

Feature ranking for multi-label classification using Markov Networks

Paweł Teisseyre

*Institute of Computer Science, Polish Academy of Sciences
Jana Kazimierza 5 01-248 Warsaw, Poland*

Abstract

We propose a simple and efficient method for ranking features in multi-label classification. The method produces a ranking of features showing their relevance in predicting labels, which in turn allows to choose a final subset of features. The procedure is based on Markov Networks and allows to model the dependencies between labels and features in a direct way. In the first step we build a simple network using only labels and then we test how much adding a single feature affects the initial network. More specifically, in the first step we use the Ising model whereas the second step is based on the score statistic, which allows to test a significance of added features very quickly. The proposed approach does not require transformation of label space, gives interpretable results and allows for attractive visualization of dependency structure. We give a theoretical justification of the procedure by discussing some theoretical properties of the Ising model and the score statistic. We also discuss feature ranking procedure based on fitting Ising model using l_1 regularized logistic regressions. Numerical experiments show that the proposed methods outperform the conventional approaches on the considered artificial and real datasets.

Keywords: feature selection, multi-label learning, Markov networks, Ising model

1. Introduction

Multi-label classification (MLC) has recently attracted a significant attention, motivated by an increasing number of applications. Examples include text categorization [1, 2, 3, 4, 5], image classification [6, 7, 8], video classification [9, 10], music categorization [11], gene and protein function prediction [12, 13, 14], medical diagnosis [15, 16], chemical analysis [17, 18], social network mining [19, 20] and direct marketing [21]. More examples can be found in [22], [23] and [24]. The key problem in multi-label learning is how to utilize label dependencies to improve the classification performance, motivated by which number of multi-label algorithms have been proposed in recent years (see [25] for extensive comparison of several methods). The recent progress in MLC is summarized in [26] and [22]. In MLC, each object of our interest (e.g. text, image, patient, etc.) is described by a vector of p features $\mathbf{x} = (x_1, \dots, x_p)^T$ and a vector of K binary labels $\mathbf{y} = (y_1, \dots, y_K)^T$. The main objective is to build a model (using some training examples) which predicts \mathbf{y} based on \mathbf{x} .

One of the trending challenges in MLC is a dimensionality reduction of the feature space [22], i.e. reducing the dimensionality of the vector \mathbf{x} . Usually only some features affect \mathbf{y} . The issue is very important as in practical applications, the dimensionality of feature space can be very large. For example in text categorization a standard approach is to use so-called *bag-of-words model* in which frequencies of occurrence of words in a corpora are taken as features. This method generates thousands of features. Moreover, one can also take into account higher degree n -grams (bigrams, trigrams, etc.) and many other types of features (e.g. stylistic features like averaged word length), which further increases the dimensionality of feature vector. Elimination of redundant features is essential for the following reasons. First, it allows to reduce the computational burden of MLC procedures. Secondly, it improves a prediction accuracy of MLC methods. Fitting many MLC models includes estimation of large number of parameters. It is well

Email address: teisseyrep@ipipan.waw.pl (Paweł Teisseyre)

URL: <http://www.ipipan.eu/~teisseyrep/> (Paweł Teisseyre)

known that fitting models with many spurious features increases the variance of estimators and thus decreases the prediction accuracy of the model (see e.g. chapter 7 in [27]). Finally, feature selection methods are used to discover dependency structure in data. This allows to understand how features affect the labels, which is particularly important in biological and medical applications. For example, in multi-morbidity (co-occurrence of two or more chronic medical conditions in one person) it is crucial to discover which characteristics of the patient influence the co-occurrence of diseases [28]. Moreover, it would be interesting to know which diseases are likely to occur simultaneously given some characteristics of the patient (for example age, gender and previous diseases). We discuss different approaches of dimensionality reduction in MLC in Section 2.

In this paper we focus on Feature Ranking (FR) methods (sometimes also called filters). Although the MLC attracted a significant attention in machine learning community, only a few works address the feature ranking problem in multi-label setting. Feature ranking (FR) methods are mainly used to assess the individual relevance of available features. More precisely, they allow to order features with respect to their relevance in predicting labels, which in turn allows to remove the least significant features and build a final classification model using the most significant features. Although usually in this approach neither the possible redundancy between features nor their joint relevancy is taken into account, the main advantage is a low computational cost, which allows to compute the importance of thousands of features relatively fast. This is crucial in many domains, like text categorization or functional genomics. Moreover, in some applications it is important to evaluate the individual relevance of features, not only their joint relevance. Some authors use FR methods as an initial step to filter out spurious features and then use more sophisticated selection methods on the remaining set of features (see e.g. Sure Independence Screening procedure proposed by [29]). We also discuss FR method, which incorporates all features simultaneously.

The FR task in multi-label setting is much more challenging than in a single-label case. In traditional classification with only one target variable, FR methods aim to model the dependence between target variable y and a single feature x_j using different variable importance measures. Then the procedure is repeated for all possible features. The most popular measures are: information gain ([30]), the chi-squared statistic and simple statistics based on univariate logistic regression ([31]), among others. On the other hand, in MLC feature x_j may affect targets y_1, \dots, y_K in different ways. First, it may happen that x_j influences only some of labels, while others are independent from x_j . More importantly, since in MLC methods dependencies between labels are usually considered, we should verify how x_j affects a given label y_k , in a presence of the remaining labels. It may happen that x_j is independent from y_k , while x_j becomes dependent on y_k , when conditioned on other labels. Finally, feature x_j can influence only the interactions between labels, while the marginal dependencies are not present. Examples of such situations are provided in Sections 3.1 and 3.3. A desirable FR method should take into account all the above aspects.

The main limitation of recent FR methods is that they require problem transformation methods: Binary Relevance (BR) or Label Powerset (LP) transformation for evaluating the relevance of given features. Unfortunately, both transformations suffer from many serious drawbacks, discussed in more detail in Section 2. To propose a desirable FR method, we make an effort to take into account the following aspects.

- The method should not use BR or LP transformation.
- The method should take into account specificity of multi-label setting, i.e. it should measure the dependence between feature x_j and label y_k , given the remaining labels.
- The method should give interpretable results to see which labels (or interactions between labels) and how are influenced by feature.
- The computational cost of the procedure should be low.

To take into account the above postulates, we propose a novel approach which is based on Markov Networks. Markov Network (see e.g. [32], Section 8) can be represented as a graph, with node set representing random variables (in our case labels and features) and edge set representing dependencies between variables. Existing edge between two variables means that they are conditionally dependent given the rest of the graph. The main advantage of Markov Networks is that they allow to model the pairwise dependencies between labels and features in a direct way. Although, Markov Networks have already been applied in MLC (see e.g. [33] or [34]), they have not been used as a feature ranking method. Our approach is based on the following idea. We initially build a Markov Network containing only labels, which allows to model the dependencies among the labels. In the second step, we test how much adding a

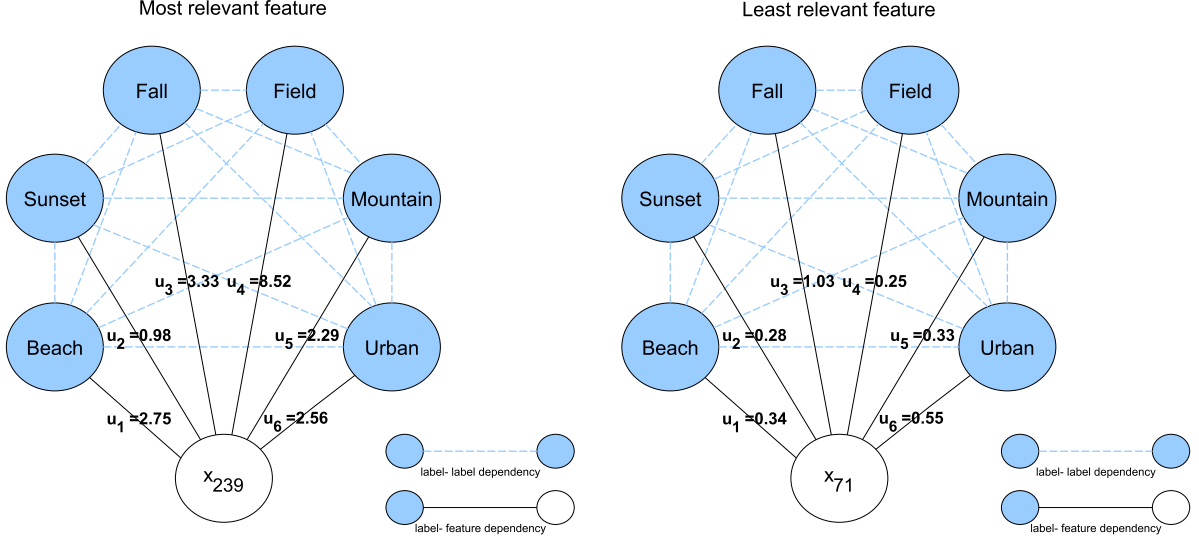


Figure 1: Markov networks corresponding to the most (x_{239}) and the least (x_{71}) significant features for scene dataset. The numbers over edges are scores statistics describing importances of features.

single feature x_j affects the initial network. This allows to test the dependence strength between a given feature x_j and a given label y_k , conditioning on the remaining labels. The procedure is repeated for all available features, which yields the final ranking. Specifically, in our method we use the Ising model ([35], [36]) which is a simple example of Markov Network. It turns out that for the Ising model, building an initial network containing only labels can be done relatively simply, especially for moderate number of labels. Please see Section 3.5 for deeper justification of using the Ising model. In a second step we propose to use the score statistic [37], which is very computationally efficient in this case. Namely, it is not necessary to refit an initial network when we add feature x_j . This allows to test a significance of added features very quickly which is crucial in FR methods. The details of the procedure are given in next Sections. Figure 1 shows networks corresponding to the most and the least significant features for scene dataset, in which the task is to predict six labels (beach, sunset, field, fall, mountain, and urban). Numbers over edges u_1, \dots, u_6 are the score statistics which reflect the conditional dependences between feature x_j and labels (given the remaining labels). The higher the value of the score statistic, the larger is the influence of x_j on the given label, in the presence of remaining labels. The score statistics for a given feature x_j are added together, which gives an importance measure for x_j . The final ranking is based on these importances. We also discuss FR procedure based on fitting Ising model using l_1 regularized logistic regressions.

The rest of the article is organized as follows. In Section 2 we discuss the existing related work. In Section 3 we present feature importance measures based on the Ising model. We describe two versions of the Ising model: the first one assumes constant interactions between labels, the second one considers feature-dependent interactions. We discuss some theoretical properties of the score statistic and justify using the Ising model. In addition we discuss a version of the Ising model which incorporates all features simultaneously and describe the estimation procedure based on l_1 regularized logistic regressions. Section 4 contains the formal description of our feature ranking procedures as well as feature selection procedure. We present the results of experiments in Section 5. Section 6 concludes with a summary. The technical proofs are provided in Appendix.

2. Related work

Dimensionality reduction of the feature space is one of the current challenges in MLC [22]. There are different approaches to reduce the dimensionality of feature space. The two main groups are: feature selection methods

(among which one can distinguish: feature ranking methods, wrappers, embedded methods) and feature transformation methods. Feature selection methods aim to identify a small subset of features which influence labels. Feature transformation methods aim to find functions of features, that can replace the original ones. In this paper we focus on feature ranking methods (FR), which belong to the first group. FR methods produce a list of features, ordered with respect to their relevance. The final model is built using the most relevant features from the list.

Let us first review the existing FR methods in MLC. The popular approach is to use Binary Relevance (BR) transformation (by considering classification tasks corresponding to separate labels) and to evaluate the relevance of each feature for each of the labels independently ([11], [38]). The scores corresponding to different labels are then combined, which yields the global ranking of features. To evaluate the relevance of features in the tasks, various feature importance measures are used, among which the chi-squared statistic and information gain are the most popular ones ([39]). The major drawback of this approach is that possible dependencies between labels are not utilized. The combinations of BR transformation with the chi-squared statistic and information gain will be referred to as *br chi2* and *br ig*, respectively.

The second popular group of methods is based on LP transformation ([23]) which reduces the multi-label problem to single-label problem with many classes by considering each combination of labels as a distinct meta-class. LP transformation combined with the chi-squared statistic has been used in music classification [11]. This approach requires discretization of features which may lead to loss of some information. [40] proposed to combine LP method and information gain (mutual information), whose estimation is in general a challenging task. They used Kozachenko-Leonenko estimator of entropy [41], which is based on nearest neighbours method. The approach was also successfully used to assess the relevance of subset of features, not only the individual significance of features, which is a big advantage. The limitation is that the presence of points having the exact same feature values may harm the estimation of entropy based on nearest neighbours. It turns out that the method based on information gain usually outperforms the chi-squared-based approach [40]. Feature selection based on information gain was also described in [42]. The combinations of LP transformation with the chi-squared statistic and information gain will be referred to as *lp chi2* and *lp ig*, respectively. Although, LP transformation is very simple, it suffers from many serious drawbacks. First of all, the number of possible meta-classes can be very large, even larger than the number of observations. As a result some meta-classes can be represented by a small amount of data and the performance of learning algorithm can be degraded. [43] proposed the Pruned Problem Transformation (PPT) to improve the LP; patterns with too rarely occurring labels are simply removed from the training set by considering label sets with predefined minimum occurrence τ . This modification was also used by [40]. The main limitations of this approach are: loss of class information due to removing some observations and the necessity of choosing the optimal value of τ . Apart from the above drawbacks, in LP-based methods we loose information about dependency structure, i.e. about which labels and how are influenced by a given feature.

Finally, let us also discuss other methods used for dimensionality reduction. Wrappers allow to assess subsets of features using some criterion function, e.g. prediction error on validation set. To avoid fitting models on all possible subsets, usually some search strategies are used, e.g. forward selection or backward elimination. The main limitation of wrappers is a significant computational cost, due to training large number of classifiers. Another important group of methods are so-called embedded feature selection procedures, in which the selection of features is an integral element of the learning process. Examples from this group are: multi-label version of decision trees proposed by [44] in which the useful features are chosen during building the tree or methods based on l_1 regularization [45, 33].

The other important group are feature transformation methods which aim to identify functions of features that can replace the original ones, e.g. Principal Component Regression [46] or Partial Least Squares Regression [47, 48]. Recently Partial Least Squares method has been successfully used in MLC [49]. Let us also mention about using Canonical Correlation Analysis in multi-label learning [50]. The comprehensive list of feature transformation methods in MLC is given in [51].

