

Feature selection using Stochastic Gates

Yutaro Yamada^{1*} Ofir Lindenbaum^{2*} Sahand Negahban¹
Yuval Kluger^{2,3,4†}

¹Statistics Department; ²Applied Mathematics Program;

³Computational Biology and Bioinformatics;

⁴Department of Pathology;

Yale University, New Haven, CT, USA

[†]Corresponding author. E-mail: yuval.kluger@yale.edu

Address: 333 Cedar St, New Haven, CT 06510, USA

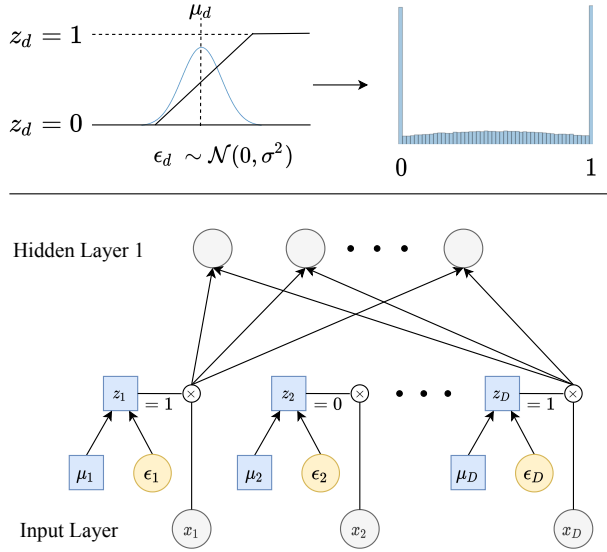
* These authors contributed equally.

Abstract

Feature selection problems have been extensively studied in the setting of linear estimation (e.g. LASSO), but less emphasis has been placed on feature selection for non-linear functions. In this study, we propose a method for feature selection in neural network estimation problems. The new procedure is based on probabilistic relaxation of the ℓ_0 norm of features, or the count of the number of selected features. Our ℓ_0 -based regularization relies on a continuous relaxation of the Bernoulli distribution; such relaxation allows our model to learn the parameters of the approximate Bernoulli distributions via gradient descent. The proposed framework simultaneously learns either a nonlinear regression or classification function while selecting a small subset of features. We provide an information-theoretic justification for incorporating Bernoulli distribution into feature selection. Furthermore, we evaluate our method using synthetic and real-life data to demonstrate that our approach outperforms other commonly used methods in both predictive performance and feature selection.

1 Introduction

Feature selection is a fundamental task in machine learning and statistics. Selecting a subset of relevant features may result in several potential benefits: reducing experimental costs [1], enhancing interpretability [2], speeding up computation, reducing memory and even improving model generalization on unseen data [3]. For example, in biomedical studies, machine learning can provide effective diagnostics or prognostics models. However, the number of features (e.g., genes or proteins) often exceeds the number of samples. In this setting, feature selection can lead to improved risk assessment and provide meaningful biological insights. While neural networks are good candidates for learning diagnostics models, identifying relevant features while building compact predictive models remains an open challenge.



Algorithm 1 STG: Feature selection using stochastic gates

Input: $\mathbf{X} \in \mathbb{R}^{N \times D}$, target variables $\mathbf{y} \in \mathbb{R}^N$, regularization parameter λ , number of epochs M , learning rate γ .

Output: Trained model f_θ and parameter $\boldsymbol{\mu} \in \mathbb{R}^D$.

- 1: Initialize the model parameter θ . Set $\boldsymbol{\mu} = 0.5$.
- 2: **for** $k = 1, \dots, K$ **do**
- 3: **for** $d = 1, \dots, D$ **do**
- 4: Sample $\epsilon_d^{(k)} \sim N(0, \sigma^2)$
- 5: Compute $z_d^{(k)} = \max(0, \min(1, \mu_d + \epsilon_d^{(k)}))$
- 6: **end for**
- 7: **end for**
- 8: Compute the loss $\hat{L} = \frac{1}{NK} \sum_{n,k} L(f_\theta(\mathbf{x}_n \odot \mathbf{z}^{(k)}), y_n)$
- 9: Compute the regularization $R = \lambda \sum_d \Phi(\frac{\mu_d}{\sigma})$
- $\boldsymbol{\theta} := \boldsymbol{\theta} - \gamma \nabla_{\boldsymbol{\theta}} \hat{L}$ and $\boldsymbol{\mu} := \boldsymbol{\mu} - \gamma \nabla_{\boldsymbol{\mu}} (\hat{L} + R)$
- 10: Repeat M epochs

Figure 1: Top left: Each stochastic gate z_d is drawn from the STG approximation of the Bernoulli distribution (shown as the blue histogram on the right). Specifically, z_d is obtained by applying the hard-sigmoid function to a mean-shifted Gaussian random variable (step 5 in algorithm 1). Bottom left: The z_d stochastic gate is attached to the x_d input feature, where the trainable parameter μ_d controls the probability of the gate being active. Right: Pseudocode of our algorithm for feature selection. See the supplementary material for a discussion of σ and λ 's selection.

Feature selection methods are classified into three major categories: filter methods, wrapper methods, and embedded methods. **Filter methods** attempt to remove irrelevant features prior to learning a model. These methods filter features using a per-feature relevance score that is created based on statistical measures [4, 5, 6, 7, 8, 9]. Wrapper methods [10, 11, 12, 13, 14] use the outcome of a model to determine the relevance of each feature. **Wrapper methods** require recomputing the model for each subset of features and, thus, become computationally expensive, especially in the context of deep neural networks [15, 16, 17]. **Embedded methods** aim to remove this burden by learning the model while simultaneously selecting the subset of relevant features. The Least Absolute Shrinkage and Selection Operator (LASSO) [18] is a well-known embedded method, whose objective is to minimize the loss while enforcing an ℓ_1 constraint on the weights of the features. LASSO is scalable and widely used [19, 20, 21], but it is restricted to the domain of linear functions and suffers from shrinkage of the model parameters. It seems natural to extend the LASSO using neural networks; however, gradient descent on an ℓ_1 regularized objective neither performs well in practice nor sparsifies the input layer [22, 23, 24].

To overcome these limitations, we develop a fully *embedded feature selection* method for nonlinear models. Our method improves upon the LASSO formulation by: a) capturing nonlinear interactions between features via neural network modeling and b) employing an ℓ_0 -like regularization using gates with weights parametrized by a smooth variant of a Bernoulli distribution. These two improvements are jointly formulated as a fully differentiable neural network that provides a solution to the important long-standing problem of feature selection for nonlinear functions.

Specifically, our contributions are as follows:

- We identify the limitations of the logistic-distribution-based Bernoulli relaxation [25, 26, 27] in feature selection and present a Gaussian-based alternative termed stochastic gate (STG), which is better in terms of model performance and consistency of feature selection.

- We develop an embedded nonlinear feature selection method by introducing the stochastic gates to the input layer (the feature space) of a neural network.
- We justify our probabilistic approach by analyzing the constrained Mutual Information maximization objective of feature selection.

We demonstrate the advantages of our method for classification, regression, and survival analysis tasks using numerous examples.

Notation: Vectors are denoted by **bold** lowercase letters \mathbf{x} and random vectors as bold uppercase letters \mathbf{X} . Scalars are denoted by lower case letters y , while random variables are uppercase Y . A set is represented by a script font \mathcal{S} . For example the n^{th} vector-valued observation is denoted as \mathbf{x}_n whereas X_d represents the d^{th} feature of the vector-valued random variable \mathbf{X} . Let $[n] = 1, 2, \dots, n$. For a set $\mathcal{S} \subset [D]$ let the vector $\mathbf{s} \in \{0, 1\}^D$ be the characteristic function for the set. That is $s_i = 1$ if $i \in \mathcal{S}$ and 0 otherwise. For two vectors \mathbf{x} and \mathbf{z} we denote $\mathbf{x} \odot \mathbf{z}$ to be the element-wise product between \mathbf{x} and \mathbf{z} . Thus, if we let $\mathbf{s} \in \{0, 1\}^D$ be the characteristic vector of \mathcal{S} , then we may define $\mathbf{x}_{\mathcal{S}} = \mathbf{x} \odot \mathbf{s}$. The ℓ_1 norm of \mathbf{x} is denoted by $\|\mathbf{x}\|_1 = \sum_{i=1}^D |x_i|$. Finally, the ℓ_0 norm of \mathbf{x} is denoted by $\|\mathbf{x}\|_0$ and counts the total number of non-zero entries in the vector \mathbf{x} .

