

Minimum Description Length Penalization for Group and Multi-Task Sparse Learning

Paramveer S. Dhillon

*Department of Computer and Information Science
University of Pennsylvania
Philadelphia, PA 19104, USA*

DHILLON@CIS.UPENN.EDU

Dean P. Foster

*Department of Statistics
Wharton School, University of Pennsylvania
Philadelphia, PA 19104, USA*

FOSTER@WHARTON.UPENN.EDU

Lyle H. Ungar

*Department of Computer and Information Science
University of Pennsylvania
Philadelphia, PA 19104, USA*

UNGAR@CIS.UPENN.EDU

Editor: Francis Bach

Abstract

We propose a framework MIC (Multiple Inclusion Criterion) for learning sparse models based on the information theoretic Minimum Description Length (MDL) principle. MIC provides an elegant way of incorporating arbitrary sparsity patterns in the feature space by using two-part MDL coding schemes. We present MIC based models for the problems of grouped feature selection (MIC-GROUP) and multi-task feature selection (MIC-MULTI). MIC-GROUP assumes that the features are divided into groups and induces two level sparsity, selecting a subset of the feature groups, and also selecting features within each selected group. MIC-MULTI applies when there are multiple related tasks that share the same set of potentially predictive features. It also induces two level sparsity, selecting a subset of the features, and then selecting which of the tasks each feature should be added to. Lastly, we propose a model, TRANSFEAT, that can be used to transfer knowledge from a set of previously learned tasks to a new task that is expected to share similar features. All three methods are designed for selecting a small set of predictive features from a large pool of candidate features. We demonstrate the effectiveness of our approach with experimental results on data from genomics and from word sense disambiguation problems.¹

Keywords: feature selection, minimum description length principle, multi-task learning

1. Introduction

Classical supervised learning algorithms use a set of feature-label pairs to learn mappings from the features to the associated labels. They generally do this by considering each classification task (each possible label) in isolation and learning a model for that task. Learning models independently for different tasks often works well, but when the labeled data is limited and expensive to obtain, an attractive alternative is to build shared models for multiple related tasks (Caruana, 1997; Ando

1. Preliminary versions of this work appeared in Dhillon et al. (2008, 2009); Dhillon and Ungar (2009).

and Zhang, 2005). For example, when one is trying to predict a set of related responses (“tasks”), be they multiple clinical outcomes for patients or growth rates for yeast strains under different conditions, it may be possible to “borrow strength” by sharing information between the models for the different responses. Inductive transfer by building shared models can also be valuable when we have a disproportionate amount of labeled data for “similar” tasks. In such a case, building separate models for each task often gives poor predictive accuracies on tasks which have little data.

As a running example, we consider the problem of disambiguating word senses based on their context. Here, each observation of a word (e.g., a sentence containing the word “fire”) is associated with multiple labels corresponding to each of the different possible meanings (e.g., for firing a person, firing a gun, firing off a note, etc.). Given the high-dimensional nature of Word Sense Disambiguation (WSD) data, feature selection is important for both linguistic understanding and for effective prediction (Chen et al., 2006). Also, since the features that are useful for predicting one sense are likely to be useful for predicting the other senses (perhaps with a coefficient of different sign.), we propose to select features that are useful in predicting these multiple responses.

Another closely related problem is grouped feature selection; that is, enforcing sparsity at the level of groups (feature classes) (Yuan and Lin, 2006; Bach et al., 2004; Dhillon et al., 2008). In this problem the group structure is over the features rather than over the tasks. Multi-task learning (described above) can also be thought of as a special case of this “group sparsity” scenario in which a group is defined by fixing a specific feature and ranging over multiple tasks. The block-norm approach to these problems uses a combination of ℓ_1 and ℓ_2 norms as regularization terms and adds each feature into the models of either none or all of the tasks (Obozinski et al., 2009) for the multi-task case and selects either none or all the features from a given group in the case of group sparsity. However, if the regularization constant is chosen using cross-validation, these approaches overestimate the support (Wainwright, 2009) as they select more features than the correct set of sparse features that generated the data. Wainwright (2009), Liu and Zhang (2008) and Nardi and Rinaldo (2008) have showed that certain scalings of the regularization coefficient yields more sparse solutions, which have with high probability the same support as the model generating the data. Even then there are further problems with these methods; in order to obtain very sparse solutions, one has to use a large regularization parameter that leads to suboptimal prediction accuracy because this high penalty not only shrinks irrelevant features to zero, but also shrinks relevant features to zero (Zhang, 2009a). Another alternative is to threshold the obtained coefficients (Lounici, 2008), but this introduces another thresholding parameter which needs to be tuned.

Motivated by the aforementioned reasons and by recent theoretical results on ℓ_0 penalty based regularization (Zhang, 2009a,b), we consider ℓ_0 penalty based formulations in this paper. In particular we propose to solve these two related problems, simultaneous feature selection for a set of multiple related tasks and grouped feature selection for a single task, by using coding schemes inspired by the Minimum Description Length (MDL) principle. We propose a common framework for these problems which we call the Multiple Inclusion Criterion (MIC). We use a “two part” version of MDL (Grünwald, 2005) to define a cost function which is greedily minimized by our methods. Since the greedy feature selection approximates the ℓ_0 penalty, we achieve a high degree of sparsity as is desired for both scientific interpretability and for accurate prediction in domains like Genomics and Word Sense Disambiguation (WSD) which have very high dimensional data. More importantly, our methods achieve *two-level sparsity*. In multi-task learning, each feature is added into models of a (possibly empty) subset of the tasks and in group feature selection, a (possibly empty) subset of the features are selected from each group (feature class).

We also propose a similarly motivated model (TRANSFEAT) (Dhillon and Ungar, 2009) for “intra-domain” adaptation which can be used to transfer knowledge from a set of already learned tasks to a new task which is similar to the aforementioned tasks. As an example, consider the task of predicting whether a word has a given sense when one already has models for predicting senses for synonyms of that word. These models are likely to share many of the same features; that is, a model for disambiguating one sense of “discharge” is likely to use many of the same features as one for disambiguating the sense of “fire” which is its synonym. Unlike MIC where we do simultaneous feature selection, the sharing in this case takes the form of a prior. TRANSFEAT is most beneficial when the word under consideration has considerably less labeled data available than the synonyms of that word (for example) so that building a supervised learning model for that word alone does not yield high predictive accuracy.

The rest of the paper is organized as follows. In the next section, we review relevant previous work. In Section 3, we provide background on feature selection and the MDL principle. Then in Section 4 we develop the general framework used by our models and describe the MIC-MULTI and MIC-GROUP models in detail. In Section 5, we show experimental results on real and synthetic data. In Section 6, we provide some model consistency results for the MIC models. In Section 7, we discuss the TRANSFEAT model and show its effectiveness for intra-domain adaptation on real world data sets. We conclude in Section 8 with a brief summary.

2. Related Work

The main contribution of this paper is to propose a joint framework for the related tasks of simultaneous feature selection for multiple related tasks and grouped feature selection for a single task. We are not aware of any previous work that addresses these two problems together, though Obozinski et al. (2009) do mention that these two problems are related. Nonetheless, there has been much previous work on each of these problems separately.

Jebara (2004) uses maximum-entropy discrimination to select a single subset of features across multiple SVM regression or classification problems that share a common set of potential features. Several other papers work within the framework of regularized regression, taking the penalty term to be an ℓ_1 norm over features of an ℓ_q norm over the coefficients for each feature (an “ $\ell_1 - \ell_q$ ” penalty). Turlach et al. (2005) consider the case $q = \infty$, while Argyriou et al. (2008) and Obozinski et al. (2009) use $q = 2$. Argyriou et al. (2008) show that the general subspace selection problem can be formulated as an optimization problem involving the trace norm. Obozinski et al. (2009) propose BBLASSO, which focuses on the case where the trace norm is not required; they instead use a homotopy-based approach to evaluate the entire regularization path efficiently (Efron et al., 2004). Ando and Zhang (2005) also propose a framework which uses multiple prediction problems to learn an underlying shared structural parameter on the input (feature) space and they penalize the weight vectors by ℓ_2 norm. The idea behind $\ell_1 - \ell_q$ penalties is that when $q > 1$, the cost of making a coefficient nonzero is smaller for features that are shared across more tasks. Indeed, for either $q = 2$ or $q = \infty$, these algorithms tend in practice to yield nonzero coefficients for all of the tasks associated with features that get selected.

The related problem of grouped feature selection for a single task has also been addressed previously by Yuan and Lin (2006), Bach et al. (2004), Meier et al. (2008), Zhao et al. (2008) and is known as “Group Lasso”. It is an extension of Lasso (ℓ_1 penalty) to the case of grouped structure in data and it enforces sparsity at the level of groups, that is, an entire group of features is selected. It

penalizes a (ℓ_1/ℓ_2) norm of the feature weights. An alternative formulation of Group Lasso is called Multiple Kernel Learning (MKL) (Bach et al., 2004; Bach, 2008); it penalizes the kernel Hilbert norm instead of the Euclidean norm.

Our approach is different from these methods in that we use ℓ_0 penalty-based greedy feature selection methods which minimize a cost function provided by MDL based coding schemes. MDL-based coding schemes provide much flexibility to incorporate arbitrary sparsity structures in the problem at hand. Recently, Huang et al. (2009) have also used coding schemes similar to the MDL for enforcing arbitrary structured sparsity patterns over the feature space.

3. Background

We assume a setting in which we are given n labeled data samples as $\{(x_i, y_i)\}_{i=1}^n \in \mathcal{X} \times \mathcal{Y}$ where $\mathcal{X} \in \mathbb{R}^p$ (the feature vector lives in a p dimensional space) and our goal is to find the parameter vector ($\mathbf{w} \in \mathbb{R}^p$) of a statistical model fit to the above data. Alternatively, we can represent the data and the response variables in matrix form as $\mathbf{X}_{n \times p}$ and $\mathbf{Y}_{n \times 1}$, respectively, and the p dimensional weight vector as $\mathbf{w}_{p \times 1}$. Standard linear or logistic regression models of the form $\mathbf{Y} = \mathbf{w} \cdot \mathbf{X}$ (or $P(Y = 1|X) = \frac{1}{1+e^{-\mathbf{w} \cdot \mathbf{X}}}$) fail to estimate the weight vector \mathbf{w} in the case in which $p > n$ as they require inversion of a rank deficient matrix. To overcome this problem, *regularized* versions of the linear or logistic regressions are used which penalize some norm of the weight vectors:

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \{ \|\mathbf{Y} - \mathbf{X} \cdot \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_q^q \},$$

where $\|\mathbf{w}\|_q$ represents the ℓ_q norm of \mathbf{w} and λ is a hyperparameter.

For $q = 2$, the penalized regression is known as *Ridge Regression*, which corresponds to a Bayesian maximum a posteriori estimate for w under a Gaussian prior and shrinks the weight vector but does not enforce sparsity. The ℓ_1 penalty (Lasso) is equivalent to a double exponential prior on w (Tibshirani, 1996) and enforces sparsity by driving some of the weights to zero. As q approaches 0, $\|\mathbf{w}\|_q^q$ approaches the number of non-zero values in w . Hence regularization with ℓ_0 penalty is subset selection: Choosing a small number of the original features to retain in the model. Once a coefficient is in the model, all that counts is the cost of adding it in the first place. The ℓ_0 penalty has a number of advantages, including bounded worst case risk with respect to the ℓ_1 penalty and better control of False Discovery Rate (FDR) (Lin et al., 2008). There are other problems with the ℓ_1 penalty other than being less sparse as mentioned earlier, namely that its sparsity is not explicitly controlled, and in order to obtain very sparse solutions, one has to use a large regularization parameter that leads to suboptimal prediction accuracy because this high penalty not only shrinks irrelevant features to zero, but also shrinks relevant features to zero (Zhang, 2009a). However, one virtue of the ℓ_1 penalty is computational tractability (Efron et al., 2004), in contrast to the ℓ_0 penalty, which requires subset search which is (worst case) NP-Hard (Natarajan, 1995). In practice, approximate greedy algorithms like forward stepwise feature selection yield accurate, highly sparse solutions.

In a regression model, the residual sum of squares is proportional up to an additive constant to the negative log-likelihood of \mathbf{Y} given \mathbf{X} (Bickel and Doksum, 2001). Thus, the ℓ_0 regularization can be rephrased as a penalized likelihood criterion as follows:

$$\text{score} = -2 \log P(\mathbf{Y} | \hat{w}_q) + F \cdot q,$$

where q is the number of features in the model, $P(\mathbf{Y}|\hat{\mathbf{w}}_q)$ is the likelihood of the data given a model containing q features and F is a free parameter that controls the amount of weight placed on the ℓ_0 norm. Various penalties have been proposed for F , including

- $F = 2$, corresponding approximately to the AIC (Akaike Information Criterion) (Akaike, 1973),
- $F = \log n$, giving the BIC (Bayesian Information Criterion) (Schwartz, 1978),
- $F = 2 \log p$, giving to RIC (Risk Inflation Criterion—similar to a “Bonferroni correction”) (Foster and George, 1994).

As discussed in next subsection, each of these penalties can also be derived by using the Minimum Description Length (MDL) principle under different coding schemes.

3.1 Minimum Description Length (MDL) Principle Preliminaries

MDL (Rissanen, 1978, 1999) is a principle for model selection which treats the best model as the one which maximally compresses a digital representation of the observed data. We can envision a “Sender” who wants to transmit some data to a “Receiver” using as few bits as possible. For an illustrative example of the MDL principle, consider the case of simple linear regression. Assume that both the Sender and Receiver know the $n \times p$ data matrix \mathbf{X} , and the Sender wants to convey the values in the $n \times 1$ response matrix \mathbf{Y} . The naïve way to do this would be to send the raw values for each of the n observations of \mathbf{Y} . However, a more efficient way to send this information would be to describe a regression model $\hat{\mathbf{w}}$ for \mathbf{Y} given \mathbf{X} and then to send the residuals $\mathbf{Y} - \mathbf{X} \cdot \hat{\mathbf{w}}$, which have a much narrower distribution and would require fewer bits to encode.

To minimize description length, then, Sender should choose $\hat{\mathbf{w}}^*$ such that

$$\hat{\mathbf{w}}^* = \underset{\hat{\mathbf{w}}}{\operatorname{argmin}}\{\mathcal{D}(\mathbf{Y}|\hat{\mathbf{w}}) + \mathcal{D}(\hat{\mathbf{w}})\}, \quad (1)$$

where the first term is the description length of the residuals about the model, and the second term is the description length of the model itself. In other words, the first term represents the fit of the model to data; as the model fits better this term shrinks. The second term represents the complexity of the model; it grows as the model becomes more complex.

This version of the MDL principle is known as “*Two part MDL*” (Grünwald, 2005); the exact meaning of both these terms is described in the following sections.

In the next sections we show how to choose the two quantities in Equation 1 by incorporating the “structure” (characteristics) of the problem at hand. We then introduce search algorithms to approximately find a sparse set of weights $\hat{\mathbf{w}}$ that minimize Equation 1.

3.1.1 CODING THE DATA: $\mathcal{D}(\mathbf{Y}|\hat{\mathbf{w}})$

The Kraft inequality in information theory (Cover and Thomas, 2006) implies that for any probability distribution $\{p_i\}$ over a finite or countable set, there exists a corresponding code with codeword length $\lceil -\lg p_i \rceil$ (The logarithm is base 2). Moreover, these code lengths are optimal in the sense of minimizing the expected code length with respect to $\{p_i\}$. Also, if the Sender and Receiver agree on a model (e.g., linear regression), then they have a probability distribution over the residuals ϵ , so they will agree to use a code for the residuals with length:

$$\mathcal{D}(\mathbf{Y}|\hat{\mathbf{w}}) = -\lg P(\boldsymbol{\varepsilon}|\hat{\mathbf{w}}) = -\lg P(\mathbf{Y}|\hat{\mathbf{w}}), \quad (2)$$

that is, the negative log-likelihood of the data given the model. We dropped the ceiling on $-\lg P(\mathbf{Y}|\hat{\mathbf{w}})$ since we use “idealized” code lengths (Barron et al., 1998).