3. Feature importance measures based on the Ising model

Before formal description of our method, let us introduce some basic notations for the multi-label learning. For the convenience of a reader, vectors and matrices are written in bold. Let $\mathbf{y} = (y_1, \dots, y_K)^T$ be a label vector containing K binary labels and let $\mathbf{x} = (x_1, \dots, x_p)^T$ be a set of p input features. By \mathbf{y}_{-k} we denote vector \mathbf{y} with k -th label removed. Further, let \mathbf{Y} ($n \times K$) and \mathbf{X} ($n \times p$) be matrices containing instances of \mathbf{y} and \mathbf{x} , respectively, in rows. Analogously,

let \mathbf{Y}_k be k -th column of \mathbf{Y} and \mathbf{Y}_{-k} be a matrix \mathbf{Y} with k -th column removed. Similarly, let \mathbf{X}_j be j -th column of \mathbf{X} . Finally, the superscript (i) will correspond to i -th instance, e.g. $\mathbf{X}^{(i)}$ is i -th row of \mathbf{X} and $\mathbf{X}_j^{(i)}$ is i -th instance of j -th column (feature) of \mathbf{X} . The main task in multi-label learning is to build a model based on training data (\mathbf{X}, \mathbf{Y}) which predicts unknown labels for some new objects. The main goal of FR methods is to evaluate the relevance of features x_1, \dots, x_p in predicting labels y_1, \dots, y_K , based on data (\mathbf{X}, \mathbf{Y}) . For simplicity, we denote by $y_k \sim x_1, \dots, x_p$ a classification problem in which y_k is a response (target) variable and x_1, \dots, x_p are input features.

3.1. Ising model with constant interaction terms

We start from a simple model in which interactions between labels do not depend on features. To assess how the individual feature x_j influences the joint distribution of labels we use the Ising model

$$P(y_1, \dots, y_K | x_j) = \frac{1}{N(x_j)} \exp \left[\sum_{k=1}^K a_k x_j y_k + \sum_{k < l} \beta_{k,l} y_k y_l \right], \quad (1)$$

where $a_k, \beta_{k,l} \in \mathbb{R}$ are parameters and

$$N(x_j) = \sum_{\mathbf{y} \in \{0,1\}^K} \exp \left[\sum_{k=1}^K a_k x_j y_k + \sum_{k < l} \beta_{k,l} y_k y_l \right] \quad (2)$$

is normalizing constant which ensures that the exponential functions sum up to 1. Note that the normalizing term depends on x_j but does not depend on \mathbf{y} . It is assumed that $\beta_{k,l} = \beta_{l,k}$. Parameters a_k describe the individual contribution of the labels, whereas $\beta_{k,l}$ correspond to interactions between labels. Note that our model (1) is identical to CORRLog model used in [34]. Number of authors consider unconditional version of (1) to model $P(y_1, \dots, y_n)$ (e.g. [45]). In statistical literature, the unconditional version of (1) is referred to as auto-logistic model [52, 53]. Model (1) describes a simple Markov Network (or more specifically Conditional Random Field, [54]) in which vertices correspond to labels and a given feature whereas edges correspond to dependencies. Labels y_k and y_l are conditionally independent given x_j and all other labels if and only if $\beta_{k,l} = 0$ (no edge between y_k and y_l). The advantage of the above model is that it indicates which labels are influenced by feature x_j . The feature x_j is not relevant when $a_1 = \dots = a_K = 0$ (no edges between x_j and the rest of the graph). A natural way to assess the relevance of x_j would be to estimate parameters in model (1) (using e.g. maximum likelihood approach) and then to perform some statistical test to verify whether $a_k \neq 0$. However it is difficult to estimate unknown parameters in (1) directly by maximizing the joint conditional log-likelihood since the probability in (1) includes the normalizing term, which requires summation of 2^K terms for each data point and makes it intractable for direct maximization. Instead we use simple procedure via node-wise regressions suggested in [45]. First it is easy to verify (see Appendix A.1 for the proof) that

$$\log \left[\frac{P(y_k = 1 | x_j, \mathbf{y}_{-k})}{P(y_k = 0 | x_j, \mathbf{y}_{-k})} \right] = \sum_{l: l \neq k} \beta_{k,l} y_l + a_k x_j. \quad (3)$$

It follows from (10) that in order to estimate parameter vector

$$\boldsymbol{\theta}_k = (\beta_{k,1}, \dots, \beta_{k,k-1}, \beta_{k,k+1}, \dots, \beta_{k,K}, a_k)^T \in \mathbb{R}^K,$$

it suffices to fit logistic model $y_k \sim \mathbf{y}_{-k}, x_j$ in which y_k is a response variable, whereas labels $y_1, \dots, y_{k-1}, y_{k+1}, \dots, y_K$ and feature x_j are the explanatory variables. The crucial in the above idea, is that the normalizing term $N(x_j)$ is eliminated. Now to assess the relevance of feature x_j , we propose to use the score statistic [37] to compare logistic models $y_k \sim \mathbf{y}_{-k}$ and $y_k \sim \mathbf{y}_{-k}, x_j$. Let

$$\hat{\boldsymbol{\theta}}_k = (\hat{\beta}_{k,1}, \dots, \hat{\beta}_{k,k-1}, \hat{\beta}_{k,k+1}, \dots, \hat{\beta}_{k,K}, 0)^T \in \mathbb{R}^K$$

be the maximum likelihood estimator in the smaller model $y_k \sim \mathbf{y}_{-k}$. We augment it to K - dimensional vector by setting the last coordinate to 0. Define $n \times K$ matrix $\mathbf{Z} = (\mathbf{Y}_{-k}, \mathbf{X}_j)$. In our setting, the score statistic is defined as

$$u_k(x_j) := |s^2(\hat{\boldsymbol{\theta}}_k) / v(\hat{\boldsymbol{\theta}}_k)|, \quad (4)$$

where

$$s(\hat{\theta}_k) := \sum_{i=1}^n \mathbf{X}_j^{(i)} (\mathbf{Y}_k^{(i)} - p^{(i)}(\hat{\theta}_k)),$$

$$p^{(i)}(\hat{\theta}_k) = \frac{\exp(\hat{\theta}_k^T \mathbf{Z}^{(i)})}{1 + \exp(\hat{\theta}_k^T \mathbf{Z}^{(i)})},$$

$$v(\hat{\theta}_k) = D(\hat{\theta}_k) - \mathbf{C}(\hat{\theta}_k) \mathbf{A}^{-1}(\hat{\theta}_k) \mathbf{B}(\hat{\theta}_k)$$

where

$$\mathbf{A}(\hat{\theta}_k) = \mathbf{Y}_{-k}^T \mathbf{W}(\hat{\theta}_k) \mathbf{Y}_{-k},$$

$$\mathbf{B}(\hat{\theta}_k) = \mathbf{Y}_{-k}^T \mathbf{W}(\hat{\theta}_k) \mathbf{X}_j,$$

$$\mathbf{C}(\hat{\theta}_k) = \mathbf{X}_j^T \mathbf{W}(\hat{\theta}_k) \mathbf{Y}_{-k},$$

$$D(\hat{\theta}_k) = \mathbf{X}_j^T \mathbf{W}(\hat{\theta}_k) \mathbf{X}_j$$

and $\mathbf{W}(\hat{\theta}_k)$ is $n \times n$ diagonal matrix with $p^{(i)}(\hat{\theta}_k)(1 - p^{(i)}(\hat{\theta}_k))$ on diagonal. Observe that $s(\hat{\theta}_k)$ measures the correlation between added feature and residuals obtained from the smaller model. The main advantage of using the score statistic is that $\hat{\theta}_k$, $\mathbf{W}(\hat{\theta}_k)$ and $\mathbf{A}^{-1}(\hat{\theta}_k)$ do not involve \mathbf{X}_j and thus these terms need to be calculated only once. Computing the remaining terms: $\mathbf{B}(\hat{\theta}_k)$, $\mathbf{C}(\hat{\theta}_k)$, $D(\hat{\theta}_k)$ and $s(\hat{\theta}_k)$ can be done very quickly and stably, even for thousands of features x_j . So computation of the score statistics requires fitting only the smaller model $y_k \sim \mathbf{y}_{-k}$. This is not the case for other popular statistics like the Wald statistic or the Likelihood Ratio statistic [55] which involve fitting both $y_k \sim \mathbf{y}_{-k}$ and $y_k \sim \mathbf{y}_{-k}, x_j$ models. In Section 3.2 we prove that, for relevant feature x_j ($a_k \neq 0$), the score statistic tends to infinity when sample size increases and moreover we show that the lower bound of the score statistic is an increasing function of $|a_k|$.

Observe that the larger the value of $u_k(x_j)$, the more important is a feature x_j in model $y_k \sim \mathbf{y}_{-k}, x_j$. We check the usefulness of x_j for predicting k -th label, when all remaining labels are present in the model. In other words, we test how much adding a feature x_j to labels \mathbf{y}_{-k} improves the prediction of y_k . Consider the following toy examples with one feature x_1 and two labels y_1, y_2 . The example shows that adding feature x_1 to y_2 may improve prediction of y_1 . Consider two binary labels y_1, y_2 , such that $P(y_2 = 1) = 0.5$ and binary feature x_1 , such that $P(x_1 = 1) = 0.5$ and assume that x_1 is independent from y_2 and $y_1 = I(y_2 + x_1 > 0)$ (where I is indicator function). It is seen that y_1 can be predicted by y_2 with maximal accuracy 75% and similarly y_1 can be predicted by x_1 with accuracy 75%. On the other hand when y_1 is explained by both y_2 and x_1 , the accuracy is 100%.

3.2. Properties of the score statistic

In this section we study some theoretical properties of the score statistic (4). Recall that the score statistic is used to test the significance of feature x_j in model $y_k \sim \mathbf{y}_{-k}, x_j$. The score statistic is a classical measure, proposed more than 60 years ago [37], however recently it has attracted again a significant attention in high-dimensional problems, mainly because of its low computational cost and good performance. For example, the score statistic has been successfully used for feature ranking in analysing Genome Wide Association Studies [56].

It is well known fact that when x_j is not significant, i.e. $a_k = 0$, then, under some regularity conditions, the score statistic $u_k(x_j)$ is approximately distributed as chi-squared with 1 degree of freedom, for large sample size n (see e.g. [55]). Thus in the following we will focus on the performance of $u_k(x_j)$, when x_j is significant, i.e. $a_k \neq 0$. Best of our knowledge, the properties of the score statistic under this setting has not yet been discussed.

So assume that $a_k \neq 0$ and we fit the smaller model $y_k \sim \mathbf{y}_{-k}$ from which we have an estimator

$$\hat{\theta}_k = (\hat{\beta}_{k,1}, \dots, \hat{\beta}_{k,k-1}, \hat{\beta}_{k,k+1}, \dots, \hat{\beta}_{k,K}, 0)^T,$$

with the coordinate corresponding to x_j set to 0. Recall that $\mathbf{Z} = (\mathbf{Y}_{-k}, \mathbf{X}_j)$. Let $\lambda_{\min}(\mathbf{A})$ be the minimal eigenvalue of matrix \mathbf{A} . Define $L := \max_{i,j} |\mathbf{X}_j^{(i)}|$ (to simplify a proof we assume $L > 1$), $\Lambda_{\min} := \lambda_{\min}(\mathbf{Z}^T \mathbf{Z}/n)$ and let $-G \leq a_k \leq G$. Constant G determines the range of unknown parameter a_k , corresponding to variable x_j . This constant is introduced to facilitate the proof of Theorem 1.

Theorem 1. *The following inequality holds*

$$P\left[u_k(x_j) \geq \frac{Cna_k^2}{H^4} \middle| \mathbf{Z}\right] \geq 1 - K \exp\left[-\frac{Cn(K+L^2)a_k^2}{2H^2}\right],$$

where $C = \left(\frac{\Lambda_{\min} v}{2e^3 L(K+L)^{3/2}}\right)^2$, $H = \max(1, G)$ and $v = \min_i p^{(i)}(\theta_k)(1 - p^{(i)}(\theta_k))$.

The proof of the above result is provided in Appendix A.4. Let us discuss the meaning of the above result as well as effects of different constants. It follows from the above Theorem that $u_k(x_j) \rightarrow \infty$, with probability tending to 1, as $n \rightarrow \infty$, which is a desired result as $u_k(x_j)$ should take large values when x_j is significant. Moreover, it is seen that the lower bound Cna_k^2/H^4 is an increasing function of $|a_k|$ and decreasing function of K , which is concordant with intuition. Indeed, the larger the value of a_k , the more significant is the feature x_j . It is very useful property as it allows to assess the significance of the feature x_j , without estimating the corresponding unknown coefficient a_k . For large value of K , it is more difficult to test the significance of the feature x_j . Constant v is a minimal (where minimum is taken over all training examples) variance of k -th label, conditioned on the remaining labels and feature x_j . Small value of v indicates that the classes, corresponding to k -th label, are almost separable and thus the logistic model may fail. Very small value of Λ_{\min} indicates that columns of matrix \mathbf{Z} are almost linearly dependent, which may harm the fitting of logistic model. We show in Lemma 4 that $\Lambda_{\min} > 0$ ensures that the likelihood function corresponding to larger model is strictly concave. In addition observe that the lower bound Cna_k^2/H^4 is a decreasing function of v and Λ_{\min} , which is again intuitive: the more difficult the problem, the more challenging is identification of a significant variable. Finally, constants H and G are introduced for technical reasons, to facilitate the proof.

The following example illustrates the above theoretical result on artificial data. Consider one feature x_1 and ten labels y_1, \dots, y_{10} . We generate labels y_2, \dots, y_{10} independently, from binomial distribution with success probability 0.5. Then we generate y_1 , from (10), with $\beta_{k,l} = 0.1$, and x_1 drawn from standard Gaussian distribution. The simulations are repeated 50 times. Figure 2 (a) shows smoothed histograms of the score statistics when x_1 is relevant ($a_1 = 1$) and irrelevant ($a_1 = 0$). In the latter case, the values of the score statistics remain close to zero. Figure 2 (b) shows the score statistics w.r.t. increasing value of a_1 (coefficient corresponding to x_1), for different sample sizes. It is clearly seen that the larger the value of the coefficient, the larger the value of the score statistic.