2 Problem Setup and Background

Let $\mathcal{X} \subset \mathbb{R}^D$ be the input domain with corresponding response domain \mathcal{Y} . Given realizations from some unknown data distribution $P_{X,Y}$, the goal of embedded feature selection methods is to simultaneously select a subset of indices $\mathcal{S} \subset \{1, \dots, D\}$ and construct a model $f_{\boldsymbol{\theta}} \in \mathcal{F}$ that predicts Y based on the selected features $\mathbf{X}_{\mathcal{S}}$.

Given a loss L , the selection of features $\mathcal{S} \subset [D]$, and choice of parameters $\boldsymbol{\theta}$ can be evaluated in terms of the following risk:

$$R(\boldsymbol{\theta}, \mathbf{s}) = \mathbb{E}_{X,Y} L(f_{\boldsymbol{\theta}}(\mathbf{X} \odot \mathbf{s}), Y), \quad (1)$$

where we recall that $\mathbf{s} = \{0, 1\}^D$ is a vector of indicator variables for the set \mathcal{S} , and \odot denotes the point-wise product. Embedded feature selection methods search for parameters $\boldsymbol{\theta}$ and \mathbf{s} that minimize $R(\boldsymbol{\theta}, \mathbf{s})$ such that $\|\mathbf{s}\|_0$ is small compared to D .

2.1 Feature Selection for Linear Models

We first review the feature selection problem in the linear setting for a least squares loss. Given observations $\{\mathbf{x}_n, y_n\}_{n=1}^N$, a natural objective derived from (1) is the constrained empirical risk minimization

$$\min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^N (\boldsymbol{\theta}^T \mathbf{x}_n - y_n)^2 \quad \text{s.t.} \quad \|\boldsymbol{\theta}\|_0 \leq k. \quad (2)$$

Since the above problem is intractable, several authors replace the ℓ_0 constraint with a surrogate function, $\Omega(\boldsymbol{\theta}) : \mathbb{R}^D \rightarrow \mathbb{R}_+$, designed to penalize the number of selected features in $\boldsymbol{\theta}$. A popular choice for Ω is the ℓ_1 norm, which yields a convex problem and more precisely the LASSO optimization [18]. Computationally efficient algorithms for solving the LASSO problem have been proposed [18, 28, 29]. While the original LASSO focuses on the constrained optimization

problem, the regularized least squares formulation, which is often used in practice, yields the following minimization objective:

$$\min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{n=1}^N (\boldsymbol{\theta}^T \mathbf{x}_n - y_n)^2 + \lambda \|\boldsymbol{\theta}\|_1. \quad (3)$$

The hyperparameter λ trades off the amount of regularization versus the fit of the objective¹. The ℓ_1 -regularized method is effective for feature selection and prediction; however, it achieves this through shrinkage of the coefficients and is restricted to linear models. To avoid shrinkage, non-convex choices for Ω have been proposed [30]. As demonstrated in several studies [31, 32, 33], non-convex regularizers perform well both theoretically and empirically in prediction and feature selection.

Our goal is to develop a regularization technique that both avoids shrinkage and performs feature selection while learning a nonlinear function. To allow nonlinearities, Kernel methods have been considered [34], but scale quadratically in the number of observations. An alternative approach is to model $f_{\boldsymbol{\theta}}$ using a neural network with ℓ_1 regularization on the input weights [22, 23, 24]. However, in practice, introducing an ℓ_1 penalty into gradient descent does not sparsify the weights and requires post-training thresholding. Below, we present our method that applies a differentiable approximation of an ℓ_0 penalty on the first layer of a neural network.

3 Proposed Method

To implement an ℓ_0 regularization to either linear or nonlinear models, we introduce a probabilistic and computationally efficient neural network approach. It is well known that an exact ℓ_0 regularization is computationally expensive and intractable for high dimensions. Moreover, the ℓ_0 norm cannot be incorporated into a gradient descent based optimization. To overcome these limitations, a probabilistic formulation provides a compelling alternative. Specifically, we introduce Bernoulli gates applied to each of the d input nodes of a neural network. A random vector $\tilde{\mathbf{S}}$ represents these Bernoulli gates, whose entries are independent and satisfy $\mathbb{P}(\tilde{S}_d = 1) = \pi_d$ for $d \in [D]$, respectively. If we denote the empirical expectation over the observations as $\hat{\mathbb{E}}_{X,Y}$, then, the empirical regularized risk (Eq. 1) becomes

$$\hat{R}(\boldsymbol{\theta}, \boldsymbol{\pi}) = \hat{\mathbb{E}}_{X,Y} \mathbb{E}_{\tilde{\mathbf{S}}} \left[L(f_{\boldsymbol{\theta}}(\mathbf{X} \odot \tilde{\mathbf{S}}), Y) + \lambda \|\tilde{\mathbf{S}}\|_0 \right], \quad (4)$$

where $\mathbb{E}_{\tilde{\mathbf{S}}} \|\tilde{\mathbf{S}}\|_0$ boils down to the sum of Bernoulli parameters $\sum_{d=1}^D \pi_d$. Note that, if we constrain $\pi_d \in \{0, 1\}$, this formulation is equivalent to the constrained version of equation (1), with a regularized penalty on cardinality rather than an explicit constraint. Moreover, this probabilistic formulation converts the combinatorial search to a search over the space of Bernoulli distribution parameters (also motivated in Section 4). Thus, the problem of feature selection translates to finding $\boldsymbol{\theta}^*$ and $\boldsymbol{\pi}^*$ that minimize the empirical risk based on the formulation in Eq. 4.

Minimization of the empirical risk via gradient descent seems like a natural way to simultaneously determine the model parameters $\boldsymbol{\theta}^*$ and Bernoulli-based feature selection parameters $\boldsymbol{\pi}^*$. However, optimization of a loss function, which includes discrete random variables, suffers from high variance (see supplementary for more details and [35]). To overcome this limitation, several authors have proposed using a continuous approximation of discrete random variables, such as the Concrete [26, 25] or Hard-Concrete (HC) [27].

¹ λ has a one-to-one correspondence to k in the convex setting via Lagrangian duality.

We observed that the HC still suffers from high variance and, thus, is not suited for the task of feature selection. Therefore, we develop an empirically superior continuous distribution that is fully differentiable and implemented only to activate or deactivate the gates linking each feature (node) to the rest of the network. Our method provides an embedded feature selection algorithm with superior results in terms of both accuracy and capturing informative features compared with the state-of-the-art.

3.1 Bernoulli Continuous Relaxation for Feature Selection

Feature selection requires stability in the selected set of features. The use of logistic distributions such as the Concrete [26, 25] and HC [27] induces high variance in the approximated Bernoulli variables due to the heavy-tailedness, which often leads to inconsistency in the set of selected features. To address such limitations, we propose a Gaussian-based continuous relaxation for the Bernoulli variables \tilde{S}_d for $d \in [D]$. We refer to each relaxed Bernoulli variable as a stochastic gate (STG) defined by $z_d = \max(0, \min(1, \mu_d + \epsilon_d))$, where ϵ_d is drawn from $\mathcal{N}(0, \sigma^2)$ and σ is fixed throughout training. This approximation can be viewed as a clipped, mean-shifted, Gaussian random variable as shown in the left part of Fig. 1. Furthermore, the gradient of the objective with respect to μ_d can be computed via the chain rule, which is commonly known as the reparameterization trick [36, 37].