Consider a forward stepwise-regression setting in which we have already added $q - 1$ features to our model (including the intercept term), and we are deciding whether to include an extra q^{th} feature. Let Y_i denote the i^{th} row of \mathbf{Y} and \hat{w}_q , a linear regression model with all q features, then,²:

$$\begin{aligned} \mathcal{D}(\mathbf{Y}|\hat{\mathbf{w}}) &= -\lg \prod_{i=1}^n P(Y_i|\hat{w}_q) \\ &= -\sum_{i=1}^n \lg \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{1}{2\sigma^2} (Y_i - \mathbf{X}_i \cdot \hat{w}_q)^2 \right) \right] \\ &= \frac{1}{2\ln 2} \left[n \ln(2\pi\sigma^2) + \frac{(\mathbf{Y} - \mathbf{X} \cdot \hat{w}_q)^2}{\sigma^2} \right], \end{aligned}$$

σ^2 is unknown in practice, but it can be estimated as:

$$\hat{\sigma}^2 = \frac{(\mathbf{Y} - \mathbf{X} \cdot \hat{w}_{q-1})^2}{n}.$$

Note that this is the ML (Maximum Likelihood) estimate for σ^2 which Sender uses, as ignoring the model-coding cost, maximizing likelihood is equivalent to minimizing description length. Some statisticians, in practice, use the unbiased estimate $\hat{\sigma}^2 = \frac{(\mathbf{Y} - \mathbf{X} \cdot \hat{w}_{q-1})^2}{n-q}$. In the experiments presented in this paper, we estimate $\hat{\sigma}^2$ without the current q^{th} feature in model, in order to prevent overfitting.

We can write the final expression for $\mathcal{D}(\mathbf{Y}|\hat{\mathbf{w}})$, incorporating $\hat{\sigma}^2$ as:

$$\mathcal{D}(\mathbf{Y}|\hat{\mathbf{w}}) = \frac{n}{2\ln 2} \left[\ln \left(\frac{2\pi \times (\mathbf{Y} - \mathbf{X} \cdot \hat{w}_{q-1})^2}{n} \right) + \left(\frac{\mathbf{Y} - \mathbf{X} \cdot \hat{w}_q}{\mathbf{Y} - \mathbf{X} \cdot \hat{w}_{q-1}} \right)^2 \right].$$

3.1.2 CODING THE MODEL: $\mathcal{D}(\hat{\mathbf{w}})$

Just as $\mathcal{D}(\mathbf{Y}|\hat{\mathbf{w}})$ depends on the model for the residuals that Sender and Receiver choose, so their coding scheme for $\hat{\mathbf{w}}$ itself will reflect their prior expectations.³ When the number of features p is large (say, > 1000), Sender will likely only want to transmit a few of them that are most relevant, and hence the $\hat{\mathbf{w}}$ will contain mostly zeros. So, the first step in coding $\hat{\mathbf{w}}$ could be to say where the non-zero entries are located; if only a few features enter the model, this can be done relatively efficiently by listing the indices of the features in the set $\{1, 2, \dots, p\}$. This requires $\lceil \lg p \rceil$ bits or approximately $\lg p$ bits.

The second step is to encode the numerical values of those coefficients. Rissanen (1983) suggested the basic approach for doing this by creating a discrete grid over some possible parameter

2. Note that in the following notation the square of vector, that is, \mathbf{Z}^2 means $\mathbf{Z}^T \mathbf{Z}$.

3. By the Kraft inequality, we can interpret $2^{-\mathcal{D}(\hat{\mathbf{w}})}$ as a prior over possible models \mathbf{w} . In fact, this is done explicitly in the Minimum Message Length (MML) principle which is a Bayesian analogue of MDL, which chooses the model $\hat{\mathbf{w}}$ with maximum $P(\mathbf{w}|\mathbf{Y})$, that is, it chooses the model that minimizes $-\lg P(\mathbf{w}|\mathbf{Y}) = -\lg P(\mathbf{w}) - \lg P(\mathbf{Y}|\mathbf{w}) + \text{const}$.

values, and use a code for integers to specify which grid point is closest. A simple way to approximate the value of a particular coefficient \hat{w} is to encode an integer version of its z-score relative to the null-hypothesis value w_0 (which in our case is 0):

$$\left\langle \frac{\hat{w} - w_0}{SE(\hat{w})} \right\rangle = \left\langle \frac{\hat{w}}{SE(\hat{w})} \right\rangle,$$

where $\langle x \rangle$ means the closest integer to x and SE represents standard error. The z-score can be coded with the idealized universal code for the positive integers of Rissanen (1983), in which the cost to code $i \in 1, 2, 3, \dots$ is

$$\lg^* i + b,$$

where $\lg^* i = \lg i + \lg \lg i + \lg \lg \lg i + \dots$ so long as the terms remain positive, and $b \approx \lg 2.865 \approx 1.516$ is the constant such that

$$\sum_{i=1}^{\infty} 2^{-(\lg^* i + b)} = 1.$$

We require the \lg^* instead of a simple \lg because the number of bits Sender uses to convey the integer i will vary, and she needs to tell the Receiver how many bits to expect. The number of bits is itself an integer than can be coded, hence the iteration of logarithms.

In fact, in practice it is unnecessary to allow our integer code to extend to arbitrarily large integers. We are interested in features near the limit of detectability and we expect our z-scores to be roughly in the range ~ 2 to ~ 4 , since if they were much higher, the true features would be obvious and would not require sensitive feature selection. We could thus impose some maximum possible z-score Z that we might ever want to encode (say, 1000) and assume that all of our z-scores will fall below it. In this case, the constant c can be reduced to a new value c_Z , now only being large enough that,

$$\sum_{i=1}^Z 2^{(-\lg^* i + c_Z)} = 1. \quad (3)$$

In particular $c_{1000} \approx 1.199$. In our implementation in this paper, we avoid computing the actual values of our z-scores and instead assume a constant 2 bits per coefficient. The reason behind choosing 2 bits over using a more conservative penalty like BIC (Bayesian Information Criterion) ($\lg n$) bits is that using a fewer number of bits allows us to select even those features which provide marginal benefit. This is important since our goal is to build sparse models with better predictive accuracy rather than identifying the correct set of sparse features. We explain later in the section on Model Consistency that using a liberal penalty of 2 bits to code a coefficient allows us to achieve finite *risk-inflation* (Foster and George, 1994).

Combining the cost of the residuals with the cost of the model gives the following formula for the description length as a function of number of features that we include in the model:

$$-\lg P(\mathbf{Y}|\hat{\mathbf{w}}) + q(\lg p + 2), \quad (4)$$

where q is the number of features in the model and p is the total number of candidate features.

The above formula represents the simplest possible coding scenario and we will refer to it later in the paper as “**Baseline Coding Scheme**” when we propose more complex coding schemes for the problems of simultaneous feature selection for a set of multiple related tasks and grouped feature selection for a single task.

4. Multiple Inclusion Criterion (MIC)

MIC is a general framework for ℓ_0 penalty based greedy feature selection which minimizes a cost function provided by the Minimum Description Length (MDL) principle. MIC provides an elegant way of incorporating arbitrary sparsity patterns in the feature space by using MDL coding schemes customized to the problem at hand. In this section, we describe how MIC can be used to provide statistically efficient models for the problems of simultaneous feature selection for multiple related tasks and grouped feature selection for a single task. To do that, we first introduce some more notation and follow up on the MDL introduction in Section 3.1.

For the problem of simultaneous feature selection for a set of related tasks (which is addressed using MIC-MULTI) we assume a set of h regression or classification tasks which can potentially share a set of p features and a total of n labeled training examples. The task is to learn a set of joint (“shared”) models for all the h tasks. We represent the feature, response and the weight matrices as $\mathbf{X}_{n \times p}$, $\mathbf{Y}_{n \times h}$ and $\mathbf{w}_{p \times h}$ respectively. Additionally, for simplicity of analysis we assume a linear regression setting of the form⁴ $\mathbf{Y} = \mathbf{w} \cdot \mathbf{X} + \boldsymbol{\varepsilon}$ with a Gaussian noise term $\boldsymbol{\varepsilon}_{n \times h}$. Note that the noise on the responses ($\boldsymbol{\varepsilon}$) may be correlated; for instance, if our responses consist of temperature measurements at various locations, taken with the same thermometer, then if our thermometer drifted high at one location, it will have been high at the other location also. Thus, we take the rows of $\boldsymbol{\varepsilon}$ to have non-zero covariance:

$$\boldsymbol{\varepsilon}_i \sim \mathcal{N}_h(0, \boldsymbol{\Sigma}),$$

where $\boldsymbol{\varepsilon}_i$ is the i^{th} row of $\boldsymbol{\varepsilon}$ and $\boldsymbol{\Sigma}$ is an arbitrary $h \times h$ covariance matrix.

Similarly, for the related problem of grouped feature selection (which is addressed using MIC-GROUP) also, we have a total of p candidate features which are further divided into K groups (equal or unequal). Again, we assume the availability of a (fixed number) n of labeled training examples. Just as above we can represent the feature, response and weight matrices as $\mathbf{X}_{n \times p}$, $\mathbf{Y}_{n \times 1}$ and $\mathbf{w}_{p \times 1}$ respectively.

Let S represent the total description length (TDL) of the MDL message that is exchanged between the Sender and the Receiver. In the case of MIC-MULTI, S is the combined message length for all h tasks and hence we select features for all the h tasks simultaneously to minimize S and in the case of MIC-GROUP it can either be the combined message length for all the features within a given group (feature class) (MIC-GROUP(I)) or the message length of a given feature (MIC-GROUP-SC). Thus, when we evaluate a feature for addition into the model, we want to maximize the reduction of TDL by adding that feature to our model. More formally, at each iteration we greedily add those features to our model that:

$$\begin{aligned} \Delta S^i &= \Delta S_E^i - \Delta S_M^i, \\ \text{Best Feature} &= \underset{\hat{i}}{\operatorname{argmax}} \{ \Delta S^i \}, \end{aligned}$$

where $\Delta S_E \geq 0$ is the reduction in residual-error coding cost, that is, the first term on right hand side in Equation 1, due to the increase in data likelihood given this new feature and $\Delta S_M > 0$ is the

4. It can be extended to the standard classification setting by replacing the squared loss with a logistic loss, but due to lack of closed form solutions for logistic regression and since correlation between residuals is inconvenient to model in classification settings, we refrain from analyzing them.

increase in model cost to encode the new feature (second term in Equation 1) and i ranges over all the p features.

In the next subsections we describe how we code the S_E and S_M terms (i.e., the residual error and model) for MIC-MULTI and MIC-GROUP in detail.

4.1 MIC-MULTI

As mentioned earlier, MIC-MULTI borrows strength across multiple tasks and hence selects a joint set of features for related tasks (Dhillon et al., 2009).

4.1.1 CODING THE MODEL

MIC-MULTI borrows strength across responses by efficiently specifying the feature-response pairs in the $p \times h$ matrix \hat{w} . The naïve approach would be to put each of the ph coefficients in a linear order and specify the index of the desired coefficient using $\lg(mh)$ bits. But we can do better. If we expect nearly all the responses to be correlated with the predictive features, we could give all the responses nonzero coefficients (using $2h$ bits to code each of the h response coefficients) and simply specify the feature that we are talking about by using $\lg p$ bits, as in Section 3.1.2. From now on we will refer to this approach as FULL-MIC-MULTI (*fully dependent MIC-MULTI*) coding scheme, as it assumes that a selected feature will be added in the models of all the tasks, in much the same way as BBLASSO (Obozinski et al., 2009). Another limiting case is the one when we do feature selection for all the tasks independently (the baseline “Independent” Coding Scheme); the coding scheme in that case takes the form given in Equation 4.

However, these assumptions are usually unrealistic; each feature is generally neither correlated with almost all the responses nor with none of the responses, but is rather correlated with a few of them. A more flexible coding scheme would allow us to specify only the subset of the responses to which we want to give nonzero coefficients. For instance, suppose we are considering feature number 2609; and, of the $h = 20$ responses, we think that only $\{3, 7, 14, 17\}$ should have nonzero coefficients with the current feature. Then, we can use $\lg p$ bits to specify our feature (number 2609) once, and then we can list the particular responses that have nonzero coefficients with feature 2609, thereby avoiding paying the cost of $\lg(mh)$ four times to specify each coefficient in isolation.

A standard practice in information theory literature to code a subset of size h is to first specify how many $k \leq h$ elements the subset contains and then which of the $\binom{h}{k}$ possible subsets with k elements we are referring to Cover and Thomas (2006). In particular, we choose to code k using $\lg^* k + c_h$ bits, with c_h as defined in Equation 3. We then need $\lg \binom{h}{k}$ additional bits to specify the particular subset. We refer to this code as *partially dependent MIC-MULTI* or simply PARTIAL-MIC-MULTI.

The total cost (S_M^i) to code the model of a feature for MIC-MULTI is composed of three parts as follows:

$$S_M^i = \ell_H + \ell_I + \ell_\theta,$$

where ℓ_H is the number of bits needed to specify the subset k of the h tasks models in which to include the feature; ℓ_I is the number of bits used to describe which feature is being added and ℓ_θ is the description length of the coefficients of non-zero features.

We have already described the cost for ℓ_H above; it is equal to:

$$\ell_H = \lg^* k + c_h + \lg \binom{h}{k}.$$

For ℓ_θ , we use a cost of 2 bits per coefficient, the motivation for which was described earlier in Section 3.1.2. For ℓ_I , which specifies the size of the code for the given feature, we use $\lg p$ bits, which is equivalent to a uniform prior over the features,⁵ that is, each feature is equally likely to be selected. This can be accomplished by simply keeping a linear array of features and coding the indices of the features with nonzero coefficients.

Thus, we can represent the total model cost for MIC-MULTI as:

$$S_M = \left(\lg^* k + c_h + \lg \binom{h}{k} \right) + (\lg p) + (2k). \quad (5)$$

4.1.2 CODING THE DATA

Let \mathbf{E} be the residual error ($\mathbf{Y} - \mathbf{X} \cdot \mathbf{w}$) matrix, and as mentioned above, let ϵ_i , $i = 1, 2, \dots, n$ denote the i^{th} row of the error and let Σ be its $h \times h$ covariance matrix. The model likelihood under the Gaussian assumption⁶ can be written as:

$$P(Y_i | \hat{w}_q) = \frac{1}{\sqrt{(2\pi)^h |\Sigma|}} \exp \left(-\frac{1}{2} \epsilon_i^T \Sigma^{-1} \epsilon_i \right), \quad (6)$$

$$\begin{aligned} S_E &= -\lg \prod_{i=1}^n P(Y_i | \hat{w}_q) \\ &= \frac{1}{2 \ln 2} \left[n \ln \left((2\pi)^h |\Sigma| \right) + \sum_{i=1}^n (Y_i - X_i \cdot \hat{w}_q)^T \Sigma^{-1} (Y_i - X_i \cdot \hat{w}_q) \right] \end{aligned} \quad (7)$$

with subscript i denoting the i^{th} row. Since Σ is in fact unknown, we estimate it using maximum likelihood (ML):

$$\hat{\Sigma}_F = \frac{1}{n} (\mathbf{Y} - \mathbf{X} \cdot \hat{w}_{q-1})^T (\mathbf{Y} - \mathbf{X} \cdot \hat{w}_{q-1}),$$

where the subscript F stands for “full covariance”, and we use \hat{w}_{q-1} to get ML estimate, instead of \hat{w}_q to prevent overfitting, as we mentioned in Section 3.1.1.