3.3. Ising model with feature-dependent interaction terms

In real applications, it is often the case that interactions between labels depend on features. In order to model this situation, we expand model (1) as

$$P(y_1, \dots, y_K | x_j) = \frac{1}{N(x_j)} \exp\left[\sum_{k=1}^K a_k x_j y_k + \sum_{k < l} (\beta_{k,l} + b_{k,l} x_j) y_k y_l\right], \quad (5)$$

where $b_{k,l}$ describes how strong feature x_j influences the interactions between y_k and y_l . The price for considering feature-dependent interactions is larger number of parameters. Analogous reasoning as in (10) leads to

$$\log \left[\frac{P(y_k = 1 | x_j, \mathbf{y}_{-k})}{P(y_k = 0 | x_j, \mathbf{y}_{-k})} \right] = \sum_{l: l \neq k} \beta_{k,l} y_l + \sum_{l: l \neq k} b_{k,l} x_j y_l + a_k x_j. \quad (6)$$

It follows from (6) that in order to estimate $\beta_{k,l}, b_{k,l}, a_k$, it suffices to fit logistic model $y_k \sim \mathbf{y}_{-k}, x_j \mathbf{y}_{-k}, x_j$ in which y_k is a response variable, whereas $\mathbf{y}_{-k}, x_j \mathbf{y}_{-k}$ and x_j are the explanatory variables. Now to assess the relevance of feature x_j , we compare models $y_k \sim \mathbf{y}_{-k}$ and $y_k \sim \mathbf{y}_{-k}, x_j \mathbf{y}_{-k}, x_j$. Define vector $\mathbf{m} = (x_j \mathbf{y}_{-k}, x_j)^T$ and let M be $n \times K$ matrix containing realizations of \mathbf{m} in rows. Let $\hat{\theta}_k$ be an estimator based on smaller model $y_k \sim \mathbf{y}_{-k}$. Multivariate version of the score statistic is defined as

$$U_k(x_j) := |\mathbf{S}^T(\hat{\theta}_k) \mathbf{V}^{-1}(\hat{\theta}_k) \mathbf{S}(\hat{\theta}_k)|, \quad (7)$$

where

$$\begin{aligned} \mathbf{S}(\hat{\theta}_k) &:= \mathbf{M}^T (\mathbf{Y}_k - \mathbf{p}(\hat{\theta}_k)), \\ \mathbf{p}(\hat{\theta}_k) &= (p^{(1)}(\hat{\theta}_k), \dots, p^{(n)}(\hat{\theta}_k))^T, \end{aligned}$$

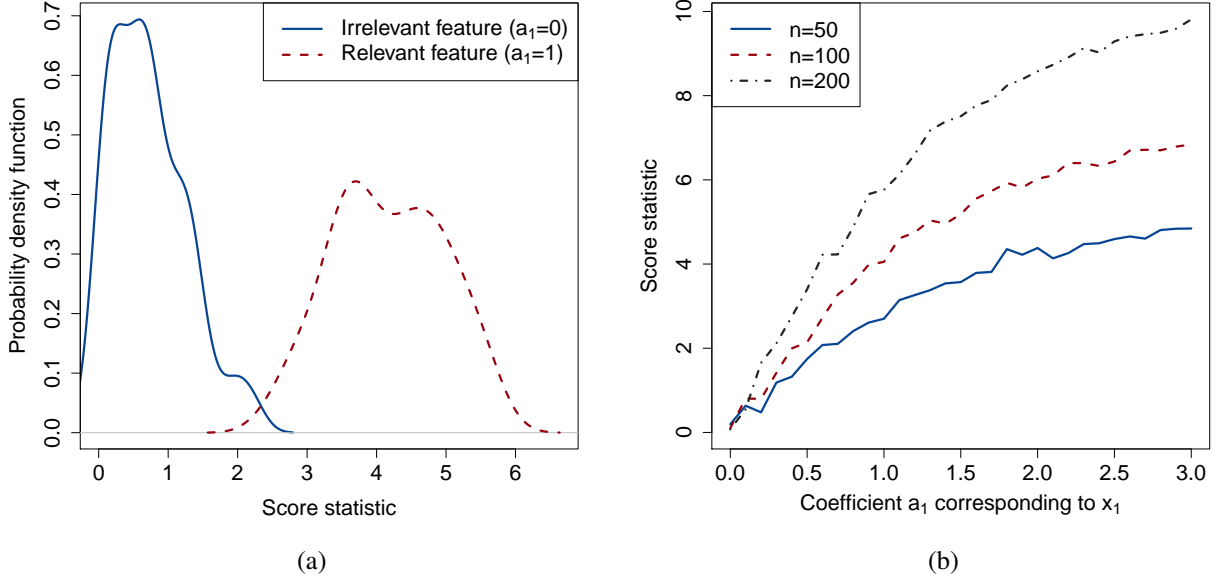


Figure 2: Performance of the score statistics for artificial data.

$$\mathbf{V}(\hat{\theta}_k) = \mathbf{D}(\hat{\theta}_k) - \mathbf{C}(\hat{\theta}_k)\mathbf{A}^{-1}(\hat{\theta}_k)\mathbf{B}(\hat{\theta}_k)$$

where

$$\mathbf{A}(\hat{\theta}_k) = \mathbf{Y}_{-k}^T \mathbf{W}(\hat{\theta}_k) \mathbf{Y}_{-k},$$

$$\mathbf{B}(\hat{\theta}_k) = \mathbf{Y}_{-k}^T \mathbf{W}(\hat{\theta}_k) \mathbf{M},$$

$$\mathbf{C}(\hat{\theta}_k) = \mathbf{M}^T \mathbf{W}(\hat{\theta}_k) \mathbf{Y}_{-k},$$

$$\mathbf{D}(\hat{\theta}_k) = \mathbf{M}^T \mathbf{W}(\hat{\theta}_k) \mathbf{M}.$$

3.4. Ising model and l_1 regularization

Model (1) allows to assess how the single feature x_j influences the joint distribution of labels. To investigate how the whole vector of features $\mathbf{x} \in R^p$ influences the joint distribution of labels we propose to use the following model

$$P(y_1, \dots, y_K | \mathbf{x}) = \frac{1}{N(\mathbf{x})} \exp \left[\sum_{k=1}^K \mathbf{a}_k^T \mathbf{x} y_k + \sum_{k < j} \beta_{k,j} y_k y_j \right], \quad (8)$$

where $\mathbf{a}_k = (a_{k,1}, \dots, a_{k,p})^T$ is a p -dimensional parameter vector and

$$N(\mathbf{x}) = \sum_{\mathbf{y} \in \{0,1\}^K} \exp \left[\sum_{k=1}^K \mathbf{a}_k^T \mathbf{x} y_k + \sum_{k < j} \beta_{k,j} y_k y_j \right] \quad (9)$$

is normalizing constant. Note that model (8) incorporates all features $\mathbf{x} = (x_1, \dots, x_p)^T$ simultaneously, not only a single feature x_j as in (1) and (5). Similar reasoning as in (10) leads to

$$\log \left[\frac{P(y_k = 1 | \mathbf{x}, \mathbf{y}_{-k})}{P(y_k = 0 | \mathbf{x}, \mathbf{y}_{-k})} \right] = \sum_{l: l \neq k} \beta_{k,l} y_l + \mathbf{a}_k^T \mathbf{x}, \quad (10)$$

which as in the case of (1) and (5) suggests that the unknown parameters can be estimated using logistic regression. Since the dimension of \mathbf{x} can be large, we use l_1 regularization to estimate parameter vector

$$\boldsymbol{\theta}_k = (\beta_{k,1}, \dots, \beta_{k,k-1}, \beta_{k,k+1}, \dots, \beta_{k,K}, \mathbf{a}_k)^T \in R^{K+p-1},$$

where $\mathbf{a}_k = (a_{k,1}, \dots, a_{k,p})^T$. So the estimate vector is obtained as

$$\hat{\boldsymbol{\theta}}_k = (\hat{\beta}_{k,1}, \dots, \hat{\beta}_{k,k-1}, \hat{\beta}_{k,k+1}, \dots, \hat{\beta}_{k,K}, \hat{\mathbf{a}}_k)^T = \arg \min_{\boldsymbol{\theta}_k \in R^{K+p-1}} \{-l(\boldsymbol{\theta}_k) + \lambda \|\boldsymbol{\theta}_k\|_1\}, \quad (11)$$

where $l_k(\cdot)$ is a logistic log-likelihood function corresponding to (11), $\lambda > 0$ is a parameter and $\|\cdot\|_1$ is l_1 norm. The above procedure was used in [45], [34] and [33]. The advantage of this approach is that we assess the joint relevance of all features, not only the individual relevance of single feature x_j , as in the case of models (1) and (5). On the other hand, this approach is much more computationally demanding than fitting models (1) and (5). We will show in Section 4 how to use the above procedure to construct the ranking of features.

3.5. Why Ising model?

The first advantage of the Ising model is that it precisely describes the dependence structure between labels and a given feature x , i.e. it indicates which labels and interactions of labels are influenced by feature x . Secondly, it follows from (10) and (6) that fitting the model can be done relatively simply. Finally, it turns out that the Ising model is a maximum entropy model, i.e. it maximizes the entropy under some constraints on the expectations of labels and interactions between labels. The details are given below. Assume that we would like to find a conditional distribution of labels $g(\mathbf{y}|x)$, which maximizes the entropy $H_g(\mathbf{y}|x) = -\sum_{\mathbf{y}} g(\mathbf{y}|x) \log(g(\mathbf{y}|x))$ under the standard constraints $g(\mathbf{y}|x) \geq 0$, $\sum_{\mathbf{y}} g(\mathbf{y}|x) = 1$ and two additional constraints:

$$\sum_{\mathbf{y}} g(\mathbf{y}|x) y_k = A_k(x), \quad k = 1, \dots, K, \quad (12)$$

$$\sum_{\mathbf{y}} g(\mathbf{y}|x) y_k y_l = B_{k,l}(x), \quad k < l, \quad (13)$$

where $A_k(x)$ and $B_{k,l}(x)$ are fixed terms dependent on x . The above constraints are very natural in multi-label setting and they simply mean that the expectations of labels as well as interactions between labels (with respect to $g(\mathbf{y}|x)$) are fixed. The following fact is proved in Appendix A.2.

Proposition 1. *Let $g(\mathbf{y}|x)$ be any probability function satisfying (12), (13) and let $p(\mathbf{y}|x)$ be probability of the form (5) also satisfying constraints (12), (13). Then $H_g(\mathbf{y}|x) \leq H_p(\mathbf{y}|x)$.*

Thus, according to the above Proposition and the principle of maximum entropy [57], the Ising distribution is the most adequate one in a situation when constraints (12), (13) are taken into account.

4. Feature ranking methods and feature selection methods

In this Section, we show how to use the Ising model described in previous Section, to construct rankings of features. The first approach is based on the Ising model with constant interaction terms and the score statistic. The second approach is based on the Ising model with feature-dependent interaction terms and the score statistic. The third approach, computationally most expensive, is based on fitting l_1 regularized logistic regressions.

4.1. Feature ranking based on the Ising model with constant interaction terms

We use the Ising model with constant interaction terms and the score statistic to construct the first feature importance measure. As a feature importance measure for feature x_j we take $\text{imp}(x_j) := \sum_{k=1}^K u_k(x_j)$. Recall that $u_k(x_j)$ is non-negative. Note that, the more important is a feature x_j for predicting consecutive labels (conditioning on the remaining labels), the greater is the measure $\text{imp}(x_j)$. In addition, the more labels are influenced by x_j , the greater is the measure $\text{imp}(x_j)$.

Based on the above feature importance measure we propose the following FR procedure which consists of two steps. We initially fit the unconditioned Ising model using only labels, which requires fitting K logistic models, with $K - 1$ input features, each. The first step is the price for avoiding LP transformation. Since the models are fitted independently, the loop can be computed in parallel. In the second step we assess whether adding input features improve the fitting of the model from the first step. The second step is very efficient and allows to assess the importance of thousands of features quickly. Thus the method is tailored to the case of large number of features and moderate number of labels. Figure 3 shows networks corresponding to these two steps in the case of three labels. The whole procedure is described by Algorithm 1. In simulation experiments we will refer to this method as *ising+score*.

Algorithm 1: Feature ranking based on the Ising Model with constant interaction terms (*ising+score*)

Data: $\mathbf{X}(n \times p)$, $\mathbf{Y}(n \times K)$
#1 step: fitting the unconditioned Ising model:
for $k \leftarrow 1$ **to** K **do**
 Fit logistic model $y_k \sim \mathbf{y}_{-k}$;
 Obtain terms not involving \mathbf{X} : $\hat{\boldsymbol{\theta}}_k$, $\mathbf{W}(\hat{\boldsymbol{\theta}}_k)$ and $\mathbf{A}^{-1}(\hat{\boldsymbol{\theta}}_k)$;
#2 step: feature ranking:
for $j \leftarrow 1$ **to** p **do**
 for $k \leftarrow 1$ **to** K **do**
 #Low computational cost:
 Calculate terms involving \mathbf{X}_j : $\mathbf{B}(\hat{\boldsymbol{\theta}}_k)$, $\mathbf{C}(\hat{\boldsymbol{\theta}}_k)$, and $D(\hat{\boldsymbol{\theta}}_k)$;
 Compute $u_k(x_j)$ using $\hat{\boldsymbol{\theta}}_k$, $\mathbf{W}(\hat{\boldsymbol{\theta}}_k)$, $\mathbf{A}^{-1}(\hat{\boldsymbol{\theta}}_k)$ and $\mathbf{B}(\hat{\boldsymbol{\theta}}_k)$, $\mathbf{C}(\hat{\boldsymbol{\theta}}_k)$, $D(\hat{\boldsymbol{\theta}}_k)$;
 Compute feature importance measure $imp(x_j) = \sum_{k=1}^K u_k(x_j)$;
Sort features with decreasing order of imp : $imp(x_{j_1}) \geq \dots \geq imp(x_{j_p})$
Output: Ordered list of features x_{j_1}, \dots, x_{j_p}

4.2. Feature ranking based on the Ising model with feature-dependent interaction terms

We use the Ising model with feature-dependent interaction terms and the score statistic to construct the second feature importance measure. As a feature importance measure for feature x_j we can take $imp(x_j) := \sum_{k=1}^K U_k(x_j)$. Alternatively, one can verify the individual contribution of each interaction by taking $imp(x_j) := \sum_{k=1}^K [u_k(x_j) + \sum_{s:s \neq k} u_k(x_j y_s)]$. Although the advantage of former version is that we assess the joint contributions of all interactions, it is computationally more demanding and less stable due to inversion of $\mathbf{V}(\hat{\boldsymbol{\theta}}_k)$. Since the performances of these two versions were similar, in next sections we present the results for the second version. The whole procedure is described by Algorithm 2. In simulation experiments we will refer to this method as *ising inter+score*. The following toy example, with two labels y_1, y_2 and one feature x_1 , shows that method based on the Ising model with constant interactions may fail, whereas the improved version with feature-dependent interactions will succeed. Consider XOR problem in which $x_1 = 0$ for $(y_1, y_2) = (1, 1)$ or $(y_1, y_2) = (0, 0)$ and $x_1 = 1$ for $(y_1, y_2) = (1, 0)$ or $(y_1, y_2) = (0, 1)$. In this case feature x_1 should be recognized as significant as it partly determines the combination of labels although it is independent from both labels. Adding x_1 to logistic model $y_1 \sim y_2$ does not improve the model fitting and will result in the score statistic close to 0. On the other hand, adding both x_1 and $x_1 y_2$ to model $y_1 \sim y_2$ improves the model fitting significantly and thus will result in large value of the score statistic.