We can now write our objective as a minimization of the empirical risk $\hat{R}(\boldsymbol{\theta}, \boldsymbol{\mu})$:

$$\min_{\boldsymbol{\theta}, \boldsymbol{\mu}} \hat{\mathbb{E}}_{X,Y} \mathbb{E}_Z [L(f_{\boldsymbol{\theta}}(\mathbf{X} \odot \mathbf{Z}), Y) + \lambda \|\mathbf{Z}\|_0], \quad (5)$$

where \mathbf{Z} is a random vector with D independent variables z_d for $d \in [D]$. Under the continuous relaxation, the expected regularization term in the objective $\hat{R}(\boldsymbol{\theta}, \boldsymbol{\mu})$ (Eq. 5) is simply the sum of the probabilities that the gates $\{z_d\}_{d=1}^D$ are active or $\sum_{d \in [D]} \mathbb{P}(z_d > 0)$. This sum is equal to $\sum_{d=1}^D \Phi\left(\frac{\mu_d}{\sigma}\right)$, where Φ is the standard Gaussian CDF. To optimize the empirical surrogate of the objective (Eq. 5), we first differentiate it with respect to $\boldsymbol{\mu}$. This computation is done using a Monte Carlo sampling gradient estimator which gives

$$\frac{1}{K} \sum_{k=1}^K \left[L'(\mathbf{z}^{(k)}) \frac{\partial z_d^{(k)}}{\partial \mu_d} \right] + \lambda \frac{\partial}{\partial \mu_d} \Phi\left(\frac{\mu_d}{\sigma}\right),$$

where K is the number of Monte Carlo samples. Thus, we can update the parameters μ_d for $d \in [D]$ via gradient descent.

Altogether, Eq. 5 is optimized using SGD over the model parameters $\boldsymbol{\theta}$ and the parameters $\boldsymbol{\mu}$, where the latter substitute the parameters $\boldsymbol{\pi}$ in Eq. 4. See Algorithm 1 for a pseudocode of this procedure.

To remove the stochasticity from the learned gates after training, we set $\hat{z}_d = \max(0, \min(1, \mu_d))$, which informs what features are selected. In our experiments, for all synthetic datasets, we observe that the coordinates of \hat{z} converge to 0 or 1. However, when the signal is weak (e.g. the class samples are not separated) training the gates until convergence may cause overfitting of the model parameters. In these cases, setting a cutoff value (e.g. 0.5) and performing early stopping is beneficial. In the supplementary material, we discuss our choice of σ .

4 Connection to Mutual Information

In this section, we use a Mutual Information (MI) perspective to show an equivalence between a constrained ℓ_0 -based optimization for feature selection and an optimization over Bernoulli

distribution parameters.

4.1 Mutual Information based objective

From an information theoretic standpoint, the goal of feature selection is to find the subset of features \mathcal{S} that has the highest Mutual Information (MI) with the target variable Y . MI between two random variables can be defined as $I(\mathbf{X}; Y) = H(Y) - H(Y|\mathbf{X})$, where $H(Y)$ and $H(Y|\mathbf{X})$ are the entropy of $p_Y(Y)$ and the conditional entropy of $p_{Y|\mathbf{X}}(Y|\mathbf{X})$, respectively [38]. We can then formulate the task as selecting \mathcal{S} such that the mutual information between $\mathbf{X}_{\mathcal{S}}$ and Y is maximized:

$$\max_{\mathcal{S}} I(\mathbf{X}_{\mathcal{S}}; Y) \quad \text{s.t. } |\mathcal{S}| = k, \quad (6)$$

where k is the hypothesized number of relevant features.

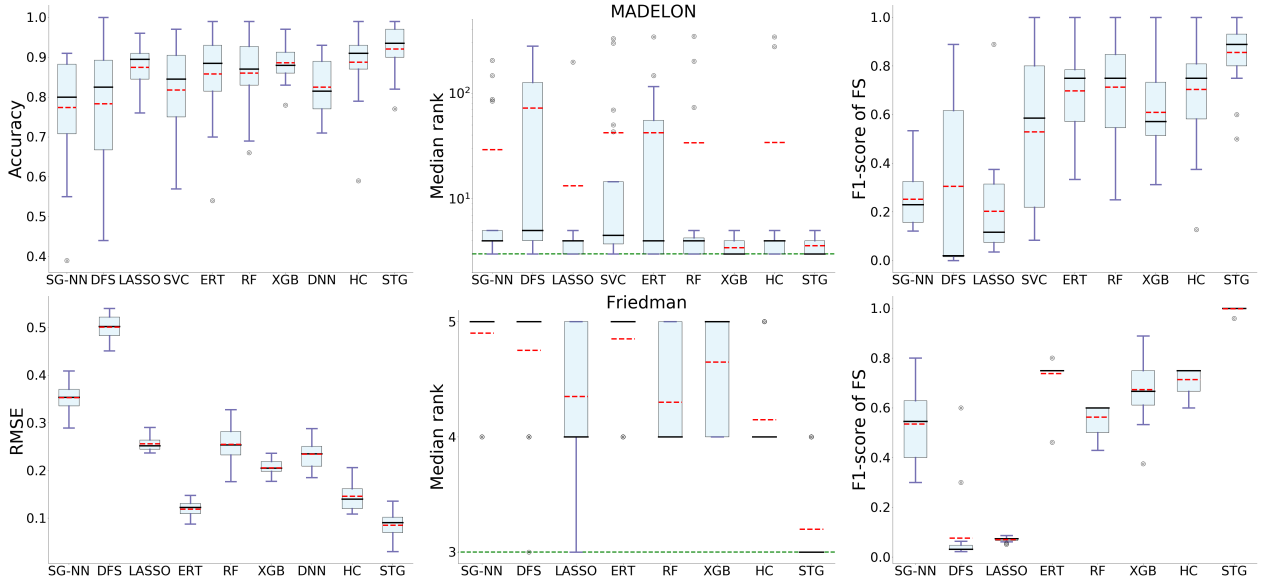


Figure 2: Evaluation of the proposed method using synthetic data. Top row: classification using the MADELON dataset with 5 informative and 495 nuisance features. Bottom row: regression using a modified version of the Friedman dataset, which also consists of 5 informative and 495 nuisance features. Left column: Accuracy/root mean squared error (RMSE). Middle column: median rank of informative features. Right column: F1-score that measures success in retrieving the informative features. We also evaluated the accuracy and MSE of a neural network with no feature selection (DNN). Black bars represent the medians, and dashed red lines are the means. In the middle column, dashed green lines are the optimal median ranks.

4.2 Introducing randomness

Under mild assumptions, we show that one can replace the deterministic search over the set \mathcal{S} (or corresponding indicator vector \mathbf{s}) with a search over the parameters of the distributions that model \mathbf{s} . Our proposition is based on the following two assumptions:

Assumption 1: There exists a subset of indices \mathcal{S}^* with a cardinality equal to k such that for any $i \in \mathcal{S}^*$ we have $I(X_i; Y|\mathbf{X}_{\setminus\{i\}}) > 0$.

Assumption 2: $I(\mathbf{X}_{\mathcal{S}^{*c}}; Y | \mathbf{X}_{\mathcal{S}^*}) = 0$.

Discussion of assumptions: Assumption 1 states that including an element from \mathcal{S}^* improves prediction accuracy. This assumption is equivalent to stating that feature i is strongly relevant [39, 40]. Assumption 2 simply states that \mathcal{S}^* is a superset of the Markov Blanket of the variable Y [40]. The assumptions are quite benign. For instance, they are satisfied if \mathbf{X} is drawn from a Gaussian with a non-degenerate covariance matrix and $Y = f(\mathbf{X}_{\mathcal{S}^*}) + w$, where w is noise independent of \mathbf{X} and f is not degenerate. With these assumptions in hand, we may present our result.

Proposition 1. *Suppose that the above assumptions hold for the model. Then, solving the optimization (6) is equivalent to solving the optimization*

$$\max_{\mathbf{0} \leq \boldsymbol{\pi} \leq \mathbf{1}} I(\mathbf{X} \odot \tilde{\mathbf{S}}; \mathbf{Y}) \quad \text{s.t.} \quad \sum_i \mathbb{E}[\tilde{S}_i] \leq k, \quad (7)$$

where the coordinates \tilde{S}_i are drawn independently at random from a Bernoulli distribution with parameter π_i .

Due to length constraints, we leave the proof of this proposition and how it bridges the MI maximization (6) and risk minimization (2) to the supplementary material.

5 Related Work

The three most related works to this study are [27], [41] and [42]. In [27], they introduce the Hard-Concrete (HC) distribution as a continuous surrogate for Bernoulli distributions in the context of model compression. The authors demonstrate that applying the HC to all of the weights leads to fast convergence and improved generalization. They did not evaluate the HC for the task of feature selection where stability of the selection is an important property.

In [41], the Concrete distribution is used to develop a framework for interpreting pre-trained models. Their method is focused on finding a subset of features given a particular sample and, therefore, is not appropriate for general feature selection. In [42], the Concrete distribution is used for feature ranking. The method is not fully embedded and requires model retraining to achieve feature selection.