In practice, we find that estimating all the h^2 entries of the covariance matrix can lead to overfitting. Therefore we use shrunken estimates of the form $\hat{\Sigma}_\lambda = \lambda \hat{\Sigma}_D + (1 - \lambda) \hat{\Sigma}_F$ for $(\lambda \in [0, 1])$, which tend to work well. Note that we do not need to pay an extra coding cost for estimating Σ as we are using a *prequential* coding scheme; Σ is calculated using information that was already paid for. We describe more technical details about our implementation in the Experiments section.

S_E as described in Equation 7 and S_M as Equation 5 form the quantities $\mathcal{D}(\mathbf{Y} | \hat{\mathbf{w}})$ and $\mathcal{D}(\hat{\mathbf{w}})$ in Equation 1. Thus our objective becomes $\hat{\mathbf{w}}^* = \underset{\hat{\mathbf{w}}}{\operatorname{argmin}} \{S_E + S_M\}$.

4.1.3 COMPARISON OF VARIOUS MIC-MULTI CODING SCHEMES

In this section we discussed three MDL based information-theoretic approaches to multitask feature selection, namely FULL-MIC-MULTI, Baseline “Independent” Coding Scheme and PARTIAL-MIC-MULTI. In general, the negative log-likelihood portion of *Independent* may differ from that of

5. The uniform code gives the worst-case minimax optimal code lengths (Grünwald, 2005).

6. As mentioned earlier, we are considering linear regression for simplicity of analysis and ease of modeling the correlation between residuals.

the other two methods, because *Full* and *Partial* can use a non-diagonal covariance estimate like $\hat{\Sigma}_F$ or $\hat{\Sigma}_\lambda$, while *Independent* only operates on one response at a time, and thus implicitly uses $\hat{\Sigma}_D$. However, since we generally use $\hat{\Sigma}_\lambda$, as mentioned earlier, for *Full* and *Partial*, and since $\hat{\Sigma}_\lambda$ approximates the diagonal $\hat{\Sigma}_D$ (for λ close to 1), the real difference comes from the coding schemes.

The coding costs for these three methods are compared in Table 1 for $p = 2000$ features, $h = 20$ responses, and for various values of k , the number of responses to which we add the current feature under consideration. FULL-MIC-MULTI is only allowed to take $k = 0$ or $k = h$, so it has h nonzero coefficients in all three rows of the table. However, if the extra $h - k$ coefficients correspond to non-predictive features, the extra reduction in residual-coding cost that FULL-MIC-MULTI enjoys over the other methods is likely to be small. As expected, each coding scheme is cheapest in the case for which it was designed; however, the MIC-MULTI methods are never excessively expensive, unlike *Independent* for $k = h$.

k	PARTIAL-MIC-MULTI	FULL-MIC-MULTI	Baseline (Independent)
1	$\log p + c_h + \log h + 2$ [18.4]	$\log p + 2h$ [51.0]	$\log p + 2$ [13.0]
$\frac{h}{4}$	$\log p + \log^* \left(\frac{h}{4}\right) + c_h + \log \binom{h}{h/4} + \frac{h}{2}$ [39.8]	$\log p + 2h$ [51.0]	$\frac{h}{4} \log p + \frac{h}{2}$ [64.8]
h	$\log p + \log^* h + c_h + 2h$ [59.7]	$\log p + 2h$ [51.0]	$h \log p + 2h$ [259.3]

Table 1: Costs in bits for each of the three schemes to code a model with $k = 1$, $k = \frac{h}{4}$, and $k = h$ nonzero coefficients. $p \gg h \gg 1$, $\ell_I = \lg p$, $\ell_{\hat{\theta}} = 2$, and for $h \in \{5, \dots, 1000\}$, $c_h \approx 1$. Examples of these values for $p = 2,000$ and $h = 20$ appear in brackets; the smallest of the costs appears in bold. *Note:* The costs are given per feature.

4.2 MIC-GROUP

MIC-GROUP is the algorithm for grouped feature selection, when features fall into *groups* or *classes* (Dhillon et al., 2008; Yuan and Lin, 2006; Bach et al., 2004). For example, genes can be divided into gene classes based on what pathway they occur in or features of a word can be grouped based on whether they are based on specific neighbouring words, parts of speech, or more global document properties. More generically, starting from any set of features, one can generate new classes of features by using projections such as principle components analysis (PCA) or non-negative matrix factorization (NNMF), transformations such as log or square root, and interactions (products of features) (Dhillon et al., 2010). The problem of grouped feature selection (MIC-GROUP) is very closely related to the problem of simultaneous feature selection for a set of related tasks (MIC-MULTI) as has also been pointed out by Obozinski et al. (2009). The multi-task problem we described earlier can also be thought of as a grouped feature selection scenario in which a group is defined by fixing a specific feature and ranging over multiple tasks. Our MIC based models for these two problems also follow the same intuition; in (MIC-MULTI) the tendency is to add a given feature into models of more and more tasks⁷ and similarly in (MIC-GROUP) the tendency is to add more and more features from the same group as the whole rationale behind doing grouped

7. The $\lg \binom{h}{k}$ part of the model cost is only small when k is small or it is very large, that is, $k \approx h$ as $\binom{h}{k} = \binom{h}{h-k}$.

feature selection is based on the fact that some feature groups contain highly predictive features than others.

4.2.1 CODING SCHEMES FOR MIC-GROUP

Since the problem of grouped feature selection is similar to the problem of simultaneous feature selection for a set of related tasks, we can propose a coding scheme which is analogous to the coding scheme for MIC-MULTI. For example, in this case also we can code the data as $P(Y_i|\hat{w}_q) = \frac{1}{\sqrt{(2\pi)^{h_{single}}|\Sigma|}} \exp\left(\frac{1}{2}\epsilon_i^T \Sigma^{-1}\epsilon_i\right)$ in a similar fashion as Equation 6 where h_{single} is the number of features in a given group (feature class) and we will estimate the covariance matrix, which represents covariance between different features in the same group (feature class), in a similar way as we did for MIC-MULTI, that is, by Maximum Likelihood Estimation. Remember, that in this case S_E term will be the message length for all the features within a given feature class. In a similar fashion, the number of bits to code the model can be represented as $S_M = \left[\lg^* k + c_h + \lg\binom{h_{single}}{k}\right] + \log p + 2k$, which corresponds to Equation 5.⁸ The other mechanics of the coding scheme will also be the same as for MIC-MULTI as this time we are trying to find a best subset of size k in a group (feature class) of size h_{single} and so we do a stepwise greedy search as earlier. From now on we refer to this coding scheme as MIC-GROUP (I).

Although this coding scheme works very well in practice, but it turns out that we are not exploiting the full flexibility that MDL based coding offer us. So, we propose a new coding scheme, which is computationally more efficient than MIC-GROUP (I), as it does not require a stepwise search for subset selection, though the predictive accuracy of both these coding schemes is comparable. We call this new computationally efficient coding scheme which uses ‘‘Switch Coding’’ as MIC-GROUP-SC and it is explained in detail below.

Coding the data with MIC-GROUP-SC (MIC-GROUP-Switch Coding): In this new coding scheme S_E is the message length for a single feature and ΔS_E represents the increase in likelihood of the data by adding that feature to the model.

Let \mathbf{E} be the residual error ($\mathbf{Y} - \mathbf{X}\cdot\mathbf{w}$) matrix as earlier, and let $\epsilon_i, i = 1, 2, \dots, n$ denote the i^{th} row of the error and let σ be variance of the Gaussian noise. The model likelihood can be written as:

$$\begin{aligned}
 P(Y_i|\hat{w}_q) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\epsilon_i^2}{2\sigma^2}\right), \\
 S_E &= -\lg \prod_{i=1}^n P(Y_i|\hat{w}_q).
 \end{aligned}
 \tag{8}$$

This equation is similar to the corresponding equation for MIC-MULTI except that here we only have a single response (task). In this case also, the variance σ^2 is estimated using the Maximum Likelihood principle.

Coding the model with MIC-GROUP-SC: This is where we differ from MIC-GROUP (I) and we use a coding scheme better suited to the group structure of the features. The intuition behind this coding scheme is that once we have selected (at least) one feature from a given group then it should

8. For simplicity of analysis and ease of comparison with coding schemes for MIC-MULTI we are assuming that all groups are of the same size h_{single} , though in reality the groups may be of unequal size and the same coding scheme still holds.

become easier to select more features from the group. In other words, the cost of adding more features from the same feature class should be low. The total cost of the model is composed of three parts:

$$S_M^i = \ell_C + \ell_I + \ell_\theta,$$

where ℓ_C is the number of bits used to code the index of the group of the evaluated feature, ℓ_I is the number of bits to code the index of the evaluated feature (within that particular group) and ℓ_θ is the cost to code the coefficient of the evaluated feature.

Our coding scheme assumes a uniform prior over all the groups; that is, each group is equally likely to contain beneficial features.⁹ So ℓ_C is $\lg K$ where K is the total number of groups (feature classes) in the data. Now, if a feature gets selected from a group (feature class) from which we had previously selected features, then we can save some bits by using “switch” coding and coding ℓ_C using only $1 + \lg Q$ bits where Q is the total number of groups (feature classes) included in the model till that point of time and 1 bit is used to represent that this group (feature class) has previously produced beneficial features. (Think of keeping an indexed list of length Q of the feature classes that have been selected). This is where our method wins over other methods and we do not need to code the same feature class over and over again if it has produced beneficial features in the past. Therefore ℓ_C is (Note that we added 1 bit to $\lg K$ also to ensure that the group whose index starts with 1 is not confused with the “switch”.):

$$\ell_C = \begin{cases} 1 + \lg K & \text{if the feature class is not in the} \\ & \text{model} \\ 1 + \lg Q & \text{if the feature class is already in} \\ & \text{the model.} \end{cases}$$

To code ℓ_I we again assume a uniform prior over all the features within that particular group. This corresponds to $\lg m_i$ bits where m_i is the total number of features in the feature class of which the i^{th} feature is a part of. This is pretty similar to RIC (Risk Inflation Criterion) style (Foster and George, 1994) coding or the widely use Bonferroni penalty. Finally, to code ℓ_θ we use 2 bits per coefficient, the motivation for which was described earlier. Therefore, the model cost per feature can be represented as:

$$S_M^i = (\ell_C) + (\lg m_i) + 2. \quad (9)$$

As mentioned earlier, this coding scheme is computationally cheaper than MIC-GROUP (I) as it does not require a subset search every time a feature is added to the model and it provides comparable predictive accuracy to MIC-GROUP (I). Note that just analogous to MIC-MULTI it is possible to come up with a new coding scheme called FULL MIC-GROUP(I) which just like its MIC-MULTI counterpart would add all the features from a given group (feature class) into the model. The MIC-GROUP schemes presented here are the most general setting and are analogous to PARTIAL MIC-MULTI for the multi response (task) scenario.

Just as with MIC-MULTI, here we optimize the objective $\hat{\mathbf{w}}^* = \underset{\hat{\mathbf{w}}}{\operatorname{argmin}}\{S_E + S_M\}$ with the S_E and S_M terms as described by Equations 8 and 9 respectively.

9. This is actually a pretty good assumption as mentioned earlier. The uniform code gives the worst-case minimax optimal code lengths (Grünwald, 2005) and hence it is reasonable to use it if the data distribution is *completely unknown* or if *no* distribution is assumed.

4.3 Algorithms and Implementation Details

In this subsection we outline the algorithms for MIC-MULTI and MIC-GROUP and also explain some details of the search strategy that we used for efficient subset search in case of MIC-MULTI.

4.3.1 ALGORITHMS

The algorithm for MIC-MULTI is as described in Algorithm 1. We provide algorithm for the most general case, that is, PARTIAL MIC-MULTI as the other two cases, that is, *Full* and *Independent* are the special cases of this scenario.

Algorithm 1 PARTIAL MIC-MULTI

```

1: Include the intercept (feature number 1) in all  $h$  response models.
2:  $remaining\_features = \{2, \dots, p\}$ .
3:  $keep\_adding\_features = true$ .
4: while  $keep\_adding\_features$  do
5:   for  $j$  in  $remaining\_features$  do
6:     // Find the best subset of response models to which to add feature  $j$ .
7:     for  $k = 1$  to  $h$  do
8:       Try including feature  $j$  in the best  $k$  response models. (We greedily assume that the best
9:        $k$  responses are the union of the best  $k - 1$  responses with the remaining response that,
10:      if included, would most increase likelihood.)
11:      Compute  $\Delta S_{jE}^k$ , the decrease in data residual cost, and  $\Delta S_{jM}^k$ , the resulting increase in
12:      model-coding cost, relative to not including feature  $j$  in any response models.
13:    end for
14:    Let  $k_j$  be the value of  $k$  that maximizes  $\Delta S_{jE}^k - \Delta S_{jM}^k$ .
15:     $\Delta S_j := \Delta S_{jE}^{k_j} - \Delta S_{jM}^{k_j}$ .
16:  end for
17:  Let  $j^*$  be the feature  $j$  that maximizes  $\Delta S_j$ , the reduction in TDL for adding feature  $j$ .
18:  if  $\Delta S_{j^*} > 0$  then
19:    Add feature  $j^*$  to the appropriate  $k_{j^*}$  response models.
20:     $remaining\_features = remaining\_features - \{j^*\}$ .
21:  else
22:     $keep\_adding\_features = false$ .
23:  end if
24: end while

```

The algorithm for MIC-GROUP-SC is described in Algorithm 2. The algorithm makes multiple passes through data and at each iteration adds the best feature to the model. It stops when no feature provides better ΔS than in the previous iteration. Since, it can be the case that it is not worth adding a single feature from a particular group (feature class) but it is still beneficial to add multiple features from that class. So, a clever search strategy that we found helpful with MIC-GROUP-SC was to use a mixed forward-backward greedy stepwise strategy in which one continues the search past the stopping criterion given in the algorithm and then sequentially removes the “worst” features from the now overfit model by making a “Backward” pass. In practice, we found this search strategy helpful. A similar hybrid forward-backward strategy was also used by Zhang (2009a).

Algorithm 2 MIC-GROUP-SC

```

1: flag = True; // flag for indicating when to stop
2: model = {}; // initially no features in model
3: prev_max = 0; // keeps track of the value of  $\Delta S_E$  in the previous iteration
4: while {flag == True} do
5:   for {i = 1 to p} do
6:     Compute  $\Delta S_E^i$ ; // Increase in likelihood by adding feature 'i' to the model
7:     Compute  $\Delta S_M^i$ ; // Number of extra bits required to code the  $i^{th}$  feature
8:      $\Delta S^i := \Delta S_E^i - \Delta S_M^i$ ;
9:   end for
10:   $i_{max} := \operatorname{argmax}_i \{\Delta S^i\}$ ; //The best feature in the current iteration
11:   $current\_max := \max_i \{\Delta S^i\}$ ; //The best penalized likelihood change in the current iteration
12:  if {current_max > prev_max} then
13:    model := model  $\cup$  { $i_{max}$ }; // Add the current feature to model
14:    prev_max := current_max;
15:  else
16:    flag := False;
17:  end if
18: end while

```

Note that we do not provide algorithm for MIC-GROUP (I) as it is pretty similar to MIC-MULTI with minor notational modifications as mentioned in the previous subsection.

4.3.2 STEPWISE SEARCH METHOD

Since MIC-MULTI requires subset search over the set of possible tasks in which to consider a feature for addition, so a discussion of our greedy search strategy is warranted.

For each feature, we evaluate the change in TDL (Total Description Length) that would result from adding that feature to the model with the optimal number of associated tasks. We add the best feature and then recompute the changes in TDL for the remaining features.¹⁰ This continues until there are no more features that would reduce TDL if added. The number of evaluations of features for possible addition is thus $O(pp_s)$, where p_s is the number of features eventually added.

To select the optimal number k of task models in which to include a given feature, we again use a stepwise-style search. In this case, we evaluate the reduction in TDL that would result from adding the feature to each task, add the feature to the best task, recompute the reduction in TDL for the remaining tasks, and continue.¹¹ However, unlike a normal stepwise search, we continue this process until we have added the feature to all h task models. The reason for this is two-fold. First, because we want to borrow strength across tasks, we need to avoid overlooking cases where the correlation of a feature with any single task is insufficiently strong to warrant addition, yet the

10. Remember that TDL changes due to the increase in likelihood and the additional model coding cost due to the added feature.