4.3. Feature ranking based on the Ising model and l_1 regularization

The third feature importance measure is based on the Ising model and l_1 regularization. Let $\hat{\mathbf{a}}_k = (\hat{a}_{k,1}, \dots, \hat{a}_{k,p})^T$ be the estimate vector which optimizes function (11). Observe that using l_1 regularization encourages sparsity, i.e. some coefficients $\hat{a}_{k,j}$ will be exactly zero. We define the feature importance measure as $imp(x_j) := \sum_{k=1}^K |\hat{a}_{k,j}|$, i.e. we aggregate the coefficients describing the influence of feature x_j on labels, in a presence of the remaining labels and remaining features. Note that $\hat{a}_{k,j}$ depends on parameter λ in (11); in experiments we take $\lambda = 0.0001\lambda_{\max}$, where λ_{\max} is a value of λ for which all coordinates of $\hat{\boldsymbol{\theta}}_k$ are exactly 0. The whole procedure is described by Algorithm 3. In simulation experiments we will refer to this method as *ising+l1*.

Algorithm 2: Feature ranking based on the Ising Model with feature-dependent interaction terms (*ising inter+score*)

Data: $\mathbf{X}(n \times p)$, $\mathbf{Y}(n \times K)$

#1 step: fitting the unconditioned Ising model:

for $k \leftarrow 1$ **to** K **do**

 Fit logistic model $y_k \sim \mathbf{y}_{-k}$;

 Obtain terms not involving \mathbf{X} : $\hat{\boldsymbol{\theta}}_k$, $\mathbf{W}(\hat{\boldsymbol{\theta}}_k)$ and $\mathbf{A}^{-1}(\hat{\boldsymbol{\theta}}_k)$;

#2 step: feature ranking:

for $j \leftarrow 1$ **to** p **do**

for $k \leftarrow 1$ **to** K **do**

 Calculate terms involving \mathbf{X}_j : $\mathbf{B}(\hat{\boldsymbol{\theta}}_k)$, $\mathbf{C}(\hat{\boldsymbol{\theta}}_k)$, and $D(\hat{\boldsymbol{\theta}}_k)$;

 Compute $u_k(x_j)$ and $u_k(x_j y_1), \dots, u_k(x_j y_{k-1}), u_k(x_j y_{k+1}), \dots, u_k(x_j y_K)$ using $\hat{\boldsymbol{\theta}}_k$, $\mathbf{W}(\hat{\boldsymbol{\theta}}_k)$, $\mathbf{A}^{-1}(\hat{\boldsymbol{\theta}}_k)$ and

$\mathbf{B}(\hat{\boldsymbol{\theta}}_k)$, $\mathbf{C}(\hat{\boldsymbol{\theta}}_k)$, $D(\hat{\boldsymbol{\theta}}_k)$;

 Compute feature importance measure $\text{imp}(x_j) := \sum_{k=1}^K [u_k(x_j) + \sum_{s:s \neq k} u_k(x_j y_s)]$;

Sort features with decreasing order of imp : $\text{imp}(x_{j_1}) \geq \dots \geq \text{imp}(x_{j_p})$

Output: Ordered list of features x_{j_1}, \dots, x_{j_p}

Algorithm 3: Feature ranking based on the Ising Model and l_1 regularization (*ising+l1*)

Data: $\mathbf{X}(n \times p)$, $\mathbf{Y}(n \times K)$

#1 step: Fitting the Ising model using l_1 regularized logistic regressions

for $k \leftarrow 1$ **to** K **do**

 Compute $\hat{\boldsymbol{\theta}}_k = (\hat{\beta}_{k,1}, \dots, \hat{\beta}_{k,k-1}, \hat{\beta}_{k,k+1}, \dots, \hat{\beta}_{k,K}, \hat{\mathbf{a}}_k)^T = \arg \min_{\boldsymbol{\theta}_k \in \mathbb{R}^{p+K-1}} \{-l_k(\boldsymbol{\theta}_k) + \lambda \|\boldsymbol{\theta}_k\|_1\}$,

 where $\hat{\mathbf{a}}_k = (\hat{a}_{k,1}, \dots, \hat{a}_{k,p})^T \in \mathbb{R}^p$;

#2 step: feature ranking:

for $j \leftarrow 1$ **to** p **do**

 Compute feature importance measure $\text{imp}(x_j) := \sum_{k=1}^K |\hat{a}_{k,j}|$;

Sort features with decreasing order of imp : $\text{imp}(x_{j_1}) \geq \dots \geq \text{imp}(x_{j_p})$

Output: Ordered list of features x_{j_1}, \dots, x_{j_p}

4.4. Feature selection

Feature ranking procedures, presented above, allow to order features with respect to their importances, i.e. they produce an ordered list of features x_{j_1}, \dots, x_{j_p} , where x_{j_1} is the most relevant feature, whereas x_{j_p} is the least relevant feature. Below we describe how we choose the final subset of features based on the ranking. Assume that we have multi-label classifier $C(x_1, \dots, x_p)$ which takes as an input features x_1, \dots, x_p and returns multi-label output. This classifier is used as a final model. In simulation experiments, classifier chains [58, 59] were used as a final classifier $C(\cdot)$. First, we split our data into training and validation sets (in simulation experiments: 70% for training and 30% for validation). Training data is used to obtain ranking of features: x_{j_1}, \dots, x_{j_p} and build classifiers. We train classifiers $C(x_{j_1}), C(x_{j_1}, x_{j_2}), \dots, C(x_{j_1}, \dots, x_{j_L})$, where $L < p$ and choose a subset $\{x_{j_1}, \dots, x_{j_s}\}$ ($s \leq L$), for which classifier $C(x_{j_1}, \dots, x_{j_s})$ achieves the maximal value of some evaluation measure calculated on validation set. In simulation experiments we use subset accuracy as an evaluation measure. Observe that the final classifier $C(\cdot)$ is built on subsets of features whose size does not exceed L (we set $L = 0.2p$). This is very natural approach in the case of large number of possible features p and when the final classifier $C(\cdot)$ requires a significant computational effort (and thus cannot be easily trained on all possible features). Moreover, usually in real data only a small set of features influences the values of labels. In the desired feature ranking, the relevant features should precede the spurious ones, and in a consequence the final classifier is built on the subset of most relevant features. When the ranking of features is poor (i.e. there are many spurious features on the top of the list), the resulting classifier will perform poorly.

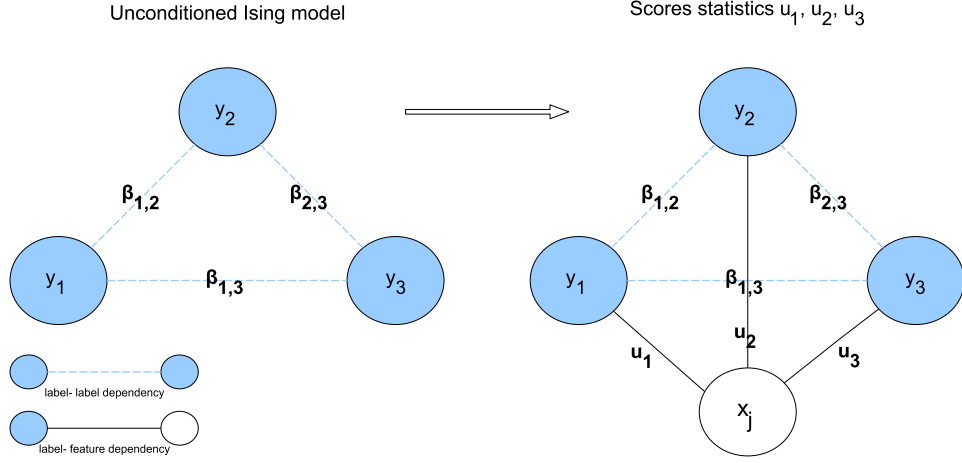


Figure 3: Example scheme of feature importance assessment for 3 labels.

5. Experimental results

In this section we evaluate the effectiveness of the proposed methods by comparing its performance against conventional FR methods. We consider two versions of our method: the first one is based on the Ising model with constant interactions (denoted by *ising+score*); the second one is based on the Ising model with feature-dependent interactions (*ising inter+score*). Moreover we consider a method which is based on the Ising model and l_1 regularization (*ising+l1*). We use two state-of-the-art methods based on BR transformation combined with chi squared statistic (*br chi2*) and information gain (*br ig*). In addition we use two conventional methods based on LP transformation combined with chi squared statistic (*lp chi2*) or information gain (*lp ig*). We also experimented with very simple OneR filter [60], but the performance was disappointing, so we do not present the results here. We carried out experiments on both artificial and real datasets.

5.1. Evaluation measures

To evaluate the performance of feature ranking methods on artificial data, we use a form of ROC curves, constructed in the following way (the similar evaluation can be found in [33]). Let i_1, \dots, i_p be the ranking of features from given feature ranking method (where i_1 corresponds to a feature recognized as the most significant by an algorithm) and t be a set of true relevant features. Let $TPR(k) := |\{i_1, \dots, i_k\} \cap t|/k$, and $FPR(k) := |\{i_1, \dots, i_k\} \setminus t|/|t^C|$, where $|\cdot|$ is set cardinality and t^C is set complement. So, $TPR(k)$ indicates how many relevant features are among top k ones whereas $FPR(k)$ indicates how many redundant features are among top k ones. Now, ROC curve is defined as $(FPR(k), TPR(k))$, $k = 1, \dots, p$. Observe that $AUC = 1$ corresponds to perfect ordering of features, i.e. all relevant features precede spurious ones in the ranking. On the other hand $AUC \approx 0.5$ corresponds to random ordering of features. Each curve is smoothed over 20 simulations. This type of evaluation is used to provide an attractive visualization. Note that our ROC curves differ from standard ROC curves used to evaluate the performance of classification models, although the idea is very similar. Our ROC curves are used to evaluate the quality of rankings, not the classification performance.

For real data, we cannot produce ROC curves described above as the relevant features are not known. Thus, in the case of real data, we use standard evaluation measures, described below. Let $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_K)^T$ be a vector of predicted labels and $\mathbf{y} = (y_1, \dots, y_K)^T$ be a vector of true labels. We consider the following evaluation measures

$$\text{Subset accuracy}(\mathbf{y}, \hat{\mathbf{y}}) = I[\mathbf{y} = \hat{\mathbf{y}}],$$

$$\text{Hamming measure}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{K} \sum_{k=1}^K I[y_k = \hat{y}_k],$$

$$\text{Jaccard measure}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{\sum_{k=1}^K I(y_k = 1 \text{ and } \hat{y}_k = 1)}{\sum_{k=1}^K I(y_k = 1 \text{ or } \hat{y}_k = 1)}$$

The measures are averaged over all instances in test set. The higher the above measures, the better the performance. The measures demonstrate different aspects of multi-label classification performance. Subset accuracy corresponds to subset 0 – 1 loss and measures the correctness of joint prediction for all labels; Hamming measure corresponds to Hamming loss and measures averaged number of correct predictions; Jaccard measure indicates how many labels are correctly predicted as 1 among those equal 1 or predicted as 1.

5.2. Artificial data

5.2.1. Correct specification

The first two datasets are generated under correct specification, i.e. we generate data from the Ising model. The data generation scheme is as follows. We fix the dimension of features $p = 50$, sample size $n = 1000$, number of labels $K = 10$. Features are generated independently from Gaussian distribution with zero mean and identity covariance matrix. Labels are generated from the following Ising model

$$P(y_1, \dots, y_K | \mathbf{x}) = \frac{1}{N(\mathbf{x})} \exp \left[\sum_{k=1}^K \mathbf{a}_k^T \mathbf{x} y_k + \sum_{k < j} \beta_{k,j} y_k y_j + \sum_{k < j} \mathbf{b}_{k,j}^T \mathbf{x} y_k y_j \right], \quad (14)$$

where $\mathbf{a}_k = (a_{k,1}, \dots, a_{k,p})^T$ and $\mathbf{b}_{k,j} = (b_{k,j,1}, \dots, b_{k,j,p})^T$ are p -dimensional parameter vectors. Model (14), from which we generate data, incorporates all features $\mathbf{x} = (x_1, \dots, x_p)^T$ simultaneously, not only single feature x as in (1) and (5). We consider the following two settings.

ArtData1 Let $t = \{1, \dots, 10\}$ be a set of true relevant features. We set $a_{k,s} = 0.2$, for $s \in t$ and $a_{k,s} = 0$ for $s \notin t$; $\beta_{k,j} = 0.1$; $b_{1,2,s} = b_{2,1,s} = 0.2$, for all $s \in t$, $b_{1,2,s} = b_{2,1,s} = 0$, for $s \notin t$ and $b_{k,j,s} = 0$, for $k, j \notin \{1, 2\}$.

ArtData2 Let $t = \{1, \dots, 10\}$ be a set of true relevant features. We set $a_{k,s} = 0$ for all s ; $\beta_{k,j} = 0.1$; $b_{1,2,s} = b_{2,1,s} = 0.2$, for all $s \in t$, $b_{1,2,s} = b_{2,1,s} = 0$, for $s \notin t$ and $b_{k,j,s} = 0$, for $k, j \notin \{1, 2\}$.

In both datasets the first 10 features are significant. In *ArtData2*, significant features do not influence the labels directly but only the interactions between labels, which makes identification of them much more challenging. Given feature vector for i -th observation $\mathbf{X}^{(i)}$ and parameters defined above, we use Gibbs sampling to generate labels, where we iteratively generate $\mathbf{Y}_k^{(i)}$, $k = 1, \dots, K$ from Bernoulli distribution with probability $P(\mathbf{Y}_k^{(i)} = 1 | \mathbf{X}^{(i)}, \mathbf{Y}_{-k}^{(i)})$ and take the last value of the sequence. The number of repetitions in Gibbs sampling is set to 30. The above procedure is repeated for all $i = 1, \dots, n$.

5.2.2. Incorrect specification

Obviously, data generation scheme presented in the previous section favours the methods based on the Ising model. It is interesting to investigate the performance of discussed methods, under incorrect specification, i.e. when data generation scheme is not related to the Ising model. For this purpose we generated two datasets closely related to the ones proposed in [40]. We consider the following two settings.