Bernoulli relaxation techniques that are based on logistic distributions (e.g. Concrete/Gumbel-Softmax and HC) are not suitable for feature selection. Specifically, the use of the Concrete/Gumbel-Softmax distribution ranks features but retains all of them (no feature selection). In contrast to our Gaussian-based relaxation of Bernoulli distributions (STG), the logistic-based HC yields high-variance gradient estimates. For model sparsification, this high variance is not problematic because the sparsity pattern within the network does not matter as long as the method achieves enough sparsity as a whole. For feature selection based on the HC approach, however, the subsets of selected features at different runs vary substantially. Thus, the stability of the HC-based feature selection is poor; see Section 8. Furthermore, higher gradient variance will also result in a slower SGD convergence, which has been demonstrated empirically in Section 6 and the supplementary material.

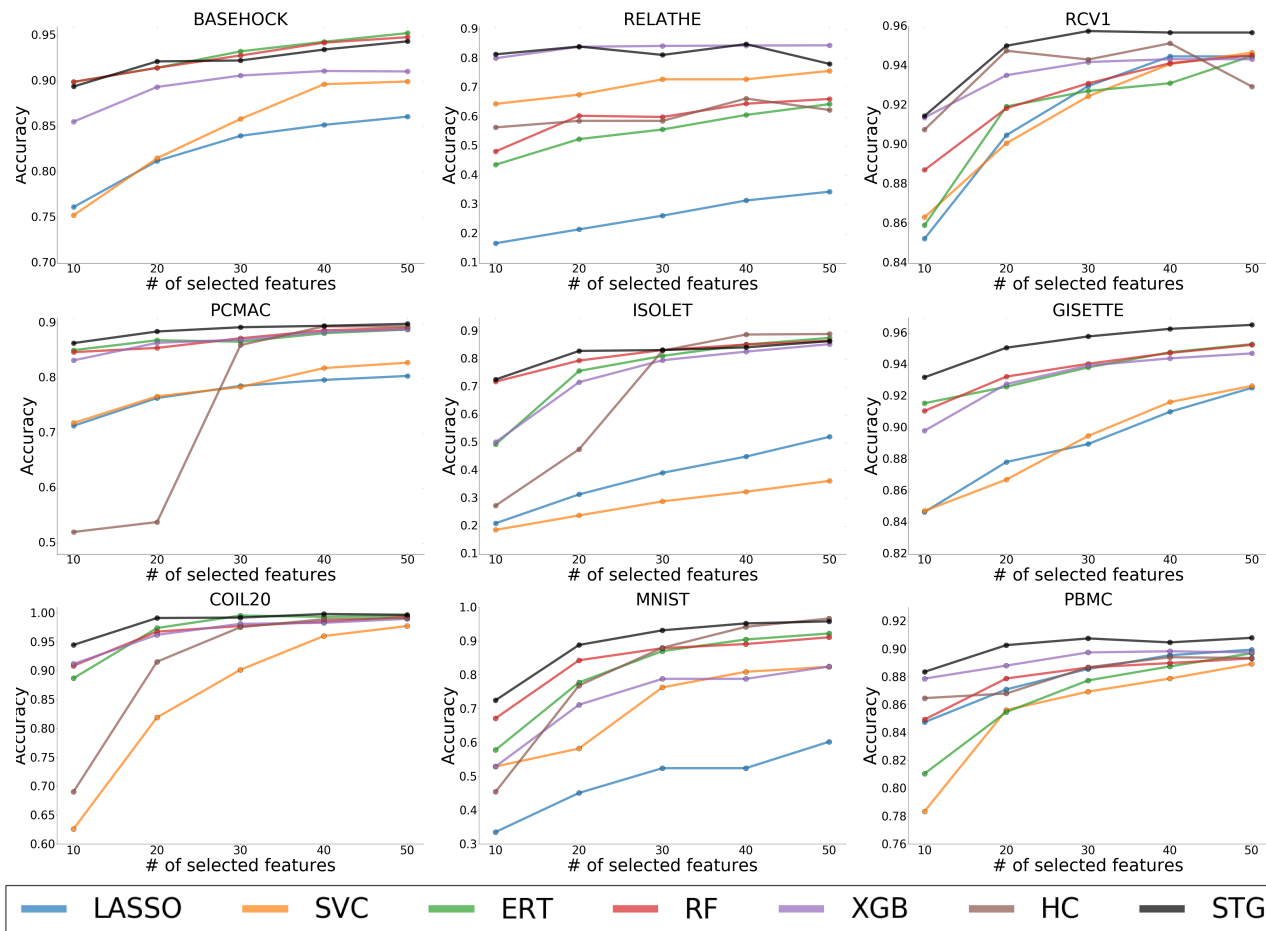


Figure 3: Classification accuracy vs. number of selected features. Descriptions of the 9 datasets appear in Table 1.

6 Experiments

Here, we evaluate our proposed embedded feature selection method. We implemented ² it using both the STG and HC [27] distributions and tested on several artificial and real datasets. We compare our method with several classification and regression algorithms including embedded methods such as LASSO [18], linear support vector classification (SVC) [43], deep feature selection (DFS) [22] and group-sparse regularization for deep neural networks (SG-NN) [23]. Our method is also compared with leading tree-based wrapper methods - extremely randomized trees (ERT) [44], random forests (RF) [45] and XGBoost (XGB) [46]. See the supplementary material for details on the hyper-parameters of all methods.

6.1 Synthetic data in the $D > N$ regime

Now we present empirical results in the challenging regime where the number of features exceeds the number of samples ($D > N$). We use synthetic datasets with informative and irrelevant nuisance variables. We begin with the MADELON dataset, a hard classification problem suggested in the NIPS 2003 feature selection challenge [47]. This dataset consists of 5 informative features and 495 nuisance features. See the supplementary material for additional details on the MADELON dataset.

²<https://github.com/runopti/stg>

Next, we present a regression scenario based on a modification of the Friedman regression dataset [48]. In this dataset, all 500 variables are uniformly distributed in $[0, 1]$, and the response is defined by the following function:

$$Y = 10 \sin(X_1 X_2)^2 + 20 X_3^2 + 10 \operatorname{sign}(X_4 X_5 - 0.2) + \xi,$$

where ξ is drawn from $\mathcal{N}(0, 1)$. Then, Y is centered and divided by its maximal value.

For the above synthetic classification and regression datasets, we generate 600 samples of which we use 450 for training, 50 for validation and 100 for a test set. The hyper-parameter (controlling the number of selected features) for each method is optimized based on the validation performance. The experiment is repeated 20 times, and the accuracy/root mean squared error (RMSE), median rank and F1-score that measures feature selection performance are presented in Fig. 2. To compute the median ranks, we utilize the scores that each method assigns to the features. We then rank all features based on these scores and compute the median of the ranks of the 5 informative features. Thus, the optimal median rank in these examples is 3. The F1-score measure for feature selection is defined as $F1 = 2(\text{precision} \cdot \text{recall}) / (\text{precision} + \text{recall})$, where precision and recall are computed by comparing the selected/removed features to the informative/nuisance features. For example, a model which retains all of the features has a recall of 1 with a precision of $5/500$.

The results presented in Fig. 2 demonstrate our embedded method’s ability to learn a powerful predictive model in the regime of $D > N$, where the majority of variables are not informative. Even though our median rank performance is comparable to tree-based methods, we outperform all baseline in F1-scores. This demonstrates that our embedded method is a strong candidate for the task of finding complex relations in high dimensional data.

Table 1: Description of the real-world data used for empirical evaluation.

	BASEHOCK	RELATHE	RCV1	PCMAC	ISOLET	GISETTE	COIL20	MNIST	PBMC
Features (D)	7862	4322	47236	3289	617	5000	1024	784	17126
Train size	1594	1141	2320	1554	1248	5600	1152	60000	2074
Test size	398	285	20882	388	312	1400	288	10000	18666
Classes	2	2	2	2	26	2	20	10	2
Data type	Text	Text	Text	Text	Audio	Image	Image	Image	scRNA-seq

We encourage the reader to look at the supplementary material, where we provide additional experiments including a more challenging variant of the MADELON data, the XOR data, the two-moons data and 3 artificial regression datasets based on [52].