11. A stepwise search that re-evaluates the quality of each task at each iteration is necessary because, if we take the covariance matrix Σ to be non-diagonal, the values of the residuals for one task may affect the likelihood of residuals for other tasks. If we take Σ to be diagonal, as we do in Section 5, then an $O(h)$ search through the tasks without re-evaluation suffices.

correlations with all of the tasks are. Second, the $\log \binom{h}{k}$ term in PARTIAL MIC-MULTI’s coding cost does not increase monotonically with k , so even if adding the feature to an intermediate number of tasks does not look promising, adding it to all of them might still be worthwhile. Thus, when evaluating a given feature, we compute the description length of the model $O(h^2)$ times. Since we need to identify the optimal k for each feature evaluation, the entire algorithm requires $O(h^2 pp_s)$ evaluations of TDL.

While not shown explicitly in Algorithm 1, we use two branch-and-bound-style optimizations to cut this cost significantly in practice:

1. Before searching through subsets of responses to find the optimal subset for each feature, we make an $O(p)$ sweep through the features to compute an upper bound on the decrease in TDL that could result from adding that feature as

$$(\text{decrease in TDL if the feature is added to all } h \text{ response models}) - \log p.$$

Here, the first term is an upper bound on the benefit of adding the feature to the optimal number of response models (since adding a feature can only make a model fit better), and the second term underestimates the model cost of adding the feature, regardless of how many response models would actually be used. We sort the features in decreasing order by this upper bound, and when we reach features whose upper bounds are less than the best actual decrease in TDL observed so far, we terminate the search early.

2. For the stepwise search over responses, we can bound from above the potential benefit of adding the feature to k response models as

$$\begin{aligned} & (\text{decrease in TDL if the feature is added to all } h \text{ response models}) \\ & - \left(\log^* k + c_k + \log \binom{h}{k} + 2k \right), \end{aligned}$$

where the subtracted term represents the coding cost of including the feature in k response models. We can stop the search early when no higher value of k has an upper bound that exceeds the best reduction in TDL seen so far for any feature’s response subset.¹²

5. Experimental Results

In this section we empirically show the usefulness of our MIC based models (MIC [MULTI and GROUP]) on a variety of real world data sets pertaining to Genomics and Computational Linguistics (particularly Word Sense Disambiguation) domains. Besides this we also show results on synthetic data sets to illustrate the cases when our models are most beneficial.

It is important to note that we are interested in maximizing the predictive accuracy in these experiments rather than the totally orthogonal and antagonistic objective of identifying the correct set of sparse features (“sparsity pattern consistency”) at the expense of predictive accuracy. We would like to note that much sparser results can be obtained by using our models by using more conservative coding schemes and also by ℓ_1/ℓ_2 penalty (BBLASSO, GROUP LASSO) models by using thresholding techniques like Zhou (2009) and Lounici (2008).

12. We say “no higher value of k ” rather than “the next higher value of k ” because (10) does not decrease monotonically with k , due to the $\log \binom{h}{k}$ quantity.

5.1 MIC-MULTI

In this section, we first evaluate the MIC-MULTI approach on three synthetic data sets, each of which is designed to match the assumptions of, respectively, the PARTIAL and FULL MIC-MULTI, and Baseline (Independent) coding scheme (Equation 4). We then test the methods on two biological data sets, a Yeast Growth data set (Perlstein et al., 2007), which consists of real-valued growth measurements of multiple strains of yeast under different drug conditions, and a Breast Cancer data set (van 't Veer et al., 2002), which involves predicting prognosis, ER (Estrogen Receptor) status, and three other descriptive variables from gene-expression values for different cell lines.

We compare the three coding schemes of Section 4.1.3 against two other multitask algorithms: ANDOZHANG (Ando and Zhang, 2005) and BBLASSO (Obozinski et al., 2009), as implemented in the Berkeley Transfer Learning Toolkit (Rakhlin, 2007). We did not compare MIC-MULTI with other methods from the toolkit as they all require the data to have additional structure, such as *meta-features* (Lee et al., 2007; Raina et al., 2006), or expect the features to be frequency counts, such as for the Hierarchical Dirichlet Processes algorithm. Also, none of the neglected methods does feature selection.

For ANDOZHANG we use 5-fold CV to find the best value of the parameter that Ando and Zhang (2005) call h (the dimension of the subspace Θ , not to be confused with h as we use it in this paper). We tried values in the range $[1, 100]$ as is done in Ando and Zhang (2005).

MIC-MULTI, as presented in Section 4.1.2, is a regression algorithm, but ANDOZHANG and BBLASSO are both designed for classification. Therefore, we made each of our responses binary 0/1 values before applying MIC-MULTI with a regular regression likelihood term. Once the features were selected, however, we used logistic regression applied to just those features to obtain MIC-MULTI's actual model coefficients.

As noted in Section 4.1.2, MIC-MULTI's negative log-likelihood term can be computed with an arbitrary $h \times h$ covariance matrix Σ among the h tasks. We did not estimate all the h^2 entries of Σ as it lead to overfitting, so we instead took Σ to be diagonal.¹³

5.1.1 EVALUATION ON SYNTHETIC DATA SETS

We created synthetic data according to three separate scenarios—called *Partial*, *Full*, and *Independent*. For each scenario, we generated a matrix of continuous responses as

$$\mathbf{Y}_{n \times h} = \mathbf{X}_{n \times p} \cdot \mathbf{w}_{p \times h} + \boldsymbol{\varepsilon}_{n \times h},$$

where $p = 2000$ features, $h = 20$ responses, and $n = 100$ observations. Then, to produce binary responses, we set to 1 those response values that were greater than or equal to the average value for their column and set to 0 the rest; this produced a roughly 50-50 split between 1's and 0's because of the normality of the data. Each nonzero entry of \mathbf{w} was i.i.d. $\mathcal{N}(0, 1)$, and entry of $\boldsymbol{\varepsilon}$ was i.i.d. $\mathcal{N}(0, 0.1)$, with no covariance among the $\boldsymbol{\varepsilon}$ entries for different tasks. Each task had $p^* = 4$ beneficial features, that is, each column of \mathbf{w} had 4 nonzero entries.

The scenarios differed according to the distribution of the beneficial features in \mathbf{w} .

13. Informal experiments showed that estimating Σ as a convex combination of the full and diagonal estimates (i.e., $\hat{\Sigma}_\lambda$) also works well but we chose to use diagonal Σ (i.e., $\hat{\Sigma}_D$) due to its simplicity and to show the advantage of using a better coding scheme to code the model as by using diagonal Σ Partial and Independent methods are the same except S_M (i.e., cost of coding the model).

- In the *Partial* scenario, the first feature was shared across all 20 responses, the second was shared across the first 15 responses, the third across the first 10 responses, and the fourth across the first 5 responses. Because each response had four features, those responses (6 – 20) that did not have all of the first four features had other features randomly distributed among the remaining features (5, 6, . . . , 2000).
- In the *Full* scenario, each response shared exactly features 1 – 4, with none of features 5 – 2000 being part of the model.
- In the *Independent* scenario, each response had four random features among candidate features 1, . . . , 2000.

For the synthetic data, we report precision and recall to measure the quality of feature selection. This can be done both at a coefficient¹⁴ level (Was each nonzero coefficient in w correctly identified as nonzero, and vice versa?) and at an overall feature level (For features with *any* nonzero coefficients, did we correctly identify them as having nonzero coefficients for any of the tasks, and vice versa?). Note that Full MIC-MULTI and BBLASSO always make entire rows of their estimated w matrices nonzero and so tend to have larger numbers of nonzero coefficients. Table 2 shows the performance of each of the methods on five instances of the Partial, Full, and Independent synthetic data sets. On the *Partial* data set, PARTIAL MIC-MULTI performed the best, closely followed by BASELINE (INDEPENDENT); on the *Full* synthetic data, FULL MIC-MULTI and PARTIAL MIC-MULTI performed equally well; and on the *Independent* synthetic data, the *Baseline* algorithm performed the best closely followed by PARTIAL MIC-MULTI. It is also worth noting that the best-performing methods tended to have the best precision and recall on coefficient selection. The performance trends of the three methods are in consonance with the theory of Section 4.1.3.

The table shows that only in one of the three cases does non-MIC methods compete with MIC methods. BBLASSO on the Full synthetic data shows comparable performance to the MIC methods, but even in that case it has a very low feature precision, since it added many more spurious features than the MIC methods.

5.1.2 EVALUATION ON REAL DATA SETS

This section compares the performance of MIC-MULTI methods with ANDOZHANG and BBLASSO on a Yeast and a Breast Cancer data set. These are typical biological data sets in that only a handful of features are predictive from thousands of potential features. This is precisely the case in which MIC-MULTI outperforms other methods. MIC-MULTI not only gives better accuracy, but does so by choosing fewer features than BBLASSO’s $\ell_1 - \ell_2$ -based approach.

Yeast Data Set: Our Yeast data set comes from Perlstein et al. (2007). It consists of real-valued growth measurements of 104 strains of yeast ($n = 104$ observations) under 313 drug conditions. In order to make computations faster, we hierarchically clustered these 313 conditions into 20 groups using correlation as the similarity measure. Taking the average of the values in each cluster produced $h = 20$ real-valued responses (tasks), which we then binarized into two categories: values at least as big as the average for that response (set to 1) and values below the average (set to 0). The features consisted of 526 markers (binary values indicating major or minor allele) and 6,189 transcript levels in rich media for a total of $p = 6715$ features.

14. A coefficient is defined as the addition of a given feature to a single task. For example if a feature was added to models of 10 tasks, then 1 feature and 10 coefficients were selected.

Method	Test Error $\mu \pm \sigma$	Coefficient Precision/Recall	Feature Precision/Recall
Partial Synthetic Data Set			
TRUE MODEL	0.07 ± 0.00	$1.00 \pm 0.00/1.00 \pm 0.00$	$1.00 \pm 0.00/1.00 \pm 0.00$
PARTIAL MIC-MULTI	0.10 ± 0.00	$0.84 \pm 0.02/0.77 \pm 0.02$	$0.99 \pm 0.01/0.54 \pm 0.05$
FULL MIC-MULTI	0.17 ± 0.01	$0.26 \pm 0.01/0.71 \pm 0.03$	$0.97 \pm 0.02/0.32 \pm 0.03$
BASELINE (INDEPENDENT)	0.12 ± 0.01	$0.84 \pm 0.02/0.56 \pm 0.02$	$0.72 \pm 0.05/0.62 \pm 0.04$
BBLASSO	0.19 ± 0.01	$0.04 \pm 0.00/0.81 \pm 0.02$	$0.20 \pm 0.03/0.54 \pm 0.01$
ANDOZHANG	0.50 ± 0.02	NA	NA
Full Synthetic Data Set			
TRUE MODEL	0.07 ± 0.00	$1.00 \pm 0.00/1.00 \pm 0.00$	$1.00 \pm 0.00/1.00 \pm 0.00$
PARTIAL MIC-MULTI	0.08 ± 0.00	$0.98 \pm 0.01/1.00 \pm 0.00$	$0.80 \pm 0.00/1.00 \pm 0.00$
FULL MIC-MULTI	0.08 ± 0.00	$0.80 \pm 0.00/1.00 \pm 0.00$	$0.80 \pm 0.00/1.00 \pm 0.00$
BASELINE (INDEPENDENT)	0.11 ± 0.01	$0.86 \pm 0.02/0.63 \pm 0.02$	$0.36 \pm 0.06/1.00 \pm 0.00$
BBLASSO	0.09 ± 0.00	$0.33 \pm 0.03/1.00 \pm 0.00$	$0.33 \pm 0.17/1.00 \pm 0.00$
ANDOZHANG	0.45 ± 0.02	NA	NA
Independent Synthetic Data Set			
TRUE MODEL	0.07 ± 0.00	$1.00 \pm 0.00/1.00 \pm 0.00$	$1.00 \pm 0.00/1.00 \pm 0.00$
PARTIAL MIC-MULTI	0.17 ± 0.01	$0.95 \pm 0.01/0.44 \pm 0.02$	$1.00 \pm 0.00/0.44 \pm 0.02$
FULL MIC-MULTI	0.36 ± 0.01	$0.06 \pm 0.01/0.15 \pm 0.02$	$1.00 \pm 0.00/0.14 \pm 0.02$
BASELINE (INDEPENDENT)	0.13 ± 0.01	$0.84 \pm 0.02/0.58 \pm 0.02$	$0.83 \pm 0.02/0.58 \pm 0.03$
BBLASSO	0.35 ± 0.01	$0.02 \pm 0.00/0.43 \pm 0.02$	$0.30 \pm 0.05/0.42 \pm 0.06$
ANDOZHANG	0.49 ± 0.00	NA	NA

Table 2: Test-set accuracy, precision, and recall of MIC-MULTI and other methods on 5 instances of various synthetic data sets generated as described in Section 5.1.1. Standard errors are reported over each task; that is, with 5 data sets and 20 tasks per data set, the standard errors represent the sample standard deviation of 100 values divided by $\sqrt{100}$. *Note:* ANDOZHANG’s NA values are due to the fact that it does not explicitly select features.

Figure 1 (a) shows classification test errors from 5-fold CV on this data set. As can be seen from the table, PARTIAL MIC-MULTI performs better than BBLASSO or ANDOZHANG. BASELINE and FULL MIC-MULTI perform slightly worse than PARTIAL MIC-MULTI, underscoring the point that it is preferable to use a more general MIC coding scheme compared to FULL MIC-MULTI or BASELINE. The latter methods have strong underlying assumptions, which cannot always correctly capture sharing across tasks.

Breast Cancer Data Set: Our second data set pertains to Breast Cancer, and contains data from five of the seven data sets used in van ’t Veer et al. (2002). It contains 1171 observations for 22,268 RMA-normalized gene-expression values. We considered five associated responses (tasks); two were binary—prognosis (“good” or “poor”) and ER (Estrogen Receptor) status (“positive” or “negative”)—and three were not—age (in years), tumor size (in mm), and grade (1, 2, or 3). We binarized the three non-binary responses into two categories: Response values at least as high as the average, and values below the average. Finally we scaled the data set down to $n = 100$ and

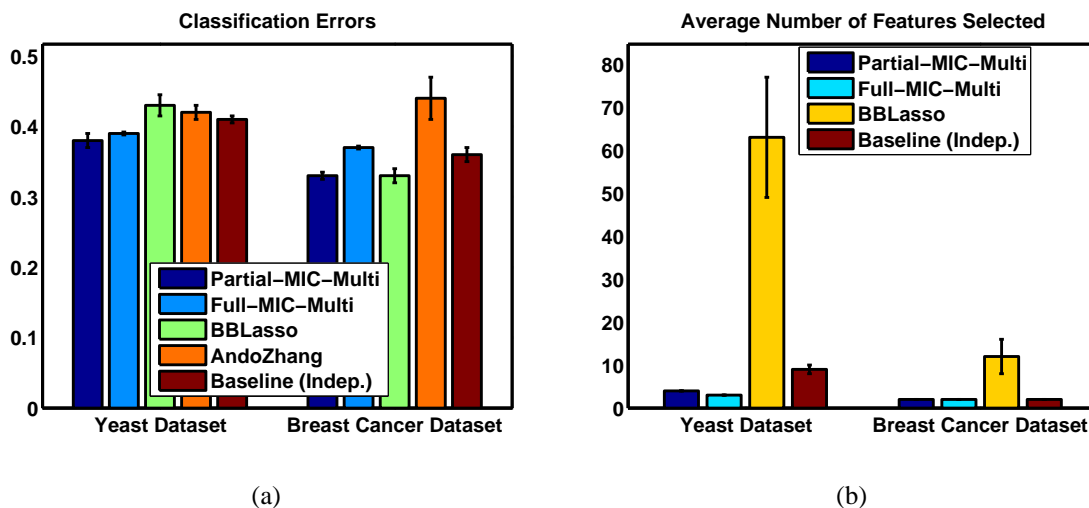


Figure 1: Accuracy and number of features selected on five folds of CV for the Yeast and Breast Cancer data sets. *Note:* 1). Remember that we are interested in better predictive accuracy and not in identifying the correct set of sparse features; we can get much sparser models if instead our objective is choosing the correct set of sparse features. 2). ANDOZHANG’s average number of features selected are not present in the graph as it does not explicitly select features. 3). These are true cross-validation accuracies and no parameters have been tuned on them.