ArtData3 We draw 5 features x_1, \dots, x_5 from uniform distribution on the $[0, 1]$ interval. Then we construct features $j = 6, \dots, 10$ as follows: $x_6 = (x_1 - x_2)/2$, $x_7 = (x_1 + x_2)/2$, $x_8 = x_3 + 0.1$, $x_9 = x_4 - 0.2$ and $x_{10} = 2x_5$. Then we add additional 40 features from the uniform distribution on the $[0, 1]$, which are independent from x_1, \dots, x_{10} . The multi-label output is build as follows

$$\begin{cases} y_1 = 1 & \text{if } x_1 > x_2 \\ y_2 = 1 & \text{if } x_4 > x_3 \\ y_3 = 1 & \text{if } y_1 + y_2 = 1 \\ y_4 = 1 & \text{if } x_5 > 0.8 \\ y_k = 0 & \text{otherwise } (k = 1, \dots, 4) \end{cases} \quad (15)$$

ArtData4 We draw 5 features x_1, \dots, x_5 from uniform distribution on the $[0, 1]$ interval. Let ϵ be drawn from Gaussian distribution with 0 mean and standard deviation equal to 0.3. We construct features $j = 6, \dots, 10$ as $x_6 =$

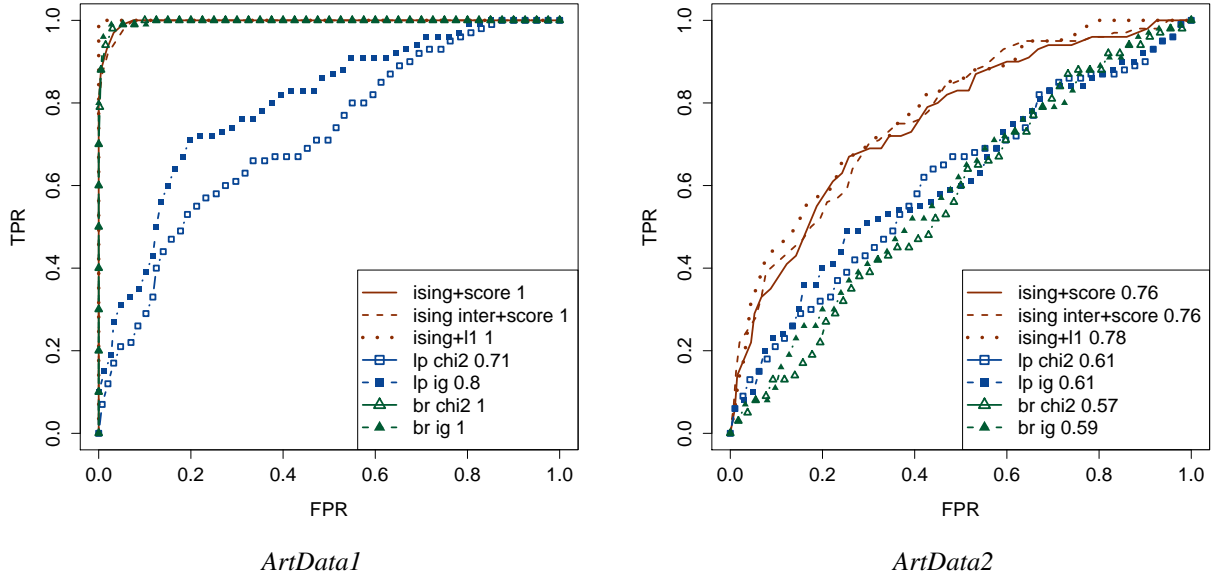


Figure 4: ROC curves for *ArtData1* and *ArtData2*. Numbers in legend correspond to AUC (AUC = 1 corresponds to perfect ordering of features).

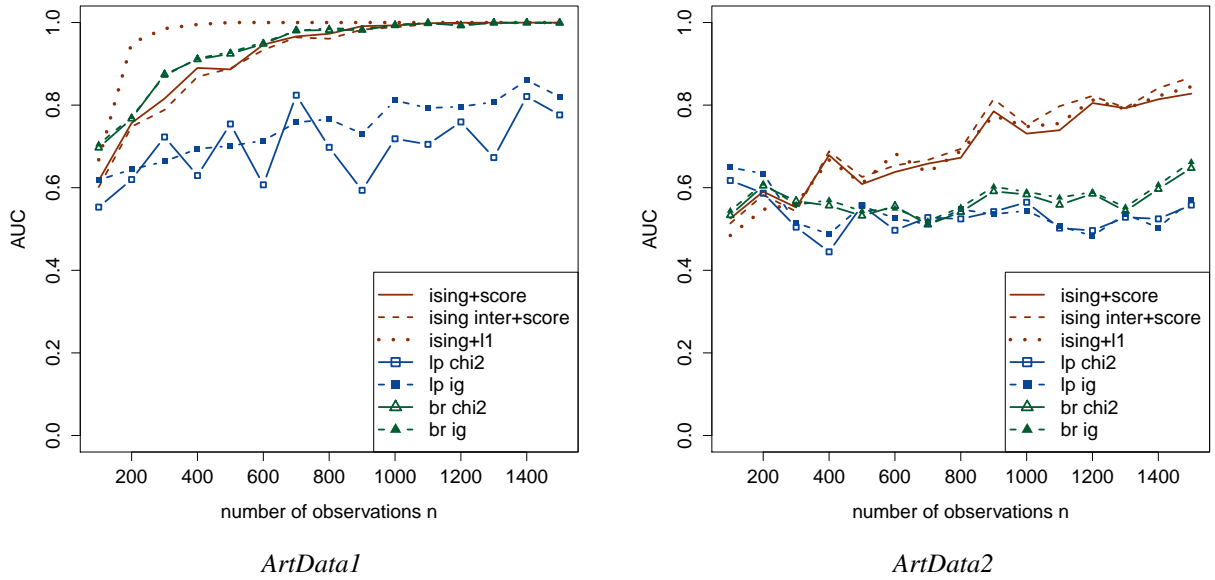


Figure 5: AUC vs sample size n for *ArtData1* and *ArtData2* (AUC = 1 corresponds to perfect ordering of features).

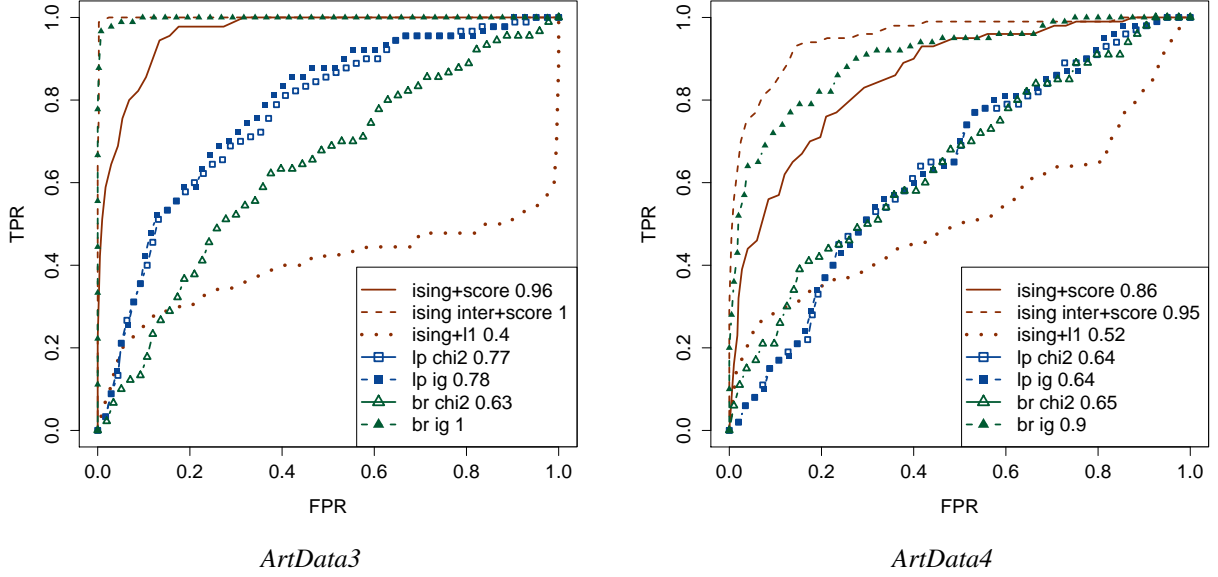


Figure 6: ROC curves for *ArtData3* and *ArtData4*. Numbers in legend correspond to AUC (AUC = 1 corresponds to perfect ordering of features).

$x_1 + \epsilon$, $x_7 = x_2 + \epsilon$, $x_8 = x_3 + \epsilon$, $x_9 = x_4 + \epsilon$, $x_{10} = x_5 + \epsilon$. Then we add additional 40 features from uniform distribution on the $[0, 1]$, which are independent from x_1, \dots, x_{10} . The multi-label output is build as follows

$$\begin{cases} y_1 = 1 & \text{if } x_1 > x_2 + \epsilon \\ y_2 = 1 & \text{if } x_4 > x_3 + \epsilon \\ y_3 = 1 & \text{if } y_1 + y_2 = 1 \\ y_4 = 1 & \text{if } x_5 + \epsilon > 0.8 \\ y_k = 0 & \text{otherwise } (k = 1, \dots, 4) \end{cases} \quad (16)$$

For *ArtData3* features $t = \{1, \dots, 5, 6, 8, 9, 10\}$ are relevant, whereas for *ArtData4*, $t = \{1, \dots, 10\}$. *ArtData3* was considered in [40], section 5.1. *ArtData4* is modification of *ArtData3*, in which some noise ϵ is introduced. To make the task more challenging we increased the total number of features from 15 [40] to $p = 50$ and decreased the number of observations from $n = 1000$ to $n = 100$. Observe that in *ArtData3* some features can be replaced by others, e.g. y_1 is determined by 2 features: x_1 and x_2 or by a single feature x_6 .

5.3. Experiment 1

The aim of the first experiment was to study the performance of the proposed feature ranking methods on artificial datasets. In the case of artificial datasets, we can assess the quality of feature ranking methods directly as we know which features are significant, i.e. which features influence the joint probability of labels. For attractive visualization, we use ROC curves described in Section 5.1. Recall that desired feature ranking will result in ROC curve significantly above the diagonal and $AUC \approx 1$.

Figure 4 shows ROC curves for *ArtData1* and *ArtData2*. It is seen that identification of true relevant features is much more difficult in the case of *ArtData2*, which is not surprising as the former one incorporates feature-dependent interactions. For *ArtData1*, all methods perform well, except methods based on LP. The proposed methods outperform conventional ones significantly for dataset *ArtData2*. For *ArtData1*, *lp ig* works slightly better than *lp chi2*, which

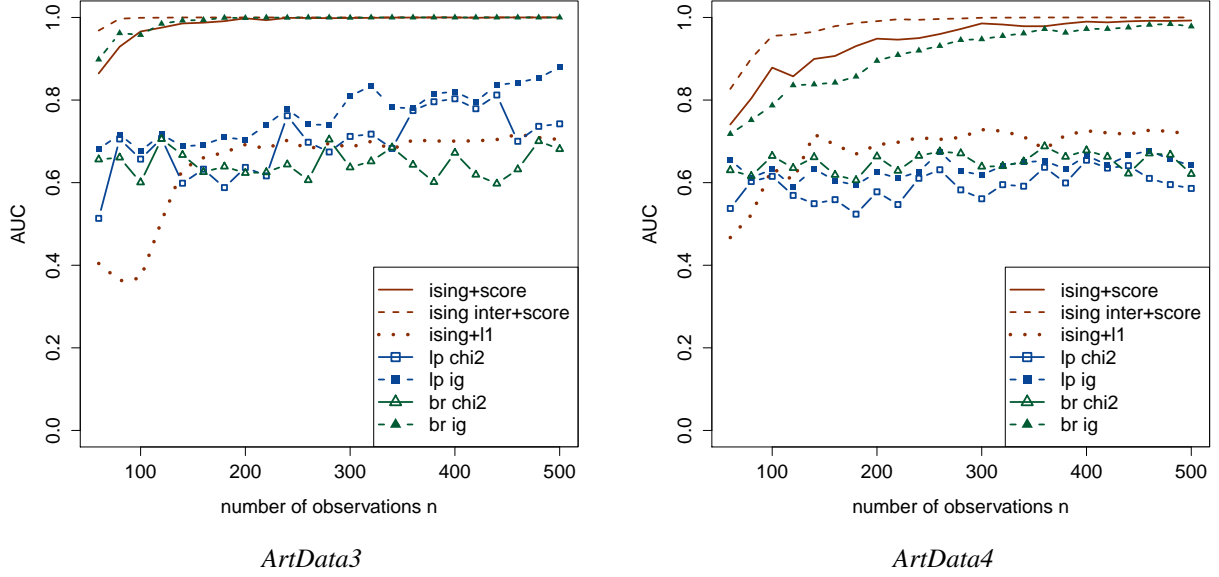


Figure 7: AUC vs sample size n for *ArtData3* and *ArtData4* (AUC = 1 corresponds to perfect ordering of features).

agrees with the conclusions of other authors [40]. In addition, we investigated an effect of data size, i.e. how number of observations affects ranking of features. Figure 5 shows AUC with respect to the sample size n . The proposed methods *ising+score* and *ising inter+score* show good performance. For *ArtData1*, AUC increases with sample size for all methods, but the results for *lp chi2* and *lp ig* are less stable. It is seen that for n large enough, the proposed methods find the correct ranking in all simulations. For *ArtData2*, methods based on BR and LP perform poorly, even for large n . AUC for proposed methods increases, but a rate of growth is smaller than for *ArtData1*.

Figure 6 shows ROC curves for *ArtData3* and *ArtData4*. It is seen that *ArtData4* is more challenging than *ArtData3* for all considered methods. For both datasets, *ising inter+score* outperforms other approaches, *br ig* is second best, whereas methods based on LP perform poorly. Surprisingly, *ising+l1* works poorly for these two datasets. Effect of sample size is shown on Figure 7. Large amount of data facilitates the task in the case of *ArtData3*. For *ArtData4*, the performance of LP does not improve even for large n .

