *9*, 432–441.
- 437 12. Kruskal, J.B. On the shortest spanning subtree of a graph and the traveling salesman problem. *Proceedings of the American Mathematical society* **1956**, *7*, 48–50.
- 438 13. Prim, R.C. Shortest connection networks and some generalizations. *Bell system technical journal* **1957**, *36*, 1389–1401.
- 439 14. Khajavi, N.T. Latent tree approximation in linear model. Acoustics, Speech and Signal Processing (ICASSP), 2017 IEEE International Conference on. IEEE, 2017, pp. 5940–5944.
- 440 15. Kadane, J.B.; Lazar, N.A. Methods and criteria for model selection. *Journal of the American statistical Association* **2004**, *99*, 279–290.
- 441 16. MacKay, D.J.C. *Information theory, inference, and learning algorithms*; Vol. 7, Citeseer, 2003.
- 442 17. Jeffreys, H. An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London. Series A. Mathematical and Physical Sciences* **1946**, *186*, 453–461.
- 443 18. Lewis-II, P.M. Approximating probability distributions to reduce storage requirements. *Information and control* **1959**, *2*, 214–225.
- 444 19. Lehmann, E.L.; Romano, J.P. *Testing statistical hypotheses*; springer, 2006.
- 445 20. Neyman, J.; Pearson, E.S. On the use and interpretation of certain test criteria for purposes of statistical inference. *Biometrika* **1928**, *20*.
- 446 21. Eguchi, S.; Copas, J. Interpreting Kullback–Leibler divergence with the Neyman–Pearson lemma. *Journal of Multivariate Analysis* **2006**, *97*, 2034–2040.
- 447 22. Scharf, L.L. *Statistical signal processing*; Vol. 98, Addison-Wesley Reading, MA, 1991.
- 448 23. Shiryaev, A.N. Probability, volume 95 of Graduate texts in mathematics, 1996.
- 449 24. Hanley, J.A.; McNeil, B.J. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* **1982**, *143*, 29–36.

- 468 25. Khajavi, N.T.; Kuh, A. The Quality of Tree Approximation from AUC Bounds. *Information Theory and*
469 *Applications Workshop* **2016**.
- 470 26. Kotz, S.; Kozubowski, T.; Podgorski, K. *The Laplace distribution and generalizations: a revisit with applications*
471 *to communications, economics, engineering, and finance*; Springer Science & Business Media, 2012.
- 472 27. Abramowitz, M.; Stegun, A.I. Handbook of mathematical functions. *Applied mathematics series* **1966**, 55, 62.
- 473 28. Provost, S.B.; Rudiuk, E.M. The exact distribution of indefinite quadratic forms in noncentral normal
474 vectors. *Annals of the Institute of Statistical Mathematics* **1996**, 48, 381–394.
- 475 29. Ha, H.T.; Provost, S.B. AN ACCURATE APPROXIMATION TO THE DISTRIBUTION OF A LINEAR
476 COMBINATION OF NON-CENTRAL CHI-SQUARE RANDOM VARIABLES. *REVSTAT-Statistical Journal*
477 **2013**, 11, 231–254.
- 478 30. Al-Naffouri, T.Y.; Hassibi, B. On the distribution of indefinite quadratic forms in Gaussian random
479 variables. *Information Theory, 2009. ISIT 2009. IEEE International Symposium on. IEEE, 2009*, pp.
480 1744–1748.
- 481 31. Cover, T.M.; Thomas, J.A. *Elements of information theory*; John Wiley & Sons, 2012.
- 482 32. Kullback, S. *Information theory and statistics*; Courier Corporation, 1968.
- 483 33. Harremoës, P.; Vajda, I. On Pairs of f -divergences and their Joint Range. *arXiv preprint arXiv:1007.0097*
484 **2010**.
- 485 34. Kavcic, A.; Moura, J.M.F. Matrices with Banded Inverses: Inversion Algorithms and Factorization of
486 Gauss-Markov Processes. *IEEE Transactions on Information Theory* **2000**, 46, 1495–1509.
- 487 35. Khajavi, N.T.; Kuh, A.; Santhanam, N.P. Spatial Correlations for Solar PV Generation and its Tree
488 Approximation Analysis. *Proc. of the Asia-Pacific Signal and Information Processing Association (APSIPA*
489 *ASC)*, 2014, pp. 1–5.
- 490 36. Laboratory, N.R.E. Measurement and Instrumentation Data Center. [Online]. Available at
491 <http://www.nrel.gov/midc/>.
- 492 37. Rasmussen, C.E.; Williams, C.K.I. Gaussian processes for machine learning. *the MIT Press* **2006**.
- 493 38. Johnson, N.P. Advantages to transforming the receiver operating characteristic (ROC) curve into likelihood
494 ratio co-ordinates. *Statistics in medicine* **2004**, 23, 2257–2266.
- 495 39. Flanders, H. Differentiation under the integral sign. *The American Mathematical Monthly* **1973**, 80, 615–627.