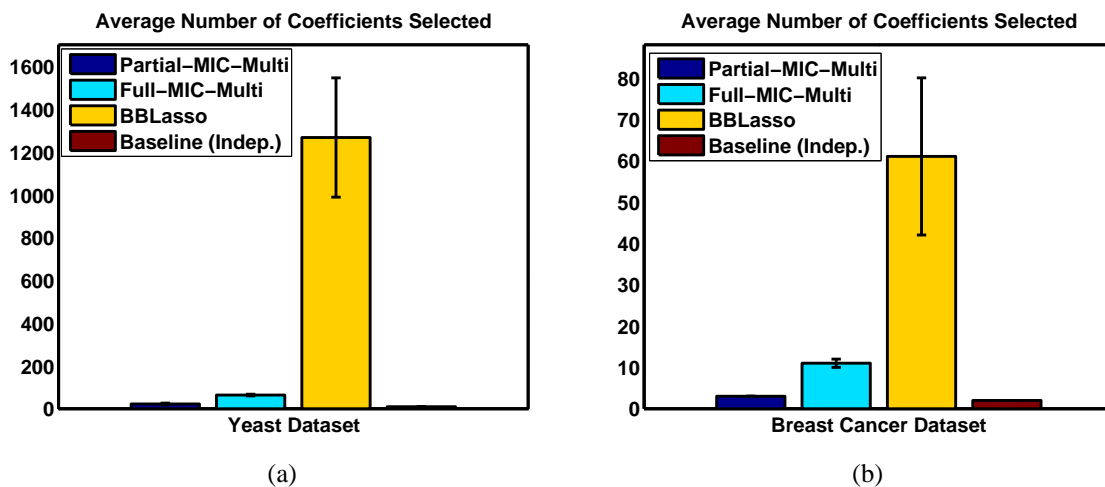


Figure 2: Number of coefficients selected on five folds of CV for the Yeast and Breast Cancer data sets. *Note:* 1). ANDOZHANG’s average number of coefficients selected are not present in the graph as it does not explicitly select features.

$p = 5000$ (the 5000 features with the highest variance), to save computational resources. Figure 1 (a) shows classification test errors from 5-fold CV on this data set. As is clear from the table, PARTIAL MIC-MULTI and BBLASSO are the best methods here. But as was the case with other data sets, BBLASSO puts in more features, which is undesirable in domains (like biology and medicine) where simpler and hence more interpretable model are sought.

The number of features and coefficients selected by all the methods are shown in Figures. 1 (b) and 2 respectively.

5.2 MIC-GROUP

In this section we demonstrate the results of the MIC-GROUP scheme on synthetic and real world data sets. For our experiments we use both the MIC-GROUP (I) and MIC-GROUP-SC (as described in Algorithm 2) methods and compare against BASELINE Feature Selection (which in this case is equivalent to a RIC penalized regression and has a coding scheme similar to Equation 4, Lasso (Tibshirani, 1996), Elastic Nets (Zou and Hastie, 2005) and Group Lasso/ Multiple Kernel Learning (Yuan and Lin, 2006; Jacob et al., 2009; Bach et al., 2004).

For Group Lasso/Multiple Kernel Learning,¹⁵ we used a set of 13 candidate kernels, consisting of 10 Gaussian Kernels (with bandwidths $\sigma = 0.5 - 20$) and 3 polynomial kernels (with degree 1-3) for each feature class as is done by Rakotomamonjy et al. (2008). In the end the kernels which have non zero weights are the ones that correspond to the selected feature classes. Since GL/MKL minimizes a mixed $\ell_1 - \ell_2$ norm so, it zeros out some groups (feature classes). However it is possible to estimate the exact support by thresholding (cross-validated) the estimated weights, as has been done by Zhou (2009) and Lounici (2008), and enforce sparsity within the groups also but as mentioned earlier our main goal is better predictive accuracy and not identifying the correct set of sparse features. The Group Lasso (Yuan and Lin, 2006; Jacob et al., 2009) and Multiple Kernel Learning are equivalent, as has been mentioned in Bach (2008), therefore we used the *SimpleMKL* toolbox (Rakotomamonjy et al., 2008) implementation for our experiments. For Lasso and Elastic Nets we used their standard LARS (Least Angle Regression) implementations (Efron et al., 2004). When running Lasso and Elastic Nets, we pre-screened the data sets and kept only the best $\sim 1,000$ features (based on their p-values), as otherwise LARS is prohibitively slow. (The authors of the code we used do similar screening, for similar reasons.) For all our experiments on Elastic Nets (Zou and Hastie, 2005) we chose the value of λ_2 (the weight on the ℓ_2 penalty term), as 10^{-6} .

We demonstrate the effectiveness of MIC-GROUP on synthetic data sets and on real data sets pertaining to Word Sense Disambiguation (WSD) (Chen and Palmer, 2005) (ONTONOTES Data Set Hovy et al., 2006) and gene expression data (Mootha et al., 2003).

5.2.1 EVALUATION ON SYNTHETIC DATA SETS

The main hypothesis is that MIC-GROUP methods are beneficial when some groups have multiple predictive features, while others lack them. MIC-GROUP is particularly effective when there are small groups which contain highly predictive features and big groups containing no predictive features.

In order to validate our hypothesis, we test MIC-GROUP on two synthetic data sets. For both the data sets, 1000 features were generated independently from a Normal Distribution $\mathcal{N}(0, 1)$, and the response vector of 100 observations Y was computed as the linear combination of a set of 7 beneficial features and Gaussian additive noise $\mathcal{N}(0, 1.7^2)$. The first data set (Set 1) had 4 groups (feature classes) of unequal sizes and 7 beneficial features, all of which lie in a small feature class of size 12. The second synthetic data set (Set 2) was generated so as to reflect the other extreme

15. There is a similar relation between MIC-GROUP and GL/MKL as it is between MIC-MULTI and BBLASSO. Both BBLASSO and GL/MKL are ℓ_1/ℓ_2 penalty based methods and try to solve the same sparsity problem as the corresponding MIC method.

case, in which all the classes are of same size, and had 100 feature classes, each of size 100. Again all 7 beneficial features were in a single feature class.

Method	Avg. Features Selected				10-Fold CV Error	
	Correct		Spurious		Set 1	Set 2
	Set 1	Set 2	Set 1	Set 2		
MIC-GROUP-SC	6.8 ± 0.1	5.6 ± 0.0	0.1 ± 0.0	0.3 ± 0.1	0.09 ± 0.02	0.27 ± 0.01
MIC-GROUP (I)	6.7 ± 0.0	5.4 ± 0.1	0.1 ± 0.1	0.2 ± 0.1	0.11 ± 0.02	0.28 ± 0.02
LASSO	5.2 ± 1.0	4.3 ± 1.2	2.2 ± 1.0	1.8 ± 0.1	0.22 ± 0.03	0.41 ± 0.02
ELASTIC NETS	6.4 ± 0.2	4.9 ± 0.7	3.3 ± 1.1	2.1 ± 1.3	0.20 ± 0.03	0.43 ± 0.02
BASELINE (RIC)	4.4 ± 1.4	3.2 ± 2.2	0.2 ± 0.1	0.0 ± 0.0	0.27 ± 0.05	0.61 ± 0.04

Table 3: The number of correct and spurious Features Selected and 10 Fold CV Test Errors averaged over 10 runs. Set 1). Unequal class sizes, Set 2). Uniform class sizes.

As can be seen from the results in Table 3, in both cases the MIC-GROUP methods outperform other competing methods.

5.2.2 EVALUATION ON REAL DATA SETS

In order to benchmark the real world performance of our MIC-GROUP, we chose two data sets pertaining to two diverse applications of feature selection methods, namely Computational Linguistics and Gene Expression Analysis. More information regarding the data and the experimental results are given below.

Word Sense Disambiguation (WSD) Data Sets: A WSD data set (ONTONOTES Hovy et al., 2006) consisting of 172 ambiguous verbs and a rich set of contextual features (Chen and Palmer, 2005) was chosen for evaluation. It consists of hundreds of observations of noun-noun collocation, noun-adjective-preposition-verb (syntactic relations in a sentence) and noun-noun combinations (in a sentence or document).

The data set had a total of 172 verbs with 40 – 45 feature classes (groups). The number of observations n for the various verbs varied from 100 to 3500 and the number of features p varied from 1000 to 11500.

As with MIC-MULTI we used MIC-GROUP to do feature selection and once we had selected the features we used logistic regression for the final classification problem. The classification test accuracies averaged over all the 172 verbs are shown in Figure 3.

Note that these accuracies are for the binary prediction problem of predicting the most frequent sense. On the entire set of 172 verbs, MIC-GROUP methods are significantly (5 % significance level (Paired t-Test)) better than the competing methods on 160/172 verbs and have the same accuracy as the best method on 4 occasions. It is also worth noting that MIC-GROUP-SC was ~ 7 times faster than MIC-GROUP (I) as we had hypothesized earlier, as for each selected feature it does a subset search within that feature’s group (feature class) to find the optimal number of features to select from the group.

Gene Set Enrichment Analysis (GSEA) Data Sets: The second real data sets that we used for our experiments were gene expression data sets from GSEA (Mootha et al., 2003). There are multiple gene expression data sets and multiple criteria on which the genes can be grouped into classes. For

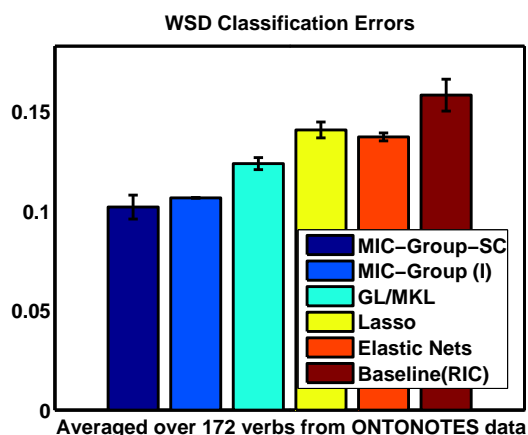


Figure 3: 10 Fold CV classification test accuracies averaged over 172 verbs. *Note:* 1). These are true cross-validation accuracies and no parameters have been tuned on them.

example, different ways of generated gene classes include C1: Positional Gene Sets, C2: Curated Gene Sets, C3: Motif Gene Sets, C4: Computational Gene Sets, C5: GO Gene Sets.

For our experiments, we used gene classes from the C1 and C2 collections. The gene sets in collection C1 consists of genes belonging to the entire human chromosome, divided into each cytogenetic band that has at least one gene. Collection C2 contained gene sets from various sources such as online pathway databases and knowledge of domain experts.

The data sets that we used and their specifications are as shown in Table 4. Though the goal of GSEA is not building classification models but identifying the groups of genes (gene families) which are over-represented when they are filtered by a certain selection procedure; however we are interested in using the transcriptional profiles and the associated group structure for classifying the phenotype, that is, ALL (Acute Lymphoblastic Leukemia) or AML (Acute Myeloid Leukemia) in case of leukemia; DMT (Diabetes Mellitus Type I) or NGT (Normal Glucose Tolerance) in case of diabetes; and determining whether the transcriptional profiles are from a M (Male) or a F (Female) for the gender data sets.

Data Set	# Observations (n)	# Features (p)	# Classes (K)
LEUKEMIA (C1)	48 (24 ALL & 24 AML)	10056	182
GENDER 1 (C1)	32 (17 F & 15 M)	15056	212
DIABETES (C2)	34 (17 NGT & 17 DMT)	15056	318
GENDER 2 (C2)	32 (17 F & 15 M)	15056	318

Table 4: GSEA Data Sets.

The results for these GSEA data sets are as shown in the Figure 4.

For these data sets also MIC-GROUP methods beat the competing methods. Here also MIC-GROUP is significantly (5% significance level, Paired t-test) better than the competing methods. It is interesting to note that MIC-GROUP methods sometimes selected substantially fewer features, but still gave better performance than other methods which goes onto show that adding all or many features from a single group contributes to a redundant signal and efficient feature selection “within” a group (feature class) is warranted.

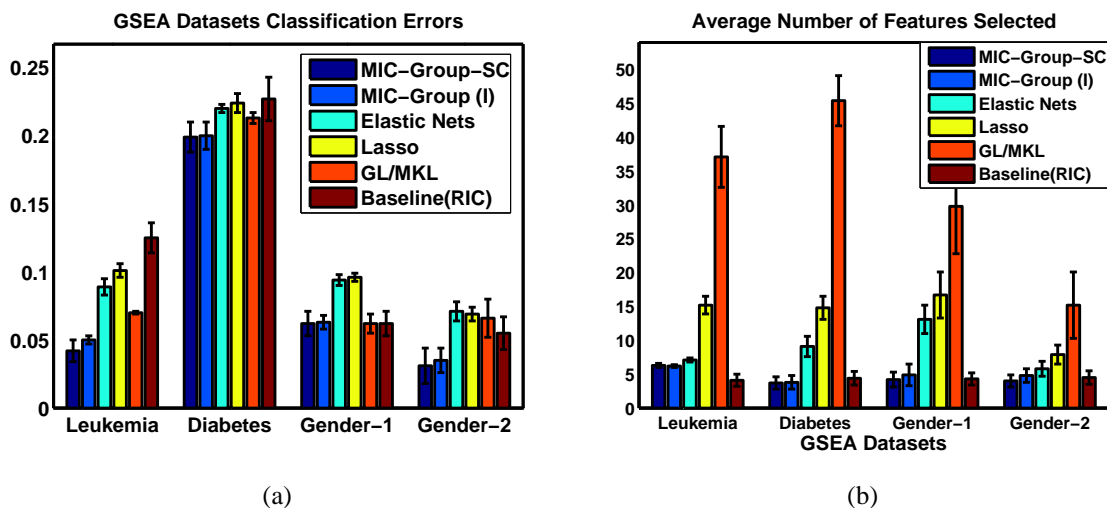


Figure 4: 10-fold CV classification test accuracies and the average number of features selected by various methods on the GSEA data sets. *Note:* 1). Remember that we are interested in better predictive accuracy and not in identifying the correct set of sparse features; we can get much sparser models if instead our objective is choosing the correct set of sparse features. 2). These are true cross-validation accuracies and no parameters have been tuned on them.

6. MIC Model Consistency

In this section we show that our MIC methods based on two part MDL and with the model coding costs as described in Section 4 are consistent. By “consistent” we mean that if the data is distributed by one of the probabilistic sources in the set of candidate model classes that our MDL based estimators consider (\mathcal{M}), then given enough data, MIC will output the true distribution generating the data. The proof of consistency is similar to the proof of classical two part MDL consistency as given in Barron and Cover (1991) and the recent improvement to that proof by Zhang (2004) by using ideas from KL-complexity. To extend these proofs to the case of MIC, we require the concepts of *Universal Codes*, *KL-Distinguishability*, *Probabilistic Sources* and the *No-Hypercompression Inequality* from information theory. (Refer to Appendix)

We first define some common notation that will be useful throughout this section. Assume we have n data samples (observations) $X_1 \in \mathcal{X}, X_2 \in \mathcal{X}, \dots, X_n \in \mathcal{X}$ and that they are distributed according to some distribution P_{True} . Further let $P_{MIC}^{(n)}$ be an arbitrary distribution on \mathcal{X} (the distribution estimated by our MIC based model). Also, as shorthand we denote $-\sum_{i=1}^n \lg Q(X_i)$ as $-\lg Q(X^n)$ throughout this section and $P^{(n)}$ denotes the marginal distribution on the first n outcomes that is induced by the probabilistic source P .