6.2 Noisy Binary XOR Classification

In the following evaluation, we consider the problem of learning a binary XOR function for classification task. The first two coordinates x_1, x_2 are drawn from a "fair" Bernoulli distribution. The response variable is set as an XOR of the first coordinates, such that $y = x_1 \oplus x_2$. The coordinates $x_i, i = 3, \dots, D$ are nuisance features, also drawn from a binary "fair" Bernoulli distribution. The number of points we generate is $N = 1,500$, of which 70 % are reserved for test and 10% of the remaining training set was reserved for validation. We compare the proposed method to four embedded feature selection methods (LASSO [18], C-support vectors (SVC) [49], deep feature selection (DFS) [22], sparse group regularized NN (SG-L1-NN) [23]). To provide more benchmarks, we also compare our embedded method against three wrapper methods (Extremely Randomized Trees (Tree) [44], Random Forests (RF) [50]) and Extreme Gradient Boosting (XGBOOST) [46].

To evaluate the feature selection performance, we calculate the Informative Features Weight Ratio (IFWR). IFWR is defined as the sum of weights W_d over the informative features divided by the sum over all weights. In the case of binary weights the IFWR is in fact a recall measure for the relevant features.

The experiment is repeated 20 times for different values of D , and the average test classification accuracy and standard deviation are presented in Fig. 4(a), followed by the IFWR in Fig. 4(b). The number of selected features affects the accuracy. Therefore, to treat all the methods in a fair manner, we tune the hyperparameter that controls the sparsity level using Optuna [51] which optimizes the overall accuracy across different D s. For instance, the wrapper methods (Tree, RF and XGBOOST) has a threshold value to retain features. We retrain them using only such features whose weight is higher than the threshold. In terms of feature ranking (see Fig. 4(c)), only the tree based methods and the proposed (STG and HC based) provide the optimal median rank (which is 1.5) for the two relevant features. Moreover, the ranking provided by STG is the most stable comparing to all the alternative methods.

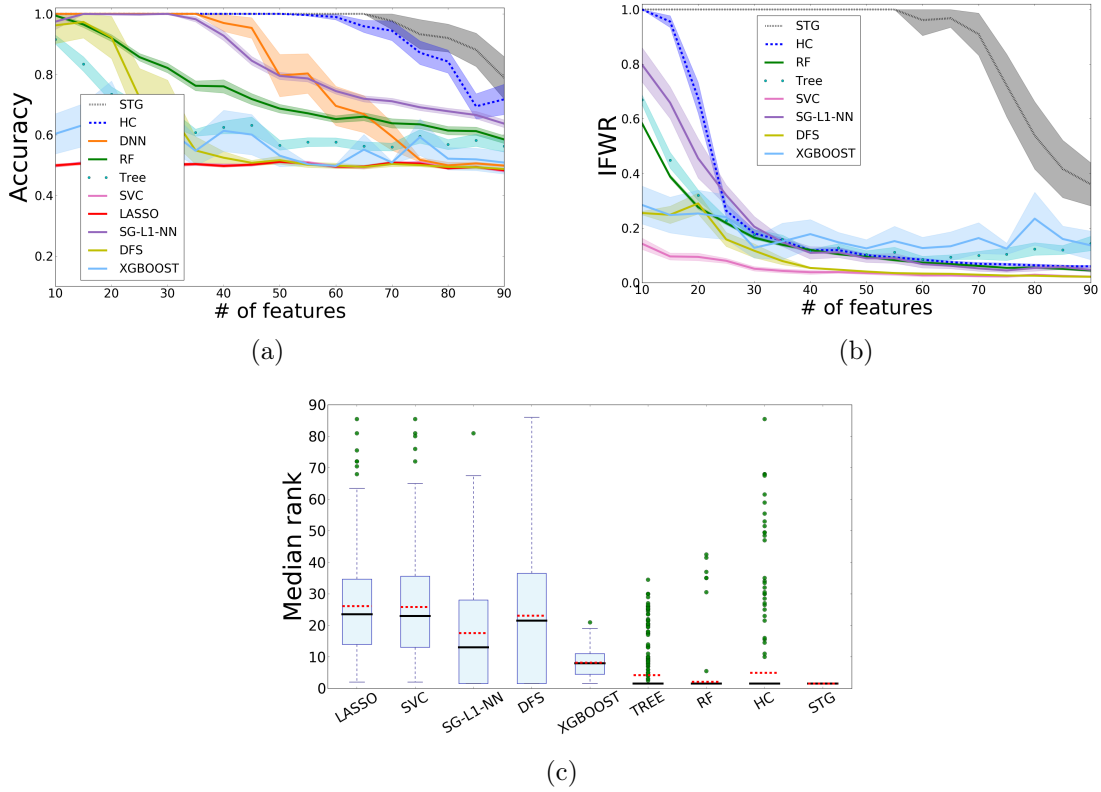


Figure 4: (a) Classification accuracy (mean and standard deviation) vs. the number of irrelevant noisy dimension (D) for the XOR problem. (b) The Informative Features Weight Ratio (IFWR), central lines are the means while shaded are represent the standard deviation. IFWR is the sum of weights W_d over the informative features divided by the sum over all weights. (c) Box plots for the median rank of the two informative features. Black line and dashed red represent the median and mean of each method. Optimal median rank in this experiment is 1.5.

6.3 Classification on real world data

We now turn our attention towards using several real-world labeled datasets to evaluate our method. Most of the datasets are collected from the ASU feature selection database available online ³. The dimensions, sample size and domains of the data are versatile and are detailed in Table 1. On all of the datasets - except MNIST [53], RCV-1 [54] and the PBMC [55] - we perform 5-fold cross validation and report the average accuracies vs. the number of features that the model uses for several baselines. For MNIST, RCV-1 and the PBMC - we used prefixed training and testing sets which are described in the supplementary material. The results are presented in Fig. 3.

Compared to the alternative linear embedded methods (i.e. LASSO and SVC), the non-linearity of our method provides a clear advantage. While there are regimes in which the tree-based methods slightly outperform our method, they require retraining a model based on the selected features; however, our method is only trained once and learns the model and features simultaneously. This experiment also demonstrates that the STG is more suited for the task of feature selection than the HC. Note that in this experiment we did not include the DFS and SG-NN, as they do not sparsify the weights and, therefore, cannot be evaluated vs. the number of selected features.

Due to lack of space we leave the real word regression experiments to the supplementary material.

7 Cox Proportional Hazard Models for Survival Analysis

A standard model for survival analysis is the Cox Proportional Hazard Model. In [56], the authors proposed DeepSurv that extends the Cox model to neural networks. We incorporate our method into DeepSurv to see how our procedure improves survival analysis based on gene expression profiles from the breast cancer dataset called METABRIC [57] (along with additional commonly used clinical variables.) See the supplementary material for more details about the dataset and experimental setup.

We compare our method Cox-STG with four other methods: a Cox model with ℓ_1 regularization (Cox-LASSO), Random Survival Forest (RSF) [58], Cox-HC, and the original DeepSurv. We evaluate the predictive ability of the learned models based on the concordance index (CI), a standard performance metric for model assessment in survival analysis; it measures the agreement between the rankings of the predicted and observed survival times. The performance of each model in terms of the CI and the number of selected features are reported in Table 2. The Cox-STG method outperforms the other baselines indicating that our approach identifies a small number of informative variables while maintaining high predictive performance.

Table 2: Performance comparison of survival analysis on METABRIC. We repeat the experiment 5 times with different training/testing splits and report the mean and standard deviation on the testing set.

	DEEPSURV	RSF	COX-LASSO	COX-HC	COX-STG
C-INDEX	0.612 (0.009)	0.626 (0.006)	0.580 (0.003)	0.615 (0.007)	0.633 (0.005)
# FEATURES	221 (ALL)	221 (ALL)	44 (0)	14 (1.72)	2 (0)

³<http://featureselection.asu.edu/datasets.php>

8 Evaluating stochastic regularization schemes

In this section, we elaborate on two aspects of our proposed method that lead to performance gains: (i) benefits of our non-convex regularization and injected noise, and (ii) advantages of the Gaussian based STG over the logistic based HC distribution in terms of feature selection performance.