5.4. Real data

We experimented with datasets from different applications. Some datasets, considered in experiments, are publicly available at <http://mulan.sourceforge.net/datasets-mlc.html>. Those are: scene, yeast, genbase, mediamill, medical, nus-wide, eurlex-dc and CAL500. We also consider Twitter dataset, analysed in [61]. The goal was to analyse a collection of tweets in English and discover its authors gender, age and personality traits: extraversion, stability, agreeableness, conscientiousness and openness. Since the original target variables are not binary, we created 7 binary labels using the original target variables in the following way. We set $y_1 = 1$ if *gender = male*; $y_2 = 1$ if *age* ≤ 34 ; $y_i = 1$, for $i = 3, 4, 5, 6, 7$ if the values of target variables extraversion, stability, agreeableness, conscientiousness and openness are greater than their medians, respectively. The details of the data sets are summarized in Table 1. As the proposed methods (in particular *ising inter+score*) are recommended for moderate number of labels, we limited the number of labels to 50, for all datasets, by taking the most frequent ones. The number of features ranges from 55 to 5000.

Dataset	Domain	#observations (n)	#features (p)	#labels (K)
scene	images	2407	294	6
yeast	biology	2417	103	14
genbase	biology	662	1186	27
bibtex	text	7395	1836	50
mediamill	video	10000	120	50
medical	video	978	1449	45
nus-wide	video	10000	500	50
eurlex-dc	video	10000	5000	50
CAL500	video	502	68	50
Twitter	text	152	55	7

Table 1: Basic statistics for the benchmark datasets.

5.5. Experiment 2

The aim of the second experiment was to study the performance of the proposed feature ranking methods on real-world datasets. For real-world datasets, relevant features are not known in advance as they were for artificially built datasets. Thus the quality of feature ranking methods cannot be evaluated directly but can be measured by the performance of a classification model based on selected features. We use feature selection procedure described in Section 4.4. To assess the quality of the considered methods we use measures described in Section 5.1: Subset accuracy, Hamming measure and Jaccard measure. We also calculate the above measures for artificial datasets, described in Section 5.2.

The considered splitting of the samples into training and test sets are the ones suggested on the website of the Mulan project. The training set is used to obtain ranking of features and then to build a multi-label classification model. For artificial datasets and Twitter dataset we randomly split the original data into training set 50% and testing set 50%. The validation set, needed to choose the final subset of features (see Section 4.4), is separated from the training set. The above performance measures are calculated on the test set. As a final classification model we use Classifier Chains [58, 59] with logistic model as a base learner. Our choice is motivated by the fact that classifier chains are among the most frequently used and successful methods in multi-label learning (see [62] for theoretical properties of classifier chains combined with logistic regression). We also experimented with nearest neighbour method, but the results were worse, even when the classifier was built using all possible features. Thus the results for nearest neighbour method are not presented.

The results are presented in Tables 2, 3 and 4. We do not show the results for *ArtData1* and *ArtData2* as the performance of classifier chains was very poor in these cases, for all considered rankings. The reason is that these two datasets were too difficult for the final classifier, even when the final classifier was trained using all possible features. Numbers printed in bold pertain to maximal values in rows (the winning method). The last row contains ranks, averaged over all datasets. Looking at the averaged ranks for Hamming and Jaccard measures, the proposed methods outperform the conventional ones, although the differences between measures are quite small. For Subset measure, *ising+ll* has the highest averaged rank. Surprisingly, *br ig* outperforms *lp ig*, which can be a consequence of the large number of classes produced by LP transformation. Probably LP transformation combined with some pruning strategies would improve the results.

To analyse the results thoroughly, we followed the two-step statistical procedure recommended in [63]. In the first step we use the Friedman test of the null hypothesis that all methods have equal performance. Friedman test is based on averaged ranks. When null hypothesis is rejected a post-hoc test is used to compare methods in a pairwise way. We use Conover post-hoc test [64]. Friedman test suggests that there are significant differences (at a standard significance level 0.05) between methods for Hamming measure (p-value=0.0009) and Subset measure (p-value=0.0002), whereas the differences for Jaccard measure are not statistically significant (p-value=0.082). Thus, we performed the post-hoc tests for Hamming and Subset measures. Results of pairwise comparisons for Hamming and Subset measures are shown in Tables 5 and 6. It is seen that there are significant differences between the proposed methods and methods based on LP transformation. The differences between the proposed methods and the ones based on BR transformation are not significant. This may be due to the fact that the number of datasets included in experiments is quite limited

from a statistical point of view. The overall performance of the proposed approaches is quite promising.

Dataset	ising+score	ising inter+score	ising+l1	lp chi2	lp ig	br chi2	br ig
scene	0.839	0.845	0.844	0.826	0.815	0.825	0.813
yeast	0.783	0.783	0.783	0.781	0.781	0.783	0.781
genbase	0.998	0.994	0.999	0.968	0.996	0.996	0.997
bibtex	0.967	0.968	0.960	0.961	0.968	0.968	0.967
enron	0.883	0.883	0.877	0.847	0.883	0.857	0.857
mediamill	0.869	0.869	0.870	0.868	0.868	0.871	0.869
medical	0.964	0.975	0.976	0.974	0.974	0.969	0.976
nus-wide	0.930	0.931	0.929	0.929	0.929	0.930	0.930
eurlex-dc	0.978	0.978	0.976	0.974	0.976	0.975	0.979
CAL500	0.560	0.558	0.556	0.553	0.554	0.564	0.564
twitter	0.665	0.680	0.613	0.660	0.660	0.618	0.648
ArtData3	0.683	0.684	0.671	0.661	0.661	0.635	0.690
ArtData4	0.591	0.607	0.604	0.567	0.567	0.585	0.620
Average rank	4.92	5.23	4.30	2.07	2.88	3.76	4.81

Table 2: Hamming measure. The average rank is the average of the ranks across all data sets. Numbers in bold pertain to maximal values in rows.

Dataset	ising+score	ising inter+score	ising+l1	lp chi2	lp ig	br chi2	br ig
scene	0.486	0.508	0.512	0.473	0.421	0.452	0.418
yeast	0.179	0.183	0.179	0.178	0.180	0.176	0.177
genbase	0.960	0.915	0.980	0.613	0.940	0.945	0.940
bibtex	0.534	0.510	0.410	0.426	0.516	0.517	0.506
enron	0.043	0.043	0.057	0.016	0.041	0.017	0.012
mediamill	0.121	0.119	0.126	0.118	0.122	0.122	0.116
medical	0.451	0.640	0.674	0.634	0.631	0.540	0.656
nus-wide	0.299	0.299	0.299	0.288	0.294	0.297	0.298
eurlex-dc	0.585	0.579	0.543	0.495	0.524	0.512	0.605
CAL500	0.000	0.000	0.000	0.000	0.000	0.000	0.000
twitter	0.066	0.092	0.092	0.053	0.053	0.079	0.079
ArtData3	0.030	0.030	0.030	0.010	0.010	0.010	0.010
ArtData4	0.010	0.010	0.010	0.000	0.000	0.000	0.010
Average rank	4.80	5.07	5.61	2.30	3.34	3.38	3.46

Table 3: Subset measure. The average rank is the average of the ranks across all data sets. Numbers in bold pertain to maximal values in rows.

5.6. Computational efficiency

The proposed procedures (*ising+score* and *ising inter+score*) require fitting K logistic models (using maximum likelihood method) with $K - 1$ input features each in the first step. This step is computationally fast for moderate number of labels K . The first step is a price for taking into account labels, when assessing the relevance of features. The second step includes computation of the score statistic, which is very simple in the case of *ising+score*: the most expensive operation is a computation of scalar products between a given feature and appropriately weighted $K - 1$ labels (see definitions of $\mathbf{B}(\hat{\theta}_k)$ and $\mathbf{C}(\hat{\theta}_k)$ in Section 3.1). In the case of *ising inter+score* we compute the score statistic for a given variable and for products of a given variable with $K - 1$ labels. This requires much more operations and can be seen as a price for taking into account feature-dependent interactions between labels. The third procedure *ising+l1* is the most computationally expensive as it requires fitting K logistic models with $K + p - 1$ input features, each, using

Dataset	ising+score	ising inter+score	ising+l1	lp chi2	lp ig	br chi2	br ig
scene	0.523	0.544	0.545	0.497	0.454	0.485	0.451
yeast	0.472	0.471	0.471	0.471	0.472	0.469	0.470
genbase	0.983	0.953	0.993	0.659	0.966	0.970	0.969
bibtex	0.154	0.123	0.000	0.026	0.136	0.135	0.116
enron	0.333	0.327	0.307	0.331	0.343	0.311	0.313
mediamill	0.419	0.415	0.411	0.414	0.418	0.419	0.412
medical	0.466	0.666	0.717	0.684	0.643	0.556	0.670
nus-wide	0.067	0.073	0.064	0.048	0.060	0.064	0.068
eurlex-dc	0.077	0.071	0.033	0.004	0.016	0.002	0.100
CAL500	0.382	0.379	0.380	0.378	0.377	0.382	0.378
twitter	0.600	0.615	0.577	0.589	0.589	0.583	0.591
ArtData3	0.507	0.507	0.495	0.481	0.481	0.431	0.520
ArtData4	0.371	0.387	0.429	0.371	0.371	0.390	0.430
Average rank	5.23	4.76	3.92	2.88	3.65	3.30	4.23

Table 4: Jaccard measure. The average rank is the average of the ranks across all data sets. Numbers in bold pertain to maximal values in rows.

	ising+score	ising inter+score	ising+l1	lp chi2	lp ig	br chi2
ising inter+score	1.000					
ising+l1	1.000	0.885				
lp chi2	0.000	0.000	0.002			
lp ig	0.007	0.001	0.142	1.000		
br chi2	0.438	0.128	1.000	0.042	0.903	
br ig	1.000	1.000	1.000	0.000	0.012	0.636

Table 5: P-values of post-hoc Conover test, used to compare all classifiers against each other with respect to Hamming measure.

	ising	ising inter+score	ising+l1	lp chi2	lp ig	br chi2
ising inter+score	1.0000					
ising+l1	0.6638	1.0000				
lp chi2	0.0001	0.0000	0.0000			
lp ig	0.0556	0.0136	0.0004	0.2898		
br chi2	0.0634	0.0162	0.0005	0.2769	1.0000	
br ig	0.0882	0.0242	0.0008	0.2151	1.0000	1.0000

Table 6: P-values of post-hoc Conover test, used to compare all classifiers against each other with respect to Subset measure.

l_1 regularization. To solve the problem we use Cyclic Coordinate Descent (CCD) algorithm proposed by [65]. CCD iterates over all $p + K - 1$ variables until convergence, the maximal number of iterations in our experiment is set to 100. Figure 8 shows how the computational time depends on the number of features p (a) and the number of labels K (b). The experiment was carried out on Work Station with Intel Core i5-3220M CPU, 2.60GHz, 12 GB RAM. All considered methods have been implemented by us in R language; the only exception is an information gain, taken from R package FSelector [66]. We generated artificial data in such a way that features were drawn from standard Gaussian distribution whereas labels were generated from binomial distribution, number of observations was $n = 500$. The curves are smoothed over 5 simulations. In the case of Figure 8 (a) we set $K = 5$ and in the case of Figure 8 (b) we set $p = 50$. Figure 8 (a) indicates that all methods depend linearly on the number of features, except *ising+l1*, which is a price for incorporating all features simultaneously. Computational times are larger for methods based on

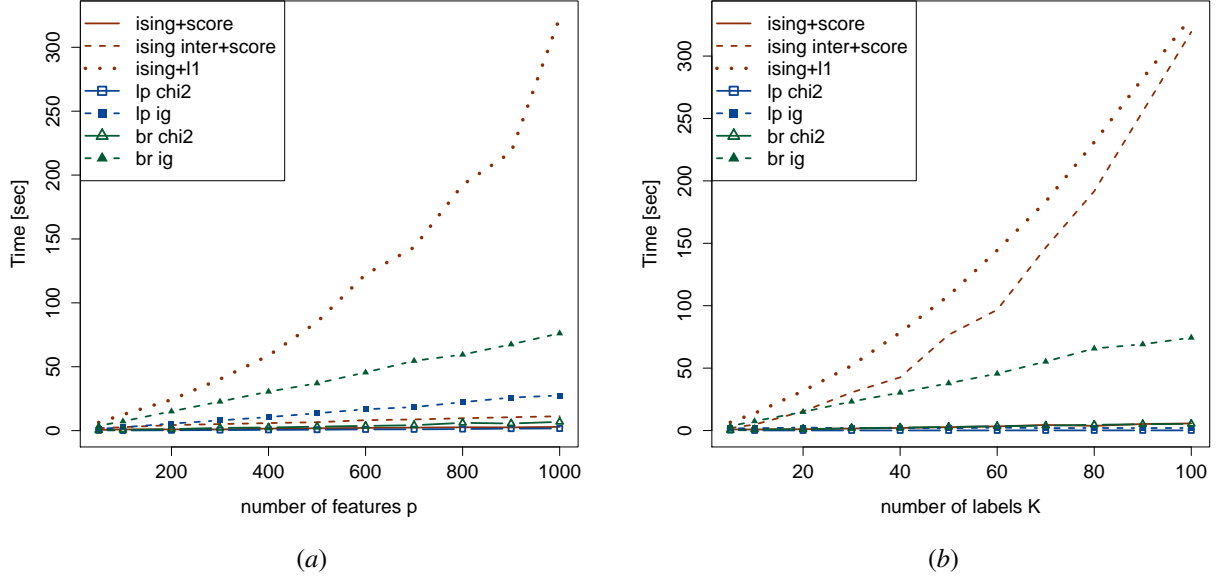


Figure 8: Computational times (in seconds) for considered FR procedures with respect to the number of features p (a) and number of labels K (b).

information gain, which is not a surprise as estimation of information gain is more challenging than computation of the chi-squared statistic or the score statistic. The proposed method *ising+score* is among the fastest ones.

Figure 8 (b) indicates that for *ising inter+score* the dependence between computational time and number of labels is quadratic, which is obvious as this method takes into account feature-dependent interactions between labels. Thus this method can be recommended for limited number of labels.

Finally, let us mention that the proposed methods as well as the conventional ones can be computed in parallel easily (the parallel versions were used in *Experiment 1* and *Experiment 2*).

6. Conclusions and future work

In this paper we propose a novel method for feature ranking in the multi-label setting. The method consists of two steps. In the first step we fit the Ising model using only labels. In the second step we test how much adding a single feature affects the initial network. It is shown that the first step can be performed simply by fitting K logistic models. The second step is based on the score statistic, which is very efficient in this case and allows to test a significance of added features very quickly which is crucial for FR methods. The final feature importance measure is based on averaged values of score statistics. We provide theoretical justification of the Ising model and the score statistic. We also consider FR procedure based on fitting the Ising model using l_1 regularized logistic regressions. This version incorporates all features simultaneously, but it is computationally expensive for large number of labels. The experiments carried out on artificial and real data show that the proposed methods can outperform the conventional ones. Thus, they can be recommended, especially for datasets with moderate number of labels and large number of features.