6.1 Consistency Results

As mentioned earlier, two part MDL has been proved to be consistent in a variety of settings (Barron and Cover, 1991; Zhang, 2004; Grünwald, 2007). Here we provide similar proofs for the case of our MIC based models.

Let \mathcal{M} be a countably infinite set of probabilistic sources and let ℓ_n be some code length function corresponding to a code over elements of \mathcal{M} . We assume that the true underlying probabilistic source belongs to the set of sources that our models consider, that is, $P_{True} \in \mathcal{M}$; Zhang (2004) and Barron and Cover (1991) also make this assumption. Also, let P_{MIC} be the probabilistic source corresponding to the two part MDL model selected by MIC and $P_{MIC}^{(n)}$ be the marginal distribution induced by this probabilistic source on the first n outcomes.

Theorem 1 *Let $\mathcal{M}^\delta = \{Q \in \mathcal{M} | KL(P_{True} || Q) \geq \delta\}$. Also, let $\ell_n(P_{True}) < \infty$, then*

$$P_{True}(P_{MIC}^{(n)} \in \mathcal{M}^\delta) \rightarrow 0 \text{ as } n \rightarrow \infty. \quad (10)$$

The theorem states that the probability that MIC selects a probabilistic source to explain the data that is KL-distinguishable from the true underlying distribution (P_{True}) approaches 0 as the number of observations increase. In other words, with overwhelming probability, P_{True} is KL-indistinguishable from $P_{MIC}^{(n)}$ as n approaches infinity.

Proof

Let \mathcal{P}_n be the distribution corresponding to the code for hypotheses with lengths ℓ_n , such that for all $Q \in \mathcal{M}$, $\mathcal{P}_n(Q) = 2^{-\ell_n(Q)}$. This follows from Kraft's Inequality (Grünwald, 2005). \mathcal{M} is countable, so \mathcal{M}^δ must also be countable; therefore we can order the elements in \mathcal{M}^δ according to increasing description length $\ell_n(Q)$ (decreasing $\mathcal{P}_n(Q)$) as Q_1, Q_2, \dots . Fix some θ ($0 < \theta < 1$) and define $\mathcal{M}_{1:N(\theta)}^\delta$ as the subset of \mathcal{M}^δ consisting of the first N distributions in \mathcal{M}^δ , where N is the smallest number such that $\sum_{j=1}^N \mathcal{P}_n(Q_j) \geq \theta$. Now, we define $\mathcal{M}_{N(\theta)+1:\infty}^\delta = \mathcal{M}^\delta \setminus \mathcal{M}_{1:N(\theta)}^\delta$, that is, $\mathcal{M}_{N(\theta)+1:\infty}^\delta = \{Q_{N+1}, Q_{N+2}, \dots\}$. It can be easily seen that,

$$\mathcal{P}_n(\mathcal{M}_{N(\theta)+1:\infty}^\delta) = \sum_{j=N+1}^{\infty} \mathcal{P}_n(Q_j) \leq 1 - \theta. \quad (11)$$

Now, for any $\mathcal{M}' \subseteq \mathcal{M}^\delta$,

$$\begin{aligned} P_{True} \left[P_{MIC}^{(n)} \in \mathcal{M}' \right] &= \\ P_{True} \left[\text{for some } Z \in \mathcal{M}' : \ell_n(P_{True}) + \ell_n(X^n | P_{True}) \geq \ell_n(Z) + \ell_n(X^n | Z) \right] \\ &\leq \sum_{Z \in \mathcal{M}'} P_{True} \left[\ell_n(P_{True}) + \ell_n(X^n | P_{True}) \geq \ell_n(Z) + \ell_n(X^n | Z) \right]. \end{aligned} \quad (12)$$

The above inequality is obtained by applying the Union Bound.

Now, by re-arranging Equation 12 and noting that the error term in two-part coding can be replaced by a log term as in Equation 2,

$$g_n(Z) = P_{True} \left[-\lg(P_{True}(X^n)) \geq -\lg(Z(X^n)) + \ell_n(Z) - \ell_n(P_{True}) \right]. \quad (13)$$

The $\ell_n(Z)$ term corresponds to the number of bits required to code the model and the $-\lg Z(X^n)$ term corresponds to the data likelihood term in the two part MDL coding scheme.

From Equations 12, 13 and 10, it follows that:

$$\begin{aligned} P_{True} \left[P_{MIC}^{(n)} \in \mathcal{M}^\delta \right] &= P_{True} \left[P_{MIC}^{(n)} \in \mathcal{M}_{1:N(\theta)}^\delta \right] + P_{True} \left[P_{MIC}^{(n)} \in \mathcal{M}_{N(\theta)+1:\infty}^\delta \right] \\ &\leq \sum_{Z \in \mathcal{M}_{1:N(\theta)}^\delta} g_n(Z) + \sum_{Z \in \mathcal{M}_{N(\theta)+1:\infty}^\delta} g_n(Z). \end{aligned} \quad (14)$$

Since, each $Z \in \mathcal{M}_{1:N(\theta)}^\delta$ is KL-Distinguishable from P_{True} and we have $g_n(Z) \leq e^{-n\delta}$ (Stein's Lemma Cover and Thomas, 2006) for all $Z \in \mathcal{M}^\delta$. Hence $\sum_{Z \in \mathcal{M}_{1:N(\theta)}^\delta} g_n(Z)$ is a sum of a finite number of Z 's, each of which is exponentially small, so we have $\lim_{n \rightarrow \infty} \sum_{Z \in \mathcal{M}_{1:N(\theta)}^\delta} g_n(Z) = 0$.

Now, we consider the second term in Equation 14. First assume that there exists $N'(\theta) \geq N(\theta) + 1$ such that the sources $\mathcal{M}_{N(\theta)+1:N'(\theta)}^\delta$ have description lengths (ℓ_n) smaller than the description length of P_{True} and the sources $\mathcal{M}_{N'(\theta)+1:\infty}^\delta$ have description lengths greater than $\ell_n(P_{True})$. Also, we have $\ell_n(P_{True}) = O(\log n)$, given all the coding schemes we have designed using MIC models. Now, by Kraft's Inequality it follows that

$$\begin{aligned} & \sum_{Z' \in \mathcal{M}_{N(\theta)+1:N'(\theta)}^\delta} 2^{-\ell_n(Z')} \leq 1 \\ \implies & \sum_{Z' \in \mathcal{M}_{N(\theta)+1:N'(\theta)}^\delta} 2^{-k \log n} \leq 1 \\ \implies & \sum_{Z' \in \mathcal{M}_{N(\theta)+1:N'(\theta)}^\delta} n^{-k} \leq 1. \end{aligned}$$

Therefore, there are at most a polynomial number n^k of elements in $\mathcal{M}_{N(\theta)+1:\infty}^\delta$ with shorter description lengths than P_{True} . Hence, as earlier $\sum_{Z' \in \mathcal{M}_{N(\theta)+1:N'(\theta)}^\delta} g_n(Z')$ is a sum of a finite number of (Z)s, each of which is exponentially small, so we have $\lim_{n \rightarrow \infty} \sum_{Z' \in \mathcal{M}_{N(\theta)+1:N'(\theta)}^\delta} g_n(Z') = 0$.

Now, we bound the remaining terms by applying the no-hypercompression inequality to each term in $\sum_{Z \in \mathcal{M}_{N'(\theta)+1:\infty}^\delta} g_n(Z)$ with $K = \ell_n(Z) - \ell_n(P_{True})$,

$$g_n(Z) = P_{True} [-\lg(P_{True}(X^n)) \geq -\lg(Z(X^n)) + \ell_n(Z) - \ell_n(P_{True})] \leq 2^{-\ell_n(Z) + \ell_n(P_{True})}. \quad (15)$$

From Equations 11 and 15, we get

$$\sum_{Z \in \mathcal{M}_{N'(\theta)+1:\infty}^\delta} g_n(Z) \leq \sum_{Z \in \mathcal{M}_{N'(\theta)+1:\infty}^\delta} 2^{-\ell_n(Z) + \ell_n(P_{True})} \leq (1 - \theta) \cdot 2^{\ell_n(P_{True})}.$$

The above holds for every $0 < \theta < 1$, so for every $\varepsilon > 0$ we can choose $\theta = 1 - \varepsilon \cdot 2^{-\ell_n(P_{True})}$ giving $\sum_{Z \in \mathcal{M}_{N'(\theta)+1:\infty}^\delta} g_n(Z) \leq \varepsilon$ for all large n . Combining this with Equation 14 we find that for all $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P_{True}[P_{MIC}^{(n)} \in \mathcal{M}^\delta] < \varepsilon. \quad \blacksquare$$

A corollary of the above theorem is that the MIC coding schemes as described in Section 4 are not an arbitrary procedure. There can be many valid codes, but we can not tweak MDL by using arbitrary codes to give the answers that we would like it to give.

Another important theoretical property that is attractive for sparse learning algorithms is *sparsistency*, which is shorthand for “sparsity pattern consistency”. In other words:

$$P[\text{supp}(\mathbf{w}^{True}) = \text{supp}(\mathbf{w}_n^{MIC})] \rightarrow 1 \text{ as } n \rightarrow \infty,$$

where $\text{supp}(\mathbf{w}) = \{\mathbf{w} : \mathbf{w}_j \neq 0\}$, \mathbf{w}^{True} is the true sparse weight vector and \mathbf{w}^{MIC} is the weight vector estimated by MIC based methods. Sparsistency implies that the learning algorithm is consistently able to identify the correct set of sparse features in the asymptotic limit.

Lasso and Group Lasso have been proved to be *sparsistent* under irrepresentable conditions that depend on the sign of the true weight vector (\mathbf{w}^{True}) (Zhao and Yu, 2006; Wainwright, 2009; Meinshausen and Bühlmann, 2006; Bach, 2008). Tropp (2004) proved that forward greedy feature selection also selects features consistently when the linear model has a zero-mean stochastic noise; Zhang (2009b) improved this result to include non-zero mean sub-Gaussian stochastic noise. However, due to the complexity of the forward greedy feature selection the sparsistency condition in this case depends only on the feature (design) matrix \mathbf{X} , unlike Lasso and Group Lasso.

Since our MIC based methods are based on forward greedy feature selection, that is, they use the MDL principle to provide a cost function which is greedily minimized by a forward search, they should be *sparsistent*. However, for ℓ_0 penalized regression, the *sparsistency* condition also depends on the information theoretic penalty in that the penalty must increase with n (the number of observations) (Wu and Zhou, 2010). For our MIC based methods this penalty is a combination of RIC, AIC (to code the coefficients) and other coding schemes which incorporate the structure of the problem at hand. The penalties for the MIC based methods as presented in this paper do not have the required dependence on n , so they are not sparsistent. However, we could modify our coding schemes slightly by using the BIC penalty ($\lg n$ bits) to code the coefficients instead of AIC to ensure sparsistency of MIC. However, we prefer that our methods are not sparsistent as in that case we achieve competitive performance with the true underlying model, that is, we get finite *risk-inflation* of about $2 \lg p$ (Foster and George, 1994) whereas if we chose sparsistency then MIC would have infinite risk-inflation. Thus, given the choice between *better model-fit* and sparsistency, we chose the former. However, if sparsistency is more important than predictive accuracy, making a small change in the coding schemes would guarantee it.

7. A Model for “Intra Domain” Adaptation: TRANSFEAT

In the previous sections we proposed MIC based methods for the related problems of simultaneous feature selection for a set of multiple related tasks (MIC-MULTI) and grouped feature selection for single task (MIC-GROUP). The focus of those methods was joint feature selection, but in many applications it is the case that some of the tasks have less data available than other tasks and building supervised learning models from the limited amount of data does not give high predictive accuracies. So, it becomes desirable to “borrow strength” for the tasks with less amount of data from the tasks with lots of data. In other words, we want to have “intra domain” adaptation or Transfer Learning (Ando and Zhang, 2005; Raina et al., 2006).

In this section, we propose a method called TRANSFEAT which addresses the above problem by transferring information between similar tasks by using a feature relevance prior. We demonstrate the effectiveness of TRANSFEAT for the problem of Word Sense Disambiguation (WSD), and show that in this domain TRANSFEAT significantly improves accuracy on tasks with less data. TRANSFEAT, could, of course, be applied to wide variety of domains, but is particularly useful for WSD

as state-of-the-art WSD systems, including the ones that use feature selection, are strongly limited by the paucity of labeled data. For example, the training set of the SENSEVAL-2 English lexical sample task has only ~ 10 labeled examples per sense (Florian and Yarowsky, 2002). Such limited data makes it difficult to build high accuracy models using standard supervised learning techniques and suggests the use of transfer learning to improve performance.

As mentioned above, TRANSFEAT learns a feature relevance prior from “similar” tasks, and gives supervised learning accuracies which are comparable to or better than state-of-the-art WSD systems. Learning this prior for feature relevance of a test task makes those features that have been selected in the models of other “similar” tasks become more likely to be selected. TRANSFEAT does this by using a MDL-based approach similar to the MIC methods presented above.

Task Setting: We are given a set of target words each having an $n \times p$ feature matrix ($\mathbf{X}_{n \times p}$), where n is the total number of observations (instances) and p is the total number of features. We have a $n \times h$ response matrix ($\mathbf{Y}_{n \times h}$) of the h sense labels for each of the n observations. The WSD task is to assign a sense to each test instance. Note that this is a multi-class problem; we have a single task, which is to predict the correct sense of the word and we have h possible choices (the word senses) for that task. So, we approach it differently from the multi-task problem (MIC-MULTI), where we predicted all tasks jointly.

Overview of TRANSFEAT: TRANSFEAT builds upon MIC-GROUP and it has several steps:

- Break the $\mathbf{Y}_{n \times h}$ matrix into h , $n \times 1$ matrices, that is, out of one multiclass (h classes) problem we make h binary class problems. The prediction problem now becomes “Is this word sense 1 or not?”, etc. The main reason for doing this is that not all senses of all words are similar to all senses of some other word. Thus, transfer learning only makes sense at level of individual word *senses* rather than at the level of whole words.
- Make separate feature matrices for these h prediction problems, because the original feature matrix $\mathbf{X}_{n \times p}$ contained features which would be useful for the multiclass problem of “What is the exact sense of the word?”, rather than for the binary problems of “Is this sense 1 or not?” and so on. We do this is by characterizing each binary problem by those features from the original p features which are positively correlated with that particular word sense.¹⁶ This gives h feature matrices $\mathbf{X}_{\{i=1, \dots, h\}}$ drawn from the original $n \times p$ feature matrix, where each of these matrices need not have the same number of features.
- Next, cluster the different word senses by using “foreground-background” clustering that puts all singleton points into a “background cluster” which we then ignore
- Learn separate MIC-GROUP-SC models for each word sense. (Remember that as mentioned in the Section about MIC-GROUP, WSD is one problem which exhibits group structure and therefore we use it as a base model on which we build TRANSFEAT.)
- For each word sense in a cluster, use TRANSFEAT to learn a feature relevance prior from the remaining word senses in that cluster that have more observations than the target word sense, on the features of that word sense. As we explain later, this feature relevance prior allows us to learn better MIC-GROUP-SC models by relaxing the *uniform prior* assumption that each

16. In general, features with positive coefficients are associated with the given sense and those with negative coefficients with other senses of that word.

group (feature class) and then each feature within that group is equally likely to be selected, that MIC-GROUP-SC makes.

- Given these better MIC-GROUP models for all the word senses, we solve the actual h class WSD problem by choosing the sense whose model gave the highest score as the most likely sense for that word.