To demonstrate these performance gains, we perform a controlled experiment in a linear regression setting. We first generate the data matrix, $\mathbf{X} \in \mathbb{R}^{N \times D}$, $D = 64$, with values randomly drawn from $\mathcal{N}(0, 1)$ and construct the response variable

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}^* + \mathbf{w}, \quad (8)$$

where the values of the noise $\mathbf{w}_i, i = 1, \dots, N$ are drawn independently from $\mathcal{N}(0, 0.5)$. As suggested by [59], we use a known sparsity $\|\boldsymbol{\beta}^*\|_0 = k$, set by $k = \lceil 0.4D^{0.75} \rceil = 10$. For each number of samples N in the range $[10, 250]$, we run 200 simulations and count the portion of correctly recovered informative features (i.e. the support of $\boldsymbol{\beta}^*$). For LASSO, the regularization parameter was set to its optimal value $\alpha_N = \sqrt{\frac{2\sigma^2 \log(D-k) \log(k)}{N}}$ [59]. For STG and HC, we set $\lambda_N = C\alpha_N$, such that C is a constant selected using a cross validated grid search in the range $[0.1, 10]$. To evaluate the effect of non-convex regularization and noise injection, we compare the STG to a deterministic non-convex (DNC) counterpart of our method (see definition below) and LASSO, which is convex. To gain insights on (ii), we also compare the HC.

We define the deterministic non-convex (DNC) objective as

$$\min_{\boldsymbol{\theta}, \boldsymbol{\mu}} \frac{1}{N} \sum_{n=1}^N (\boldsymbol{\theta}^T \mathbf{x}_n \odot \tilde{\mathbf{z}} - y_n)^2 + \lambda \sum_{d=1}^D \Phi\left(\frac{\mu_d}{0.5}\right), \quad (9)$$

where Φ is the standard Gaussian CDF. Combined with \tilde{z}_d , this non-convex regularized objective is deterministic and differentiable, and its solution can be searched via gradient descent.

As demonstrated in Fig. 5, the non-convex formulation requires less samples for perfect recovery of informative features than the LASSO. The injected noise based on the HC and STG provides a further improvement. Finally, we observe that the STG is more stable and has a lower variance than the HC, as shown by the shaded colors.

Application of the deterministic formulation is associated with a phenomenon that causes the gradient of an input feature to vanish and *never* acquire a nonzero value if it is zeroed out in an early training phase. In contrast, when we apply STG, a feature that at a certain step has a zero value is not permanently locked because the gate associated with it may change its value from zero to one at a later phase during training. This is due to the injected noise that allows our proposed method to *reevaluate* the gradient of each gate. In Fig. 5, we demonstrate this ‘‘second chance’’ effect using $N = 60$ samples and presenting the gate’s values (throughout training) for an active feature.

The advantage of the Gaussian-based STG distribution over the HC distribution stems from the heavier tail of HC, whose form is a logistic distribution. We demonstrate that the heavy-tail distribution is not suitable for feature selection due to its high variance. An ideal feature selection algorithm is expected to identify a consistent set of features across different runs (feature stability), but HC selects many different features in each run resulting in high variance or lack of stability of the selected features.

To further examine the effect of heavy tail distributions, we train two identical neural networks on MNIST but use two different distributions for the gates: Gaussian-STG and HC. Both regularization parameters are tuned to retain 6 features. In Fig. 6, we show that the

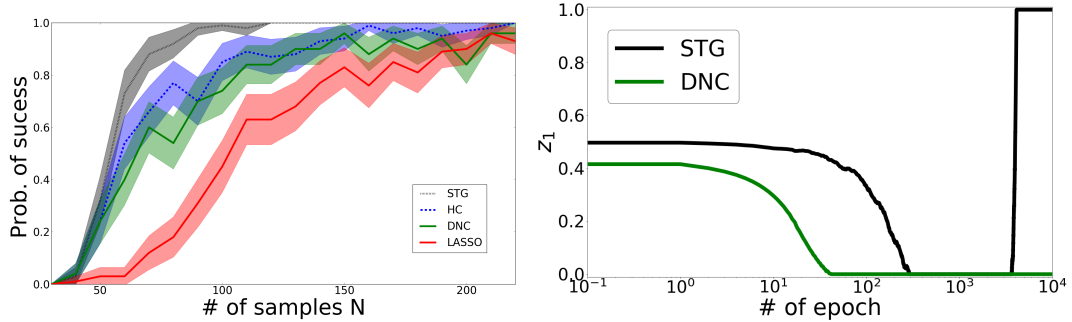


Figure 5: Feature selection in linear regression (see Section 8). The goal is to identify the subset of informative features. Top: Probability of recovering the informative features as a function of the number of samples. Comparison between STG, HC, LASSO and DNC. Bottom: The value of a gate z_1 throughout training. In STG, injected noise may lead to a “second chance” effect, which in this example occurs after 4000 epochs (black line). In the deterministic DNC setting (green line), a feature’s elimination causes its gradient to vanish for the rest of the training.

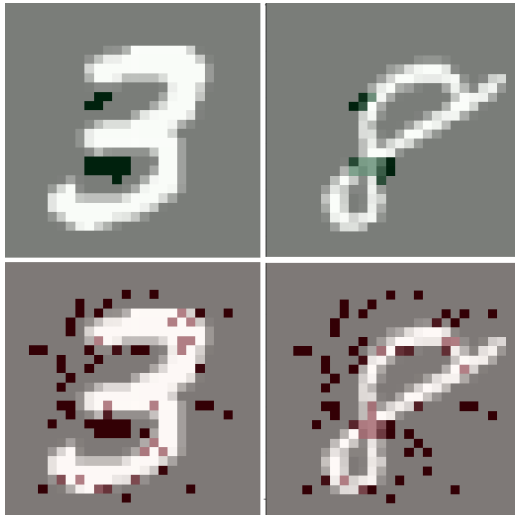


Figure 6: Comparing stability of feature selection by STG and HC. We train our method to classify 3s and 8s (from MNIST) with a regularization parameter tuned to retain ~ 6 features. We repeat the experiment using 20 random initializations. Dark pixels represent the union of selected features based on the STG (top) and HC (bottom) overlaid on top of two randomly sampled examples from each class. This demonstrates that an HC-based feature selection does not provide a stable selection of features across different runs.

selected features from the Gaussian-STG are much more consistent across 20 runs than HC. Furthermore, the variance in the number of selected features is 3.8 for HC and 1 for STG. The average accuracies of STG and HC on the test are comparable: 92.4% and 91.7%, respectively.

9 Feature Selection with Correlations

Lastly, we evaluate our proposed method using data with correlated features. In real-world, high-dimensional datasets, many features are correlated. Such correlations introduce a challenge for feature selection. For instance, in the most extreme case if there are copies of the same feature, then it is not clear which to select. As another example consider if a large subset of

features are a function of a small subset of features that we wish to identify. That large subset of seemingly useful features can confound a feature selection method. Below, we consider a number of examples in various correlated feature settings and demonstrate the strong performance of STG.

We first evaluate the proposed method in a linear setting. To introduce correlated features, we extend the linear regression experiment described in Section 8 using a correlated design matrix with a covariance matrix whose values are defined by $\Sigma_{i,j} = 0.3^{|i-j|}$. We run 100 simulations and present the probability of recovering the correct support of β^* . Fig. 7 shows that even if the features are correlated, STG successfully recovers the support with fewer samples than HC, DNC, and LASSO.

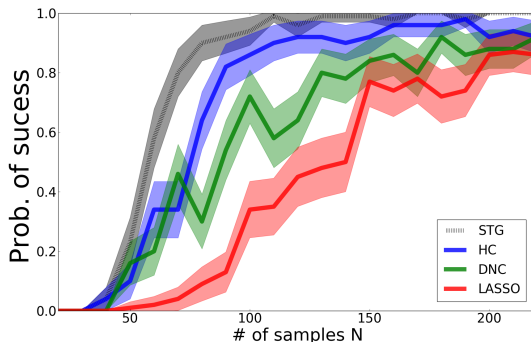


Figure 7: Feature selection in linear regression using a correlated design matrix. Probability of recovering the informative features as a function of the number of samples. Comparison between STG, HC, LASSO and DNC.

Next we evaluate our method in a non-linear setting using a variant of the MADELON dataset, which includes correlated features. Following [47], the first 5 informative features of MADELON are used to create 15 additional coordinates based on a random linear combination of the first 5. A Gaussian noise $\mathcal{N}(0, 1)$ is injected to each feature. Next, additional 480 nuisance coordinates drawn from $\mathcal{N}(0, 1)$ are added. Finally, 1% of the labels are flipped.⁴ We use 1,500 points from this dataset and evaluate the ability of STG to detect the informative features.