Future work should include generalization of the proposed approach to more general Markov Networks. In particular one can consider a generalized Ising model in which some non-linear functions of features are used instead of linear combinations. The major problem associated with more general Markov Networks is how to estimate the parameters and how to test the significance of features efficiently.

In our procedure we test the significance of feature x_j in model $y_k \sim \mathbf{y}_{-k}, x_j$ using logistic regression. This suggests that other classification models can be used, e.g. decision trees which allow to discover non-linear dependencies. Note however that the main problem with decision trees (and some other classification models) is how to verify the significance of x_j in model $y_k \sim \mathbf{y}_{-k}, x_j$, possibly without refitting the model when adding x_j (as in our procedure). The modification of the score statistic $u_k(x_j)$ would be necessary for other models.

The other interesting question is how to choose the final subset of features having their ordering. Here we use a simple approach based on validation set. It would be worthwhile to have a more sophisticated method, which does not require separating a validation set.

The limitation of many FR methods is that the features are accessed individually and the possible redundancy as well as joint relevance of features is not taken into account. On the other hand the methods, which take into account all features simultaneously (e.g. *ising+ll*), are usually slow for large number of features or labels. It would be interesting to combine methods which assess the individual relevance of the features (like *ising+score*) with those taking into account all features simultaneously (like *ising+ll*). This could be done by applying two-step procedure in which *ising+score* is used at first to filter out least significant features and then *ising+ll* is launched on the remaining set of features.

Acknowledgements

I am grateful to the Associate Editor and anonymous reviewers for their valuable comments that helped to improve the initial version of this paper.

Research of Paweł Teisseyre was supported by the European Union from resources of the European Social Fund within project 'Information technologies: research and their interdisciplinary applications' POKL.04.01.01-00-051/10-00.

Appendix A.

Appendix A.1. Proof of (10)

For simplicity, we write x instead of x_j . Using the definition of conditional probabilities we can write

$$\frac{P(y_k = 1|x, \mathbf{y}_{-k})}{P(y_k = 0|x, \mathbf{y}_{-k})} = \frac{P(y_k = 1, x, \mathbf{y}_{-k})/P(x, \mathbf{y}_{-k})}{P(y_k = 0, x, \mathbf{y}_{-k})/P(x, \mathbf{y}_{-k})} = \frac{P(y_k = 1, x, \mathbf{y}_{-k})}{P(y_k = 0, x, \mathbf{y}_{-k})} = \frac{\exp[a_k x + \sum_{l:l \neq k} a_l x y_l + \sum_{s<l:s,l \neq k} \beta_{s,l} y_s y_l + \sum_{l:l \neq k} \beta_{l,k} y_l]}{\exp[\sum_{l:l \neq k} a_l x y_l + \sum_{s<l:s,l \neq k} \beta_{s,l} y_s y_l]} = \exp[a_k x + \sum_{l:l \neq k} \beta_{l,k} y_l],$$

which ends the proof.

Appendix A.2. Proof of Proposition 1

We can write

$$\begin{aligned} H_g(\mathbf{y}|x) &= - \sum_{\mathbf{y}} g(\mathbf{y}|x) \log(g(\mathbf{y}|x)) = - \sum_{\mathbf{y}} g(\mathbf{y}|x) \log \left[\frac{g(\mathbf{y}|x)}{p(\mathbf{y}|x)} p(\mathbf{y}|x) \right] = \\ &= -KL(g, p) - \sum_{\mathbf{y}} g(\mathbf{y}|x) \log(p(\mathbf{y}|x)) \leq - \sum_{\mathbf{y}} g(\mathbf{y}|x) \log(p(\mathbf{y}|x)), \end{aligned}$$

where $KL(g, p)$ is a Kullback-Leibner divergence between g and p and the last inequality follows from $KL(g, p) \geq 0$ (see e.g. Theorem 8.6.1 in [67]). Using the definition of p and the fact that both p and g must satisfy constraints (12) and (13), we obtain

$$\begin{aligned} - \sum_{\mathbf{y}} g(\mathbf{y}|x) \log(p(\mathbf{y}|x)) &= - \sum_{\mathbf{y}} g(\mathbf{y}|x) \left[-\log(Z(x)) + \sum_{k=1}^K a_k x y_k + \sum_{k<j} (\beta_{k,j} + b_{k,j} x) y_k y_j \right] = \\ &= - \sum_{\mathbf{y}} p(\mathbf{y}|x) \left[-\log(Z(x)) + \sum_{k=1}^K a_k x y_k + \sum_{k<j} (\beta_{k,j} + b_{k,j} x) y_k y_j \right] = - \sum_{\mathbf{y}} p(\mathbf{y}|x) \log(p(\mathbf{y}|x)), \end{aligned}$$

which ends the proof.

Appendix A.3. Auxiliary facts

Let us introduce some additional notation. In the following $\|\mathbf{w}\|$ will denote Euclidean norm of vector \mathbf{w} and $\|\mathbf{w}\|_\infty$ maximum norm. In addition $\lambda_j(\mathbf{A})$ denotes j -th eigenvalue of matrix \mathbf{A} , $\lambda_{\min}(\mathbf{A})$ ($\lambda_{\max}(\mathbf{A})$) its minimal (maximal) eigenvalue.

Let $l(\cdot)$ be log-likelihood function based on larger model $y_k \sim \mathbf{y}_{-k}, x$. Let $\mathbf{s}(\cdot)$ be gradient of $l(\cdot)$. Recall that $\mathbf{Z} = (\mathbf{Y}_{-k}, \mathbf{X}_j)$ is $n \times K$ matrix. It is easy to calculate that $\mathbf{s}(\boldsymbol{\theta}_k) = \mathbf{Z}^T (\mathbf{Y}_k - \mathbf{p}(\boldsymbol{\theta}_k))$, $\mathbf{p}(\boldsymbol{\theta}_k) = (p^{(1)}(\boldsymbol{\theta}_k), \dots, p^{(n)}(\boldsymbol{\theta}_k))^T$. Since coordinates of $\mathbf{s}(\boldsymbol{\theta}_k)$ are sums whose summands are bounded by L , it follows from Hoeffding inequality that

$$P(\|\mathbf{s}(\boldsymbol{\theta}_k)\|_\infty > \delta | \mathbf{Z}) \leq K \exp \left[-\frac{2\delta^2}{nL^2} \right], \quad (\text{A.1})$$

for any $\delta > 0$.

For logistic regression, Hessian matrix of $l(\cdot)$ is equal $-\mathbf{I}(\cdot)$, where $\mathbf{I}(\cdot) = \mathbf{Z}^T \mathbf{W}(\cdot) \mathbf{Z}$.

Lemma 1. Assume that $|(\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{Z}^{(i)}| \leq 1$, for some vector $\mathbf{w} \in R^K$. Then

$$p^{(i)}(\mathbf{w})(1 - p^{(i)}(\mathbf{w})) > e^{-3} p^{(i)}(\boldsymbol{\theta}_k)(1 - p^{(i)}(\boldsymbol{\theta}_k)).$$

Proof. Observe that for \mathbf{w} such that $|(\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{Z}^{(i)}| \leq 1$ we have

$$\frac{p^{(i)}(\mathbf{w})(1 - p^{(i)}(\mathbf{w}))}{p^{(i)}(\boldsymbol{\theta}_k)(1 - p^{(i)}(\boldsymbol{\theta}_k))} = e^{(\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{Z}^{(i)}} \left[\frac{1 + e^{\boldsymbol{\theta}_k^T \mathbf{Z}^{(i)}}}{1 + e^{\mathbf{w}^T \mathbf{Z}^{(i)}}} \right]^2 \geq e^{-1} \left[\frac{e^{-\boldsymbol{\theta}_k^T \mathbf{Z}^{(i)}} + 1}{e^{-\boldsymbol{\theta}_k^T \mathbf{Z}^{(i)}} + e} \right]^2 \geq e^{-3}.$$

□

Recall that $\hat{\boldsymbol{\theta}}_k$ is an estimator of $\boldsymbol{\theta}_k$ based on model $y_k \sim \mathbf{y}_{-k}$ in which the last coordinate corresponding to x_j is set to 0.

Lemma 2. The following inequality holds

$$P[l(\boldsymbol{\theta}_k) - l(\hat{\boldsymbol{\theta}}_k) > e^{-3} \Lambda_{\min} v n d^2 / 4 | \mathbf{Z}] \geq 1 - K \exp \left[-\frac{Cn(K + L^2)a_k^2}{2H^2} \right],$$

where $d = \frac{|a_k|}{\sqrt{K + L^2}H}$, $H = \max(1, G)$.

Proof. Define set $A = \{\mathbf{w} : \|\mathbf{w} - \boldsymbol{\theta}_k\| \leq d\}$ and observe that the last coordinate of $\hat{\boldsymbol{\theta}}_k$ is set to 0, thus $\hat{\boldsymbol{\theta}}_k \notin A$, as $d \leq |a_k|$. Define function $H(\mathbf{w}) := l(\boldsymbol{\theta}_k) - l(\mathbf{w})$ and observe that $H(\mathbf{w})$ is convex (as $l(\cdot)$ is concave) and $H(\boldsymbol{\theta}_k) = 0$. Thus it suffices to show that $H(\mathbf{w}) > e^{-3} \Lambda_{\min} v n d^2 / 4$ on the boundary of A , i.e. for \mathbf{w} such that $\|\mathbf{w} - \boldsymbol{\theta}_k\| = d$, with large probability.

Using Taylor expansion and the fact that Hessian matrix of $l(\cdot)$ is equal to $-\mathbf{I}(\cdot)$, we can write

$$H(\mathbf{w}) = -(\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{s}(\boldsymbol{\theta}_k) + (\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{I}(\bar{\mathbf{w}})(\mathbf{w} - \boldsymbol{\theta}_k)/2, \quad (\text{A.2})$$

where $\bar{\mathbf{w}}$ is some point in set A .

Using Cauchy-Schwarz inequality we have

$$|(\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{Z}^{(i)}| \leq \|(\mathbf{w} - \boldsymbol{\theta}_k)\| \cdot \|\mathbf{Z}^{(i)}\| \leq d \sqrt{K + L^2} = \frac{|a_k|}{\max(1, G)} \leq 1. \quad (\text{A.3})$$

It follows from (A.3), Lemma 1 and Assumption 2 that

$$(\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{I}(\bar{\mathbf{w}})(\mathbf{w} - \boldsymbol{\theta}_k)/2 \geq e^{-3} (\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{I}(\boldsymbol{\theta}_k)(\mathbf{w} - \boldsymbol{\theta}_k)/2 \geq e^{-3} d^2 \Lambda_{\min} v n / 2. \quad (\text{A.4})$$

Observe that

$$|(\mathbf{w} - \boldsymbol{\theta}_k)^T \mathbf{s}(\boldsymbol{\theta}_k)| \leq \sqrt{K} \|\mathbf{w} - \boldsymbol{\theta}_k\| \cdot \|\mathbf{s}(\boldsymbol{\theta}_k)\|_\infty \leq \sqrt{K + L^2} d \|\mathbf{s}(\boldsymbol{\theta}_k)\|_\infty. \quad (\text{A.5})$$

Now using (A.2), (A.4) and (A.5) we can write

$$P[H(\mathbf{w}) > e^{-3}\Lambda_{\min}vnd^2/4|\mathbf{Z}] \geq P[-d\|\mathbf{s}(\theta_k)\|_{\infty} \sqrt{K+L^2} + e^{-3}\Lambda_{\min}vnd^2/2 \geq e^{-3}\Lambda_{\min}vnd^2/4|\mathbf{Z}] \geq P\left[\|\mathbf{s}(\theta_k)\|_{\infty} \leq \frac{d\Lambda_{\min}vn}{4e^3\sqrt{K+L^2}}|\mathbf{Z}\right] \geq 1 - K \exp\left[-\frac{d^2\Lambda_{\min}^2v^2n}{8e^6(K+L^2)L^2}\right] = 1 - K \exp\left[-\frac{Cn(K+L^2)a_k^2}{2[\max(1, G)]^2}\right],$$

where the last inequality follows from (A.1). This ends the proof. \square

Recall that $v(\hat{\theta}_k) = D(\hat{\theta}_k) - \mathbf{C}(\hat{\theta}_k)\mathbf{A}^{-1}(\hat{\theta}_k)\mathbf{B}(\hat{\theta}_k)$, where $\mathbf{A}, \mathbf{B}, \mathbf{C}, D$ are defined in Section 3.1.

Lemma 3. *The following inequality holds*

$$v^{-1}(\hat{\theta}_k) \geq \frac{4}{(K+L^2)L^2n}.$$

Proof. First observe that using a definition of Shur complement (see e.g. [68], Section 3.4) we have that $v^{-1}(\hat{\theta}_k) = [\mathbf{I}^{-1}(\hat{\theta}_k)]_{K,K}$, where $[A]_{K,K}$ denotes element in K -th row and K -th column of matrix A . Observe that

$$[\mathbf{I}^{-1}(\hat{\theta}_k)]_{K,K} \geq \lambda_{\min}(\mathbf{I}^{-1}(\hat{\theta}_k)) = \frac{1}{\lambda_{\max}(\mathbf{I}(\hat{\theta}_k))}.$$

Since $p^{(i)}(\hat{\theta}_k)(1 - p^{(i)}(\hat{\theta}_k)) < 0.25$, each element on the diagonal of $\mathbf{I}(\hat{\theta}_k)$ is bounded by $L^2n/4$ and thus

$$\lambda_{\max}(\mathbf{I}(\hat{\theta}_k)) \leq \sum_{j=1}^K \lambda_j(\mathbf{I}(\hat{\theta}_k)) = \sum_{j=1}^K [\mathbf{I}(\hat{\theta}_k)]_{j,j} \leq K \max_j [\mathbf{I}(\hat{\theta}_k)]_{j,j} \leq \frac{KL^2n}{4} \leq \frac{(K+L^2)L^2n}{4},$$

which ends the proof. \square

Lemma 4. *Let $\Lambda_{\min} = \lambda_{\min}(\mathbf{Z}^T\mathbf{Z}/n) > 0$. Then function $l(\cdot)$ is concave.*

Proof. Note that $\Lambda_{\min} > 0$ implies positive definiteness of $\mathbf{Z}^T\mathbf{Z}$. We have to show that $\mathbf{I}(\mathbf{w})$ is positive definite for any $\mathbf{w} \in R^K$. For any vectors $\mathbf{w}, \mathbf{c} \in R^K$ we have

$$\min_i p^{(i)}(\mathbf{w})(1 - p^{(i)}(\mathbf{w}))\mathbf{c}^T\mathbf{Z}^T\mathbf{Z}\mathbf{c} \leq \mathbf{c}^T\mathbf{I}(\mathbf{w})\mathbf{c}.$$

Since $\min_i p^{(i)}(\mathbf{w})(1 - p^{(i)}(\mathbf{w})) > 0$, positive definiteness of $\mathbf{Z}^T\mathbf{Z}$ implies positive definiteness of $\mathbf{I}(\mathbf{w})$, for any \mathbf{w} , which ends the proof. \square

Appendix A.4. Proof of Theorem 1

Using Taylor expansion of log-likelihood function we obtain

$$l(\theta_k) = l(\hat{\theta}_k) + (\theta_k - \hat{\theta}_k)^T \mathbf{s}(\hat{\theta}_k) - (\theta_k - \hat{\theta}_k)^T \mathbf{I}(\bar{\theta}_k)(\theta_k - \hat{\theta}_k)/2, \quad (\text{A.6})$$

where $\bar{\theta}_k$ is point on the line segment between θ_k and $\hat{\theta}_k$. Note that the first $K-1$ coordinates of $\mathbf{s}(\hat{\theta}_k)$ are equal zero and thus (A.6) reduces to

$$l(\theta_k) - l(\hat{\theta}_k) = a_k s(\hat{\theta}_k) - (\theta_k - \hat{\theta}_k)^T \mathbf{I}(\bar{\theta}_k)(\theta_k - \hat{\theta}_k)/2. \quad (\text{A.7})$$

Now from (A.7), non-negativity of $(\theta_k - \hat{\theta}_k)^T \mathbf{I}(\bar{\theta}_k)(\theta_k - \hat{\theta}_k)/2$ and Lemma 2 we have

$$|s(\hat{\theta}_k)| \geq \frac{\Lambda_{\min}vnd^2}{4e^3a_k} = \frac{\Lambda_{\min}vn|a_k|}{4e^3(K+L^2)[\max(1, G)]^2}, \quad (\text{A.8})$$

with probability given in Lemma 2. The assertion of the Theorem follows directly from (A.8) and Lemma 3.