We learn the feature relevance prior only from distributionally similar word *senses*; in contrast to Ando (2006) who share knowledge across “all” the senses of “all” the words. Our approach makes sense as it is difficult to find words which are similar in all their senses; however, one can often find words which have one or a few similar senses. For example, one sense of “fire” (as in “fire someone”) should share features with one sense of “dismiss” (as in “dismiss someone”), but other senses of “fire” (as in “fire the gun”) do not. Similarly, other meanings of “dismiss” (as in “dismiss an idea”) should not share features with “fire”. Similarly, the words “kill”, “capture” and “arrest,” share one similar sense. This justifies our choice of breaking down the problem down to the level of individual word senses.

Thus, knowledge can only be fruitfully transferred between the shared senses of different words, even though the models being learned are for disambiguating different senses of a single word. To address this problem, we cluster similar word senses of different words, and then use the models learned for all the word senses in the cluster with more data (observations) than the held out word sense (called “training word senses”) to put a feature relevance prior on what features will be more predictive for the held out test word sense. We hold out each word sense in the cluster once and learn a prior from the remaining word senses in that cluster. For example, we can use the models for discriminating the senses of the words “kill” and the senses of “capture”, to put a prior on what features should be included in a model to disambiguate senses of the distributionally similar word “arrest”, which has considerably less data than the other two words (ONTONOTES data set), hence enabling us to learn high accuracy models for “arrest”. If at least one sense of the word “arrest”, that we are trying to model is similar to the other word senses (for “kill” and “capture”), some of the same features should be beneficial for all of them.

7.1 TRANSFEAT Formulation

We now describe TRANSFEAT in detail and show how it can be used to learn better feature selection models by relaxing the overly simplistic assumption of the model coding schemes of MIC methods of *uniform prior* by learning a feature relevance prior.

We define a binary random variable $f_i \in \{1,0\}$ that denotes the event of the i^{th} feature being in or not being in the model for the test word sense, and model it as being from a Bernoulli distribution parameterized by θ_i :

$$p(f_i|\theta_i) = \theta_i^{f_i} (1 - \theta_i)^{1-f_i}. \quad (16)$$

Given the data for the i^{th} feature for all the training word senses, we can write: $\mathcal{D}_{f_i} = \{f_{i1}, \dots, f_{iv}, \dots, f_{it}\}$. The model likelihood (under the i.i.d assumption) can be written as:

$$p(\mathcal{D}_{f_i}|\theta_i) = \prod_{v=1}^t p(f_{iv}|\theta_i) = \prod_{v=1}^t \theta_i^{f_{iv}} (1 - \theta_i)^{1-f_{iv}},$$

and the posteriors can be calculated by putting a prior over the parameters θ_i as:

$$p(\theta_i | \mathcal{D}_{f_i}) = p(\mathcal{D}_{f_i} | \theta_i) \times p(\theta_i | a, b),$$

where a and b are the hyperparameters of the Beta Prior (the conjugate of the Bernoulli distribution). The predictive distribution of θ_i is:

$$p(f_i = 1 | \mathcal{D}_{f_i}) = \int_0^1 p(f_i = 1 | \theta_i) p(\theta_i | \mathcal{D}_{f_i}) d\theta_i.$$

Substituting from 16 in the above equation we get:

$$p(f_i = 1 | \mathcal{D}_{f_i}) = \int_0^1 \theta_i p(\theta_i | \mathcal{D}_{f_i}) d\theta_i = \mathbb{E}[\theta_i | \mathcal{D}_{f_i}].$$

Using the standard results for the mean and the posterior of a Beta distribution we obtain:

$$p(f_i = 1 | \mathcal{D}_{f_i}) = \frac{k + a}{k + l + a + b}, \tag{17}$$

where k is the number of times that the i^{th} feature is selected and l is the complement of k , that is, the number of times the i^{th} feature is not selected in the training data.

As can be seen from Equation 17, the probability that a feature is selected for the held out test word sense is a “smoothed” average of the number of times it was selected in the models for the senses of other words that are similar to it.

Using similar reasoning, we can extend the above concept to the groups (feature classes) so that the probability that a group (feature class) is selected is also a “smoothed” average of the number of times it was selected in the models for the senses of other words that are similar to it.

In light of the above reasoning, the modified model cost for MIC-GROUP for coding the i^{th} feature when previously no features have been selected from the j^{th} feature class which contains that feature can be written as follows:

$$S_M^i = -\lg p(G_j = 1 | \mathcal{D}_{G_j}) - \lg p(f_i = 1 | \mathcal{D}_{f_i}) + 2,$$

and for the case when some features have already been selected from the j^{th} feature class, we can write a modified coding cost as follows:

$$S_M^i = \min [-\lg p(G_j = 1 | \mathcal{D}_{G_j}), 1 + \lg(Q)] - \lg p(f_i = 1 | \mathcal{D}_{f_i}) + 2,$$

where the first term represents the probability of selecting at least one feature from the j^{th} feature class, the second term represents the probability of selecting the i^{th} feature, and the third term which is used to code the coefficient values remains the same as earlier.¹⁷ Note that in the case when we have previously selected features from a given feature class, the most efficient way to code the feature class is to use the minimum of the TRANSFEAT cost and the actual “switch” coding cost as described in Section 4.2. Thus TRANSFEAT replaces the implicit uniform prior of MIC-GROUP with a coding scheme which is more informed by the prior learned from similar tasks.

The detailed algorithm for TRANSFEAT is given in Algorithm 3.

17. The negative sign is due to the duality between Bayesian and Information Theoretic interpretation as mentioned earlier.

Algorithm 3 TRANSFEAT

-
- 1: Break the multiclass problem into h binary prediction problems.
 - 2: Make the feature matrices for each of these problems, that is, $\mathbf{X}_{\{i=1,\dots,h\}}$.
 - 3: Cluster the different word senses by “foreground-background” clustering.
 - 4: $total_clusters = \{1, \dots, c\}$
 - 5: $word_senses_k = s_k$ // Number of word senses in k^{th} cluster.
 - 6: **for** i in $total_clusters$ **do**
 - 7: **for** t in $word_senses_i$ **do**
 - 8: Learn separate MIC-GROUP-SC models for all the word senses. // Uniform prior assumption
 - 9: **end for**
 - 10: **end for**
 - 11: **for** i in $total_clusters$ **do**
 - 12: **for** t in $word_senses_i$ **do**
 - 13: Learn TRANSFEAT model on all word senses in the cluster which have more data (observations) than the t^{th} word sense.
 - 14: Use the revised model costs S_M output by TRANSFEAT to learn better MIC-GROUP -SC model for t^{th} word sense.
 - 15: // The uniform prior assumption of MIC-GROUP-SC has been relaxed.
 - 16: **end for**
 - 17: **end for**
 - 18: Disambiguate the word as a whole by choosing the correct sense (from h possible senses) as the one whose model gave the highest score.
-

7.1.1 CHOICE OF HYPERPARAMETERS

The hyperparameters a and b in Equation 17 control the “smoothing” of our probability estimates, that is, how strongly we want the evidence obtained from similar word senses to affect the model that we learn for the test word sense.

In all our experiments we set $a = 1$ and choose b so that in the limiting case of no transfer, that is, ($k = l = 0$ in Equation 17) the coding scheme will reduce to the baseline feature selection described in (Equation 4). Thus, we choose $b = p - 1$ where p is the total number of features/feature classes (depending on what we are coding) in the test word sense.

7.2 Experimental Results

In this section we first describe our data and similarity metric that we used; we then report the results of applying TRANSFEAT to the SENSEVAL-2 and ONTONOTES data sets.

7.2.1 SIMILARITY METRIC

Finding a good similarity metric between different word senses is perhaps one of the biggest challenges that we faced. It is also the part of this section that is specific to the problem of word sense disambiguation. There are many ways in which word senses can be judged as similar, including having similar “meanings” or similar syntactic usages. Human annotated lexicons such as Levin classes (Levin, 1993), hypernyms or synonyms according to WORDNET (Miller, 1990; Lin, 1999),

or VERBNET classes (Kipper et al., 2000; Schuler, 2006) capture different aspects of this similarity, as does INFOMAP(<http://infomap.stanford.edu>) (Raina et al., 2006), which gives distributional similarity score for words in the corpus. We choose instead to define a similarity metric based, as described below on combinations of many different aspects of the lexical and syntactic context of the word.

One might think of doing K-means clustering of the word senses based on their features, but this works poorly, as it assigns all the word senses to some cluster, while in reality, there are in practice many word senses that are not sufficiently similar to any other word sense, either semantically or syntactically and hence many word senses occur in “singleton” clusters. K-means and perhaps surprisingly, hierarchical agglomerative clustering, even after extensive use of different ‘K’ or thresholds, failed to give reasonable clusters.

We thus need a clustering method that gives tight clusters of word senses, and does not attempt to cluster those word senses which are not similar to any other word sense in the corpus. We do this using a “foreground-background” clustering algorithm as proposed by Kandylas et al. (2007). This algorithm gives highly cohesive clusters of word senses (the *foreground*) and puts all the remaining word senses in the *background*. The parameters that it takes as input are the % of data points to put in *background* (i.e., what would be the singleton clusters) and a similarity threshold which impacts the number of *foreground* clusters. We experimented with putting 20% and 33% data points in background and adjusted the similarity threshold to give us 50 – 100 *foreground* clusters. The results reported below have 20% background and 50 – 100 *foreground* clusters.

7.2.2 DESCRIPTION OF DATA

We used the SENSEVAL-2 English lexical sample data, which contains a total of 73 different words (29 nouns, 29 verbs, and 15 adjectives) and the ONTONOTES verb data (the same one used for experiments of MIC-GROUP), containing 172 verbs. The main difference between these two data sets is that SENSEVAL-2 data contains “fine grained” senses of the words and as a result tends to have more senses per word than the “coarse grained” verb senses in ONTONOTES. (See Table 5.)

Data Set	#words	#train	avg #senses per word
SENSEVAL-2 (nouns+verbs+adj.)	73	8611	10.7
ONTONOTES (only verbs)	172	See Note (<i>in caption</i>)	3.7

Table 5: Data Statistics of SENSEVAL-2 and ONTONOTES data sets. *Note:* In our experiments we used the standard test-train splits for SENSEVAL-2; ONTONOTES data does not have any standard splits so we report 10-Fold cross validation test-accuracies

7.2.3 RESULTS

We cluster the word senses based on all the features, that is, semantic+syntactic similarity features. We experimented clustering using only syntactic and only semantic features but we got the best results using the combined feature set.

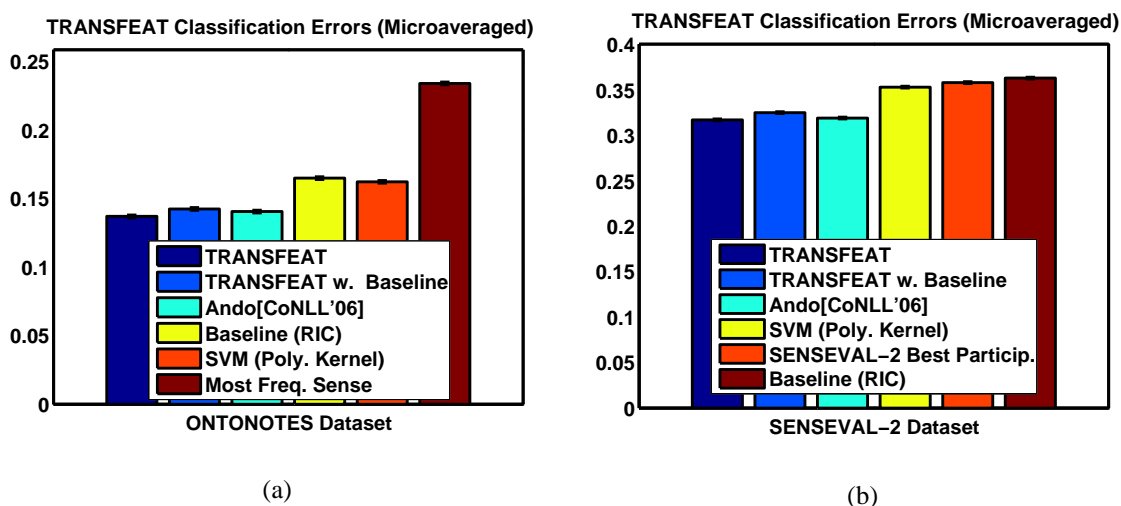


Figure 5: 10-fold CV (micro-averaged) test accuracies of various methods for ONTONOTES and SENSEVAL-2 (English Lexical Sample) data sets. *Note:* 1.) These are true cross-validation accuracies and no parameters have been tuned on them. 2.) The final accuracies reported are averaged over the entire 172 verbs. 3.) We used the standard test-train splits for SENSEVAL-2 as mentioned on the data website and as used in previous studies

All results reported are micro-averaged¹⁸ accuracies. In order to ensure fairness of comparison we compute the predicted sense for each observation by selecting the word sense model (from among the different senses for that word) with the highest score for that observation sense. As in earlier experiments we use TRANSFEAT only to select features and later we use logistic regression for classification. This “one vs all” approach to prediction in multi-class problems is widely used, although higher accuracy can sometimes be obtained by more complex pairwise comparison methods.

We use two versions of TRANSFEAT, as can be seen in Figure 5. The first version is exactly the same as mentioned in Algorithm 3, while the second version, TRANSFEAT w. Baseline, builds upon baseline feature selection (Equation 4) instead of MIC-GROUP-SC. We compare TRANSFEAT methods against baseline feature selection (Equation 4), SVM with a polynomial kernel, Ando[CoNLL'06] (Ando, 2006), computed with the standard implementation of the algorithm from the Berkeley Transfer Learning Toolkit (Rakhlin, 2007), and a simple most frequent sense baseline. For SVM we used the standard libSVM package (Chang and Lin, 2001). We used a polynomial kernel, as it gave better performance on held out data than other kernels including linear and RBF. We tuned the cost parameter ‘c’ and the degree of polynomial ‘d’ parameters of the polynomial kernel using a separate cross validation.

18. Our precision and recall are always the same as we assign exactly one sense to each instance. Hence the accuracy that we report is the same as the F-measure or ‘micro-averaged’ recall as is reported in many WSD studies.

7.2.4 ANALYSIS OF ONTONOTES RESULTS

The results for the different settings for the ONTONOTES data set averaged over the entire 172 verbs are shown in Figure 5 (a). The TRANSFEAT models are significantly better (5% significance level using a paired t-test) than all the competing methods except Ando[CoNLL'06].

Some examples will help to emphasize the point that we made earlier that transfer helps the most in cases in which the target word sense has much less data than the word senses from which knowledge is being transferred. “kill” had roughly 6 times more data than all other word senses in its cluster (i.e., “arrest”, “capture”, “strengthen”, etc.) In this case, TRANSFEAT gave 3.2 – 8.7% higher accuracies than competing methods on these three words. Both versions of TRANSFEAT do much better than Ando[CoNLL'06] on these select words even though on average over all 172 verbs the difference is slender. Similarly, for the case of word “do” which had roughly 10 times more data than the other word senses in its cluster (e.g., “die” and “save”), TRANSFEAT gave 4.1 – 6.2% higher accuracies than other methods. Transfer makes the biggest difference when the target words have much less data than the word senses they are generalizing from, but even in cases where the words have comparable amounts of data we still get a 1.5 – 2.5% increase in accuracy.

However, as one might expect, transfer learning can sometimes hurt performance; there can be so-called “negative-transfer” (Caruana, 1997). This was the case for 8 verbs out of the 172.

7.2.5 ANALYSIS OF SENSEVAL-2 RESULTS

The results for SENSEVAL-2 data set are shown in Figure 5(b). Here also TRANSFEAT does significantly better (5% significance level using a paired t-test) than the baseline feature selection method and most of the other state-of-the-art algorithms. It is worth noting that a high degree of engineering goes into the state-of-the-art SENSEVAL-2 systems. This is in contrast to TRANSFEAT, which uses information theoretic feature selection and thus has no free parameters to tune. The TRANSFEAT results are comparable to those reported in Ando (2006), which is the state-of-the-art system on SENSEVAL-2. Since Ando (2006), only mentions the overall accuracy and not the accuracy on individual words, we cannot tell whether this slender difference is statistically significant.