Fig. 8(a) shows the precision of feature selection (black line) and the number of selected features (red line) as a function of the regularization parameter λ in the range $[0.01, 10]$. We observe that there is a wide range of λ values in which our method selects only relevant features (i.e. the precision is 1). Furthermore, there is a wide range of λ values in which 5 features are selected consistently.

Next, we compare classification accuracy to Random Forest and LASSO using a 5-fold cross validation. As evident from Fig. 8(b), STG achieves the highest accuracy while using less features. Moreover, as depicted from this figure, peak performance occurs when selecting 5 features; thus, STG provides a clear indication to the true number of informative features. Both LASSO and RF on the other hand, do not provide a clear indication of the true number of relevant features.

10 Conclusion

In this paper, we propose a novel embedded feature selection method based on stochastic gates. It has an advantage over previous ℓ_1 regularization based methods in its ability to achieve a high level of sparsity in nonlinear models such as neural networks, without hurting performance.

⁴generated using `dataset.make_classification` from `scikit-learn` (<http://scikit-learn.org/>)

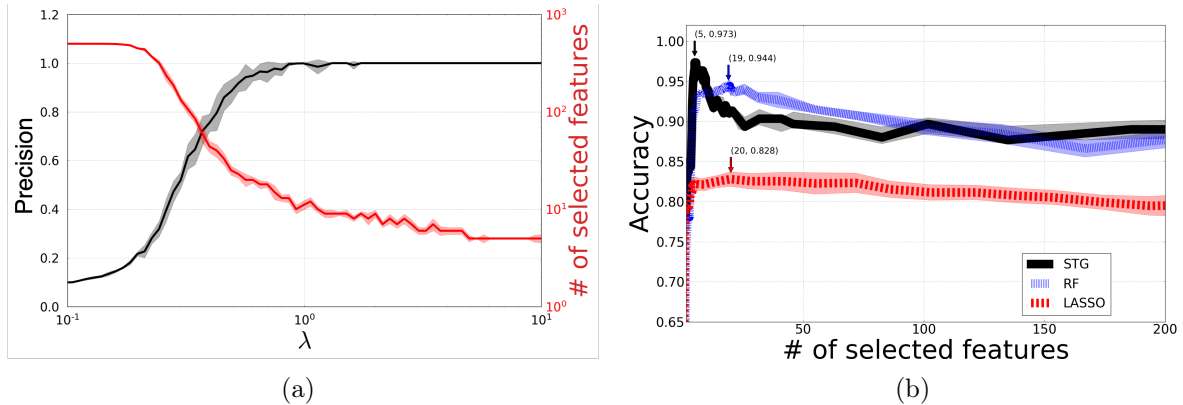


Figure 8: (a) An empirical evaluation of the effect the regularization parameter λ has on the precision of feature selection (black line) and the number of selected features (red line). The precision and the number of selected features are presented on the left and right side of the y -axis, respectively. The means are displayed as solid lines while the standard deviations are marked as shaded regions around the means. (b) Classification accuracy on the MADELON datasets. We evaluate performance using 5-fold cross validation for different number of selected features. In this dataset, the first 5 coordinates are strongly relevant features and the next 15 are weakly relevant. In that regime the proposed method outperforms RF and LASSO.

We justify our probabilistic feature selection framework from an information theoretic perspective. In experiments, we demonstrate that our method consistently outperforms existing embedded feature selection methods in both synthetic datasets and real datasets.

Acknowledgements

The authors thank Nicolas Casey and the anonymous reviewers for their helpful feedback. This work was supported by the National Institutes of Health [R01GM131642, R01HG008383, P50CA121974 and R61DA047037], National Science Foundation DMS 1723128, and the Funai Overseas Scholarship to YY.

References

- [1] F. Min, Q. Hu, and W. Zhu, “Feature selection with test cost constraint,” *International Journal of Approximate Reasoning*, vol. 55, no. 1, pp. 167–179, 2014.
- [2] M. T. Ribeiro, S. Singh, and C. Guestrin, ““why should i trust you?": Explaining the predictions of any classifier,” in *Proceedings of the 22Nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD ’16. New York, NY, USA: ACM, 2016, pp. 1135–1144. [Online]. Available: <http://doi.acm.org/10.1145/2939672.2939778>
- [3] G. Chandrashekar and F. Sahin, “A survey on feature selection methods,” *Computers & Electrical Engineering*, vol. 40, no. 1, pp. 16–28, 2014.
- [4] R. Battiti, “Using mutual information for selecting features in supervised neural net learning,” *IEEE Transactions on neural networks*, vol. 5, no. 4, pp. 537–550, 1994.

- [5] H. Peng, F. Long, and C. Ding, “Feature selection based on mutual information criteria of max-dependency, max-relevance, and min-redundancy,” *IEEE Transactions on pattern analysis and machine intelligence*, vol. 27, no. 8, pp. 1226–1238, 2005.
- [6] P. A. Estévez, M. Tesmer, C. A. Perez, and J. M. Zurada, “Normalized mutual information feature selection,” *IEEE Transactions on Neural Networks*, vol. 20, no. 2, pp. 189–201, 2009.
- [7] L. Song, A. Smola, A. Gretton, K. M. Borgwardt, and J. Bedo, “Supervised feature selection via dependence estimation,” in *Proceedings of the 24th international conference on Machine learning*. ACM, 2007, pp. 823–830.
- [8] L. Song, A. Smola, A. Gretton, J. Bedo, and K. Borgwardt, “Feature selection via dependence maximization,” *Journal of Machine Learning Research*, vol. 13, no. May, pp. 1393–1434, 2012.
- [9] J. Chen, M. Stern, M. J. Wainwright, and M. I. Jordan, “Kernel feature selection via conditional covariance minimization,” in *Advances in Neural Information Processing Systems*, 2017, pp. 6946–6955.
- [10] R. Kohavi and G. H. John, “Wrappers for feature subset selection,” *Artificial intelligence*, vol. 97, no. 1-2, pp. 273–324, 1997.
- [11] G. Stein, B. Chen, A. S. Wu, and K. A. Hua, “Decision tree classifier for network intrusion detection with ga-based feature selection,” in *Proceedings of the 43rd annual Southeast regional conference-Volume 2*. ACM, 2005, pp. 136–141.
- [12] Z. Zhu, Y.-S. Ong, and M. Dash, “Wrapper–filter feature selection algorithm using a memetic framework,” *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, vol. 37, no. 1, pp. 70–76, 2007.
- [13] J. Reunanen, “Overfitting in making comparisons between variable selection methods,” *Journal of Machine Learning Research*, vol. 3, no. Mar, pp. 1371–1382, 2003.
- [14] G. I. Allen, “Automatic feature selection via weighted kernels and regularization,” *Journal of Computational and Graphical Statistics*, vol. 22, no. 2, pp. 284–299, 2013.
- [15] A. Verikas and M. Bacauskiene, “Feature selection with neural networks,” *Pattern Recognition Letters*, vol. 23, no. 11, pp. 1323–1335, 2002.
- [16] M. M. Kabir, M. M. Islam, and K. Murase, “A new wrapper feature selection approach using neural network,” *Neurocomputing*, vol. 73, no. 16-18, pp. 3273–3283, 2010.
- [17] D. Roy, K. S. R. Murty, and C. K. Mohan, “Feature selection using deep neural networks,” in *Neural Networks (IJCNN), 2015 International Joint Conference on*. IEEE, 2015, pp. 1–6.
- [18] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288, 1996.
- [19] C. Hans, “Bayesian lasso regression,” *Biometrika*, vol. 96, no. 4, pp. 835–845, 2009.