References

- [1] R. E. Schapire, Y. Singer, Boostexter: A boosting-based system for text categorization, *Machine Learning* 39 (2-3) (2000) 135–168.
- [2] I. Katakis, G. Tsoumakas, I. Vlahavas, Multilabel text classification for automated tag suggestion, in: *Proceedings of the ECML/PKDD 2008 Discovery Challenge*, 2008.
- [3] C. D. Nguyen, T. A. Dung, T. H. Cao, Text classification for dag-structured categories, in: *Proceedings of the 9th Pacific-Asia Conference on Advances in Knowledge Discovery and Data Mining, PAKDD'05*, Springer-Verlag, Berlin, Heidelberg, 2005, pp. 290–300.
- [4] M. E. Loza, J. Fürnkranz, Efficient pairwise multilabel classification for large-scale problems in the legal domain, in: *Machine Learning and Knowledge Discovery in Databases*, Vol. 5212 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2008, pp. 50–65.
- [5] T. N. Rubin, A. Chambers, P. Smyth, M. Steyvers, Statistical topic models for multi-label document classification, *Machine Learning* 88 (1-2) (2012) 157–208.
- [6] M. Wang, X. Zhou, T.-S. Chua, Automatic image annotation via local multi-label classification, in: *Proceedings of the 2008 International Conference on Content-based Image and Video Retrieval, CIVR '08*, ACM, New York, NY, USA, 2008, pp. 17–26.
- [7] Textonboost for image understanding: Multi-class object recognition and segmentation by jointly modeling texture, layout, and context, *International Journal of Computer Vision* 81 (1) (2009) 2–23.
- [8] N. Kumar, A. C. Berg, P. N. Belhumeur, S. K. Nayar, Attribute and simile classifiers for face verification, in: *In IEEE International Conference on Computer Vision (ICCV)*, 2009.
- [9] M. R. Boutell, J. Luo, X. Shen, C. M. Brown, Learning multi-label scene classification, *Pattern Recognition* 37 (9) (2004) 1757–1771.
- [10] J. Wang, Y. Zhao, X. Wu, X.-S. Hua, A transductive multi-label learning approach for video concept detection, *Pattern Recognition* 44 (10-11) (2011) 2274–2286.
- [11] K. Trohidis, G. Tsoumakas, G. Kalliris, I. Vlahavas, Multilabel classification of music into emotions, in: *International Conference on Music Information Retrieval*, 2008, pp. 325–330.
- [12] A. Elisseeff, J. Weston, A kernel method for multi-labelled classification, in: *In Advances in Neural Information Processing Systems 14*, MIT Press, 2001, pp. 681–687.
- [13] S. Diplaris, G. Tsoumakas, P. Mitkas, I. Vlahavas, Protein classification with multiple algorithms, in: *10th Panhellenic Conference on Informatics*, 2005, pp. 448–456.
- [14] Z. Barutcuoglu, R. E. Schapire, O. G. Troyanskaya, Hierarchical multi-label prediction of gene function, *Bioinformatics* 22 (7) (2006) 830–836.
- [15] M. Lappenschaar, A. Hommersom, J. Lagro, P. Lucas, Understanding the co-occurrence of diseases using structure learning, in: *Artificial Intelligence in Medicine*, Vol. 7885 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2013, pp. 135–144.
- [16] Q. Abbas, M. E. Celebi, C. Serrano, I. Fondón GarcíA, G. Ma, Pattern classification of dermoscopy images: A perceptually uniform model, *Pattern Recognition* 46 (1) (2013) 86–97.
- [17] K. Kawai, Y. Takahashi, Identification of the dual action antihypertensive drugs using tfs-based support vector machines, *Chem-Bio Informatics Journal* 4 (2009) 44–51.
- [18] M. A. Mammadov, A. M. Rubinov, J. Yearwood, The study of drug-reaction relationships using global optimization techniques, *Optimization Methods Software* 22 (1) (2007) 99–126.
- [19] L. Tang, H. Liu, Relational learning via latent social dimensions, in: *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '09*, ACM, New York, NY, USA, 2009, pp. 817–826.
- [20] S. Peters, L. Denoyer, P. Gallinari, Iterative annotation of multi-relational social networks, in: *Proceedings of the 2010 International Conference on Advances in Social Networks Analysis and Mining, ASONAM'10*, IEEE Computer Society, 2010, pp. 96–103.
- [21] Y. Zhang, S. Burer, W. N. Street, Ensemble pruning via semi-definite programming, *Journal of Machine Learning Research* 7 (2006) 1315–1338.
- [22] E. Gibaja, S. Ventura, A tutorial on multilabel learning, *ACM Computing Surveys* 47 (3) (2015) 1–38.
- [23] G. Tsoumakas, I. Katakis, Multilabel classification: an overview, *International Journal of Data Warehouse and Mining* 3 (2007) 1–13.
- [24] K. Dembczyński, W. Waegeman, W. Cheng, E. Hüllermeier, On label dependence and loss minimization in multi-label classification, *Machine Learning* 88 (2012) 5–45.
- [25] G. Madjarov, D. Kocev, D. Gjorgjevikj, S. Džeroski, An extensive experimental comparison of methods for multi-label learning, *Pattern Recognition* 45 (9) (2012) 3084–3104.
- [26] M. Zhang, Z. Zhou, A review on multi-label learning algorithms, *IEEE Transactions on Knowledge and Data Engineering* 26 (2013) 1819 – 1837.
- [27] T. Hastie, R. Tibshirani, J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference and Prediction*, Springer, 2009. URL <http://www-stat.stanford.edu/~tibs/ElemStatLearn/>
- [28] S. Bromuri, D. Zufferey, J. Hennebert, M. Schumacher, Multi-label classification of chronically ill patients with bag of words and supervised dimensionality reduction algorithms, *Journal of Biomedical Informatics* 51 (2014) 165–175.
- [29] J. Fan, J. Lv, Sure independence screening for ultra-high dimensional feature space (with discussion), *Journal of the Royal Statistical Society B* 70 (2008) 849–911.
- [30] H. Peng, F. L., C. Ding, Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 27 (8) (2005) 1226–1238.
- [31] J. Fan, R. Samworth, Y. Wu, Ultrahigh dimensional feature selection: Beyond the linear model, *J. Mach. Learn. Res.* 10 (2009) 2013–2038.
- [32] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*, Springer-Verlag New York, 2006.
- [33] J. Cheng, E. Levina, P. Wang, J. Zhu, A sparse Ising model with covariates, *Biometrics* 70 (2014) 943–953.
- [34] W. Bian, B. Xie, D. Tao, Corlog: Correlated logistic models for joint prediction of multiple labels, in: *JMLR Proceedings*, Vol. 22, 2012, pp. 109–117.
- [35] E. Ising, Beitrag zur theorie des ferromagnetismus, *Zeitschrift für Physik* 31 (1925) 253–258.
- [36] W. Lenz, Beiträge zum verständnis der magnetischen eigenschaften in festen körpern, *Physikalische Zeitschrift* 21 (1920) 613–615.

- [37] R. C. Rao, Large sample tests of statistical hypotheses concerning several parameters with applications to problems of estimation, *Mathematical Proceedings of the Cambridge Philosophical Society* 44 (1948) 50–57.
- [38] W. Chen, J. Yan, B. Zhang, Z. Chen, Q. Yang, Document transformation for multi-label feature selection in text categorization, in: *Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on Data Mining, 2007*, pp. 451–456.
- [39] N. Spolaôr, E. A. Cherman, M. C. Monard, H. D. Lee, A comparison of multi-label feature selection methods using the problem transformation approach, *Electronic Notes in Theoretical Computer Science* 292 (2013) 135 – 151.
- [40] G. Doquire, M. Verleysen, Mutual information-based feature selection for multilabel classification, *Neurocomputing* 122 (2013) 148 – 155.
- [41] L. F. Kozachenko, N. N. Leonenko, Sample estimate of the entropy of a random vector, *Problems of Information Transmission*, 1987, 23:2, 95101 23 (1987) 9–16.
- [42] J. Lee, D.-W. Kim, Feature selection for multi-label classification using multivariate mutual information, *Pattern Recognition Letters* 34 (3) (2013) 349–357.
- [43] J. Read, A pruned problem transformation method for multi-label classification, in: *In Proc. 2008 New Zealand Computer Science Research Student Conference (NZCSRS), 2008*, pp. 143–150.
- [44] A. Clare, R. King, Knowledge discovery in multi-label phenotype data, in: L. De Raedt, A. Siebes (Eds.), *Principles of Data Mining and Knowledge Discovery*, Vol. 2168 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2001, pp. 42–53.
- [45] P. Ravikumar, M. Wainwright, J. Lafferty, High-dimensional Ising model selection using ℓ_1 -regularized logistic regression, *Annals of Statistics* 38 (2010) 1287–1319.
- [46] I. T. Jolliffe, A note on the use of principal components in regression, *Journal of the Royal Statistical Society. Series C (Applied Statistics)* 31 (3) (1982) 300–303.
- [47] H. Martens, Reliable and relevant modelling of real world data: a personal account of the development of PLS regression, *Chemometrics and Intelligent Laboratory Systems* 58 (2) (2001) 85 – 95.
- [48] S. Wold, Personal memories of the early PLS development, *Chemometrics and Intelligent Laboratory Systems* 58 (2) (2001) 83–84.
- [49] Penalized partial least square discriminant analysis with for multi-label data, *Pattern Recognition* 48 (5) (2015) 1724 – 1733.
- [50] L. Sun, S. Ji, J. Ye, Canonical correlation analysis for multilabel classification: A least-squares formulation, extensions, and analysis, *Pattern Analysis and Machine Intelligence, IEEE Transactions on* 33 (1) (2011) 194–200.
- [51] L. Sun, S. Ji, J. Ye, *Multi-Label Dimensionality Reduction*, Chapman and Hall/CRC, London, 2013.
- [52] J. E. Besag, Nearest-neighbour systems and the auto-logistic model for binary data, *Journal of the Royal Statistical Society. Series B (Methodological)* 34 (1) (1972) 75–83.
- [53] M. Zalewska, W. Niemiro, B. Samoliński, MCMC imputation in autologistic model, *Monte Carlo Methods and Applications* 16 (2010) 421–438.
- [54] J. D. Lafferty, A. MacCallum, F. C. N. Pereira, Conditional random fields: Probabilistic models for segmenting and labeling sequence data, in: *Proceedings of the Eighteenth International Conference on Machine Learning*, Vol. 22, 2001, pp. 282–289.
- [55] L. Fahrmeir, Asymptotic testing theory for generalized linear models, *Statistics* 1 (1987) 65–76.
- [56] Q. He, D. Lin, A variable selection method for genome-wide association studies, *Bioinformatics* 27 (1) (2011) 1–8.
- [57] E. T. Jaynes, Information theory and statistical mechanics, *Physical Review* 106 (1957) 620–630.
- [58] J. Read, B. Pfahringer, G. Holmes, E. Frank, Classifier chains for multi-label classification, *Machine Learning* 85 (2011) 333–359.
- [59] K. Dembczyński, W. Cheng, E. Hüllermeier, Bayes optimal multilabel classification via probabilistic classifier chains, in: *Proceedings of the twenty-seventh international conference on machine learning*, Vol. 22, 2010, pp. 109–117.
- [60] R. Holte, Very simple classification rules perform well on most commonly used datasets, *Machine Learning* 11 (1993) 63–91.
- [61] P. Przybyła, P. Teisseyre, What do your look-alikes say about you? Exploiting strong and weak similarities for author profiling, in: *Notebook for PAN at CLEF, 2015*.
- [62] P. Teisseyre, Asymptotic consistency and order specification for logistic classifier chains in multi-label learning, *Under review* (2015).
- [63] J. Demšar, Statistical comparisons of classifiers over multiple data sets, *J. Mach. Learn. Res.* 7 (2006) 1–30.
- [64] W. J. Conover, *Practical nonparametric statistics*, Wiley, New York, 1980.
- [65] J. Friedman, T. Hastie, R. Tibshirani, Regularization paths for generalized linear models via coordinate descent, *Journal of Statistical Software* 22.
- [66] P. Romanski, L. Kotthoff, FSelector: Selecting attributes, R package version 0.20 (2014).
URL <http://CRAN.R-project.org/package=FSelector>
- [67] T. M. Cover, J. A. Thomas, *Elements of Information Theory* (Wiley Series in Telecommunications and Signal Processing), Wiley-Interscience, 2006.
- [68] J. E. Gentle, *Matrix Algebra: Theory, Computations, and Applications in Statistics*, Springer, New York, 2007.