For words that had considerably fewer observations than other words in their cluster, TRANSFEAT again gave major benefits. For example, “begin” had ~ 8 times more data (on average per sense) than the other word senses in its cluster (i.e., “work” and “develop”). In this case, TRANSFEAT gave 6.1 – 7.1% improvement in accuracy over the baseline feature selection. Similarly, “leave” had ~ 2 times more data than “turn” and “strike”, and in this case TRANSFEAT gave 5.1 – 6.2% improvement in accuracy over the baseline. These improvements are considerably larger than the average improvement over all the words as reported in Figure 5(b).

For this data set there was negative transfer on 5 out of 73 words.

8. Conclusion

In this paper we presented a framework for learning sparse models based on the information theoretic Minimum Description Length (MDL) principle. We presented two models based on the MIC (Multiple Inclusion Criterion) which greedily select features using the MDL principle in the single and multi task settings respectively. Both the methods, MIC-MULTI and MIC-GROUP, induce two level sparsity; MIC-GROUP does feature selection at the level of groups and also at the level of features within each group and MIC-MULTI allows each selected feature to be added to the models

of some or all of the tasks. We showed how we can use MDL to specify *customized* coding schemes in scenarios where the problem has complex structure. We also discussed the conditions under which the MDL based methods are consistent and *sparsistent* and also showed that the MDL coding schemes are not arbitrary and have a corresponding Bayesian interpretation. Lastly, we proposed a model, TRANSFEAT which can be used to transfer a feature relevance prior to tasks which have less data available. We evaluated all three methods on a variety of domains including genomics (for both yeast and breast cancer) and natural language processing (Word Sense Disambiguation). Our methods are consistently at least as accurate as state-of-the-art methods, while producing models that are more sparse. Such sparseness is particularly important for applications such as genomics and computational linguistics, where interpretable models are valued.

Appendix A.

No Hypercompression Inequality:

$$\forall K > 0, P_{True} [-\lg P_{True}(X^n) \geq -\lg P_{Model}(X^n) + K] \leq 2^{-K}.$$

This inequality states that the probability of a code compressing the data by more than K bits, than the code corresponding to P_{True} is exponentially small in K , where K is any positive number.

The proof follows by using Markov's inequality and can be found in Grünwald (2007).

KL-Distinguishability: If the actual data was generated by the distribution P_{True} then the distribution P_{Model} is said to be *KL distinguishable* from P_{True} if its relative entropy (KL Divergence) from P_{Model} is greater than δ . In other words P_{Model} lies outside a δ ball of P_{True} in a relative entropy sense. More formally,

$$\forall \delta > 0, KL(P_{True} || P_{Model}) \geq \delta.$$

Universal Coding Schemes: Going back to the standard MDL setting which envisions a Sender and Receiver, assume that the Sender and Receiver have a set of candidate coding schemes \mathcal{L} for \mathcal{X}^n available. Both of them know that one of these available codes will give the highest compression for the sequence $X^n \in \mathcal{X}^n$. In other words:

$$L_{Optimal}(X^n) = \underset{i}{\operatorname{argmin}} \{L_i(X^n)\} \forall L_i \in \mathcal{L}.$$

However, they must decide on a code before the sender observes the actual data X^n and they do not know which is the best code. In the Bayesian terminology, this problem is similar to finding the classifier which has the Optimal Bayes Risk—that is, the classifier with the minimum possible risk among all the candidate classifiers. One thing that the Sender can do is on seeing the data (X^n), he encodes the data using $L_{Optimal}$ as described above. However, this is not feasible as the Receiver does not know what code the Sender used and so he would not be able to decode the message. Therefore it is not possible to find the best code that compresses the data and so in practice people use *universal codes* which compress the data almost as well as $L_{Optimal}$. This corresponds to a classifier whose risk is close to Bayes Risk. It has been shown that the two part MDL codes that we used in this paper to describe MIC based methods are universal codes (Grünwald, 2005; Grünwald, 2007; Rissanen, 1999). Moreover the “uniform prior” code and the “combinatorial code” we used in coding the model for MIC-MULTI and MIC-GROUP are also universal codes and minimax optimal (Grünwald, 2007).

Probabilistic Sources: Consider a sample space \mathcal{X} and let $\mathcal{X}^+ := \bigcup_{n \geq 1} \mathcal{X}^n$ denote the set of all possible samples of each length. Also, define $\mathcal{X}^0 = \{x^0\}$ where x^0 is a special sequence which is called empty sample. Finally $\mathcal{X}^* = \mathcal{X}^+ \cup \mathcal{X}^0$. Now, a probabilistic source with outcomes in \mathcal{X} is a function $P : \mathcal{X}^* \rightarrow [0, \infty)$ such that for all $n \geq 0$, all $x^n \in \mathcal{X}^n$ we have:

- $\sum_{z \in \mathcal{X}} P(x^n, z) = P(x^n)$ (compatibility condition)
- $P(x^0) = 1$

The two conditions say that the “event” that data (x^n, z) arrives is identical to the event that x^n arrives first and data z arrives afterward. Intuitively, probabilistic sources can be thought of as probability distributions over infinite sequences, but defining them as probability distributions over \mathcal{X}^∞ requires measure theory and the interested user can find the details in any advanced probability book.

References

- H. Akaike. Information theory and the extension of the maximum likelihood principle. In *2nd International Symposium on Information Theory, Budapest*, pages 261–281, 1973.
- R. Ando. Applying alternating structure optimization to word sense disambiguation. In *(CoNLL-X)*, 2006.
- R. Ando and T. Zhang. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853, 2005.
- A. Argyriou, T. Evgeniou, and M. Pontil. Convex multi-task feature learning. *Machine Learning*, 73(3):243–272, 2008. ISSN 0885-6125.
- F. Bach. Consistency of the group lasso and multiple kernel learning. *Journal of Machine Learning Research*, 9:1179–1225, 2008.
- F. R. Bach, G. R. G. Lanckriet, and M. I. Jordan. Multiple kernel learning, conic duality, and the smo algorithm. In *ICML*, 2004.
- A. R. Barron and T. M. Cover. Minimum complexity density estimation. *IEEE Transactions on Information Theory*, 37(4):1034–1054, 1991.
- A. R. Barron, J. Rissanen, and B. Yu. The minimum description length principle in coding and modeling. *IEEE Transactions on Information Theory*, 44(6):2743–, 1998.
- P. Bickel and K. Doksum. *Mathematical Statistics*. Prentice Hall, 2001.
- R. Caruana. Multitask learning. In *Machine Learning*, pages 41–75, 1997.
- C. C Chang and C.J. Lin. *LIBSVM: a library for support vector machines*, 2001. URL <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- J. Chen and M. Palmer. Towards robust high performance word sense disambiguation of english verbs using rich linguistic features. In *IJCNLP*, pages 933–944, 2005.

- J. Chen, A. I. Schein, L. H. Ungar, and M. Palmer. An empirical study of the behavior of active learning for word sense disambiguation. In *HLT-NAACL*, 2006.
- T. M. Cover and J. A. Thomas. *Elements of information theory*. Wiley-Interscience, New York, NY, USA, 2006.
- P. S. Dhillon and L. H. Ungar. Transfer Learning, Feature Selection and Word Sense Disambiguation. In *Annual Meeting of the Association of Computational Linguistics, (ACL)*, August 2009.
- P. S. Dhillon, D. P. Foster, and L. H. Ungar. Efficient Feature Selection in the Presence of Multiple Feature Classes. In *International Conference on Data Mining (ICDM)*, pages 779–784, 2008.
- P. S. Dhillon, B. Tomasik, D. P. Foster, and L. Ungar. Multi-Task Feature Selection Using The Multiple Inclusion Criterion (MIC). In *European Conference on Machine Learning (ECML)-PKDD*, Lecture Notes in Computer Science. Springer, September 2009.
- P. S. Dhillon, D. P. Foster, and L. Ungar. Feature selection using multiple streams. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 13, 2010.
- B. Efron, T. Hastie, L. Johnstone, and R. Tibshirani. Least angle regression. *Annals of Statistics*, 32:407–499, 2004.
- R. Florian and D. Yarowsky. Modeling consensus: classifier combination for word sense disambiguation. In *EMNLP '02*, pages 25–32, 2002.
- D. P. Foster and E. I. George. The risk inflation criterion for multiple regression. *The Annals of Statistics*, 22(4):1947–1975, 1994. ISSN 00905364.
- P. D. Grünwald. A tutorial introduction to the minimum description length principle. In *Advances in Minimum Description Length: Theory and Applications*. MIT Press, 2005.
- P. D. Grünwald. *The Minimum Description Length Principle (Adaptive Computation and Machine Learning)*. The MIT Press, 2007. ISBN 0262072815.
- E. H. Hovy, M. P. Marcus, M. Palmer, L. A. Ramshaw, and R. M. Weischedel. Ontonotes: The 90% solution. In *HLT-NAACL*, 2006.
- J. Huang, T. Zhang, and D. Metaxas. Learning with structured sparsity. In *ICML '09*, 2009.
- L. Jacob, G. Obozinski, and J-P. Vert. Group lasso with overlap and graph lasso. In *ICML '09*, 2009.
- T. Jebara. Multi-task feature and kernel selection for SVMs. In *Proceedings of the Twenty-first International Conference on Machine Learning*. ACM New York, NY, USA, 2004.
- V. Kandylas, S. P. Upham, and L. H. Ungar. Finding cohesive clusters for analyzing knowledge communities. In *ICDM*, pages 203–212, 2007.
- K. Kipper, H. T. Dang, and M. Palmer. Class-based construction of a verb lexicon. In *AAAI/IAAI*, pages 691–696, 2000.

- S.-I. Lee, V. Chatalbashev, D. Vickrey, and D. Koller. Learning a meta-level prior for feature relevance from multiple related tasks. In *ICML '07*, pages 489–496, 2007. ISBN 978-1-59593-793-3.
- B. Levin. *English Verb Classes and Alternations*. University of Chicago Press, 1993.
- D. Lin. Review of WordNet: an electronic lexical database by Christiane Fellbaum. The MIT Press 1998. *Comput. Linguist.*, 25(2):292–296, 1999. ISSN 0891-2017.
- D. Lin, E. Pitler, D. P. Foster, and L. H. Ungar. In defense of ℓ_0 . In *Workshop on Feature Selection, (ICML 2008)*, 2008.
- H. Liu and J. Zhang. On the ℓ_1 - ℓ_q regularized regression. Technical report, Carnegie Mellon University, 2008.
- K. Lounici. Sup-norm convergence rate and sign concentration property of lasso and dantzig estimators. *Electronic Journal of Statistics*, 2:90–102, 2008.
- L. Meier, S. van de Geer, and P. Bühlmann. The group lasso for logistic regression. *Journal of the Royal Statistical Society. Series B*, 70(1):53–71, 2008.
- N. Meinshausen and P. Bühlmann. High dimensional graphs and variable selection with the lasso. *Annals of Statistics*, 34:1436–1462, 2006.
- G. Miller. WordNet: An on-line lexical database. *Special Issue: International Journal of Lexicography*, 4(3), 1990.
- V. K. Mootha, C. M. Lindgren, K.-F. Eriksson, A. Subramanian, S. Sihag, J. Lehar, P. Puigserver, E. Carlsson, M. Ridderstrale, E. Laurila, N. Houstis, M. J. Daly, N. Patterson, J. P. Mesirov, T. R. Golub, P. Tamayo, B. Spiegelman, E. S. Lander, J. N. Hirschhorn, D. Altshuler, and L. C. Groop. PGC-1 α -responsive genes involved in oxidative phosphorylation are coordinately downregulated in human diabetes. *Nature Genetics*, 34:267 – 73, 2003. Datasets available at: <http://www.broad.mit.edu/gsea/datasets.jsp>.
- Y. Nardi and A. Rinaldo. On the asymptotic properties of the group lasso estimator for linear models. *Electronic Journal of Statistics*, 2:605–633, 2008.
- B. K. Natarajan. Sparse approximate solutions to linear systems. *SIAM Journal on Computing*, 24:227, 1995.
- G. Obozinski, B. Taskar, and M. I. Jordan. Joint covariate selection and joint subspace selection for multiple classification problems. *Statistics and Computing*, 2009.
- E. O. Perlstein, D. M. Ruderfer, D. C. Roberts, S. L. Schreiber, and L. Kruglyak. Genetic basis of individual differences in the response to small-molecule drugs in yeast. *Nat Genet*, 39, 2007.
- R. Raina, A. Y. Ng, and D. Koller. Constructing informative priors using transfer learning. In *ICML '06*, pages 713–720, New York, NY, USA, 2006. ACM. ISBN 1-59593-383-2.
- A. Rakhlin. Transfer learning toolkit, 2007. Software available at: <http://multitask.cs.berkeley.edu>.

- A. Rakotomamonjy, F. Bach, S. Canu, and Y. Grandvalet. SimpleMKL. *JMLR*, 9:2491–2521, 2008.
- J. Rissanen. Modeling by shortest data description. *Automatica*, 14:465–471, 1978.
- J. Rissanen. A universal prior for integers and estimation by minimum description length. *Annals of Statistics*, 11(2):416–431, 1983.
- J. Rissanen. Hypothesis selection and testing by the mdl principle. *The Computer Journal*, 42:260–269, 1999.
- K. K. Schuler. Verbnet: A broad coverage, comprehensive verb lexicon. In *Ph.D. Thesis, Computer and Information Sciences, University of Pennsylvania*, June 2006.
- G.: Schwartz. Estimating the dimensions of a model. *The Annals of Statistics*, 6(2):461–464, 1978.
- R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.
- J. A. Tropp. Greed is good: Algorithmic results for sparse approximation. *IEEE Trans. Inform. Theory*, 50:2231–2242, 2004.
- B.A. Turlach, W.N. Venables, and S.J. Wright. Simultaneous variable selection. *Technometrics*, 47(3):349–363, 2005.
- L. J. van 't Veer, H. Dai, M. J. van de Vijver, Y. D. He, A. A. Hart, M. Mao, H. L. Peterse, K. van der Kooy, M. J. Marton, A. T. Witteveen, G. J. Schreiber, R. M. Kerkhoven, C. Roberts, P. S. Linsley, R. Bernards, and S. H. Friend. Gene expression profiling predicts clinical outcome of breast cancer. *Nature*, 415(6871):530–536, January 2002. ISSN 0028-0836.
- M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (lasso). *IEEE Trans. Inf. Theor.*, 55(5):2183–2202, 2009. ISSN 0018-9448.
- Z. Wu and H. H. Zhou. Model selection and sharp asymptotic minimaxity. *Under Submission*, (-):–, 2010. ISSN -.
- M. Yuan and Y. Lin. Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, February 2006. ISSN 1369-7412.
- T. Zhang. On the convergence of mdl density estimation. In *COLT*, pages 315–330, 2004.
- T. Zhang. Adaptive forward-backward greedy algorithm for sparse learning with linear models. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems 21*, pages 1921–1928. Curran Associates, Inc., 2009a.
- T. Zhang. On the consistency of feature selection using greedy least squares regression. *Journal of Machine Learning Research (JMLR)*, 10:555–568, 2009b. ISSN 1532-4435.
- P. Zhao and B. Yu. On model selection consistency of lasso. *Journal of Machine Learning Research (JMLR)*, 7:2541–2563, 2006. ISSN 1532-4435.

- P. Zhao, G. Rocha, and B. Yu. Grouped and hierarchical model selection through composite absolute penalties. *Annals of Statistics*, 2008.
- S. Zhou. Thresholding procedures for high dimensional variable selection and statistical estimation. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems 22*, pages 2304–2312. 2009.
- H. Zou and T. Hastie. Regularization and variable selection via the elastic net. *Journal Of The Royal Statistical Society Series B*, 67(2):301–320, 2005.