- [20] W. Li, J. Feng, and T. Jiang, “Isolasso: a lasso regression approach to rna-seq based transcriptome assembly,” in *International Conference on Research in Computational Molecular Biology*. Springer, 2011, pp. 168–188.
- [21] F. Li, Y. Yang, and E. P. Xing, “From lasso regression to feature vector machine,” in *Advances in Neural Information Processing Systems*, 2006, pp. 779–786.
- [22] Y. Li, C.-Y. Chen, and W. W. Wasserman, “Deep feature selection: theory and application to identify enhancers and promoters,” *Journal of Computational Biology*, vol. 23, no. 5, pp. 322–336, 2016.
- [23] S. Scardapane, D. Comminiello, A. Hussain, and A. Uncini, “Group sparse regularization for deep neural networks,” *Neurocomput.*, vol. 241, no. C, pp. 81–89, Jun. 2017. [Online]. Available: <https://doi.org/10.1016/j.neucom.2017.02.029>
- [24] J. Feng and N. Simon, “Sparse-Input Neural Networks for High-dimensional Nonparametric Regression and Classification,” *ArXiv e-prints*, Nov. 2017.
- [25] C. J. Maddison, A. Mnih, and Y. W. Teh, “The concrete distribution: A continuous relaxation of discrete random variables,” *CoRR*, vol. abs/1611.00712, 2016. [Online]. Available: <http://arxiv.org/abs/1611.00712>
- [26] E. Jang, S. Gu, and B. Poole, “Categorical reparameterization with gumbel-softmax,” 2017. [Online]. Available: <https://arxiv.org/abs/1611.01144>
- [27] C. Louizos, M. Welling, and D. P. Kingma, “Learning sparse neural networks through l0 regularization,” *CoRR*, vol. abs/1712.01312, 2017.
- [28] Y. Nesterov, “Gradient methods for minimizing composite functions,” *Mathematical Programming*, vol. 140, no. 1, pp. 125–161, 2013.
- [29] J. Qian, W. Du, Y. Tanigawa, M. Aguirre, R. Tibshirani, M. A. Rivas, and T. Hastie, “A fast and flexible algorithm for solving the lasso in large-scale and ultrahigh-dimensional problems,” *BioRxiv*, p. 630079, 2019.
- [30] J. Fan and R. Li, “Variable selection via nonconcave penalized likelihood and its oracle properties,” *Journal of the American statistical Association*, vol. 96, no. 456, pp. 1348–1360, 2001.
- [31] J. Huang, H. Xie *et al.*, “Asymptotic oracle properties of scad-penalized least squares estimators,” in *Asymptotics: Particles, processes and inverse problems*. Institute of Mathematical Statistics, 2007, pp. 149–166.
- [32] L. Laporte, R. Flamary, S. Canu, S. Déjean, and J. Mothe, “Nonconvex regularizations for feature selection in ranking with sparse svm,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 25, no. 6, pp. 1118–1130, 2013.
- [33] P. Zhu, W. Zhu, W. Wang, W. Zuo, and Q. Hu, “Non-convex regularized self-representation for unsupervised feature selection,” *Image and Vision Computing*, vol. 60, pp. 22–29, 2017.
- [34] M. Yamada, W. Jitkrittum, L. Sigal, E. P. Xing, and M. Sugiyama, “High-dimensional feature selection by feature-wise kernelized lasso,” *Neural computation*, vol. 26, no. 1, pp. 185–207, 2014.

- [35] A. Mnih and D. J. Rezende, “Variational inference for monte carlo objectives,” *arXiv preprint arXiv:1602.06725*, 2016.
- [36] A. Miller, N. Foti, A. D’Amour, and R. P. Adams, “Reducing reparameterization gradient variance,” in *Advances in Neural Information Processing Systems*, 2017, pp. 3708–3718.
- [37] M. Figurnov, S. Mohamed, and A. Mnih, “Implicit reparameterization gradients,” in *Advances in Neural Information Processing Systems*, 2018, pp. 441–452.
- [38] T. M. Cover and J. A. Thomas, *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. New York, NY, USA: Wiley-Interscience, 2006.
- [39] R. Kohavi and G. H. John, “Wrappers for feature subset selection,” *Artificial intelligence*, vol. 97, no. 1-2, pp. 273–324, 1997.
- [40] G. Brown, A. Pocock, M.-J. Zhao, and M. Luján, “Conditional likelihood maximisation: a unifying framework for information theoretic feature selection,” *Journal of machine learning research*, vol. 13, no. Jan, pp. 27–66, 2012.
- [41] J. Chen, L. Song, M. J. Wainwright, and M. I. Jordan, “Learning to explain: An information-theoretic perspective on model interpretation,” *arXiv preprint arXiv:1802.07814*, 2018.
- [42] C.-H. Chang, L. Rampasek, and A. Goldenberg, “Dropout feature ranking for deep learning models,” 12 2017.
- [43] Y.-W. Chang and C.-J. Lin, “Feature ranking using linear svm,” in *Causation and Prediction Challenge*, 2008, pp. 53–64.
- [44] R. Rastogi and K. Shim, “Public: A decision tree classifier that integrates building and pruning,” *Data Mining and Knowledge Discovery*, vol. 4, no. 4, pp. 315–344, 2000.
- [45] R. Díaz-Uriarte and S. A. De Andres, “Gene selection and classification of microarray data using random forest,” *BMC bioinformatics*, vol. 7, no. 1, p. 3, 2006.
- [46] T. Chen and C. Guestrin, “Xgboost: A scalable tree boosting system,” in *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*. ACM, 2016, pp. 785–794.
- [47] I. Guyon, S. Gunn, A. Ben-Hur, and G. Dror, “Result analysis of the nips 2003 feature selection challenge,” in *Advances in neural information processing systems*, 2005, pp. 545–552.
- [48] J. H. Friedman, “Multivariate adaptive regression splines,” *The annals of statistics*, pp. 1–67, 1991.
- [49] C.-C. Chang and C.-J. Lin, “Libsvm: A library for support vector machines,” *ACM transactions on intelligent systems and technology (TIST)*, vol. 2, no. 3, p. 27, 2011.
- [50] C. Strobl, A.-L. Boulesteix, T. Kneib, T. Augustin, and A. Zeileis, “Conditional variable importance for random forests,” *BMC bioinformatics*, vol. 9, no. 1, p. 307, 2008.
- [51] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama, “Optuna: A next-generation hyperparameter optimization framework,” *KDD*, 2019.

- [52] M. Gregorová, J. Ramapuram, A. Kalousis, and S. Marchand-Maillet, “Large-scale nonlinear variable selection via kernel random features,” *arXiv preprint arXiv:1804.07169*, 2018.
- [53] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [54] D. D. Lewis, Y. Yang, T. G. Rose, and F. Li, “Rcv1: A new benchmark collection for text categorization research,” *Journal of machine learning research*, vol. 5, no. Apr, pp. 361–397, 2004.
- [55] G. X. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R. Wilson, S. B. Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu *et al.*, “Massively parallel digital transcriptional profiling of single cells,” *Nature communications*, vol. 8, p. 14049, 2017.
- [56] J. L. Katzman, U. Shaham, A. Cloninger, J. Bates, T. Jiang, and Y. Kluger, “Deepsurv: personalized treatment recommender system using a cox proportional hazards deep neural network,” *BMC Medical Research Methodology*, vol. 18, 2018.
- [57] C. Curtis, S. P. Shah, S.-F. Chin, G. Turashvili, O. M. Rueda, M. J. Dunning, D. Speed, A. G. Lynch, S. Samarajiwa, Y. Yuan *et al.*, “The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups,” *Nature*, vol. 486, no. 7403, p. 346, 2012.
- [58] H. Ishwaran, U. Kogalur, E. Blackstone, and M. Lauer, “Random survival forests,” *Annals of Applied Statistics*, vol. 2, no. 3, pp. 841–860, 9 2008.
- [59] M. J. Wainwright, “Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (lasso),” *IEEE Transactions on Information Theory*, vol. 55, no. 5, p. 2183–2202, May 2009. [Online]. Available: <http://dx.doi.org/10.1109/tit.2009.2016018>
- [60] R. J. Williams, “Simple statistical gradient-following algorithms for connectionist reinforcement learning,” *Machine Learning*, vol. 8, pp. 229–256, 1992.
- [61] H. Ishwaran and U. B. Kogalur, “Random survival forests for r,” 2007.
- [62] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” in *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, 2010, pp. 249–256.

Supplemental Material

S1 Proof of Proposition 1

We now provide a proof for proposition 1 showing the equivalence between the stochastic optimization (7) and the deterministic one (6). We start by considering \mathcal{S}' be a subset such that $\mathcal{S}^* \setminus \mathcal{S}' \neq \emptyset$. That is there exists some element in \mathcal{S}^* that is not in \mathcal{S}' . For any such set \mathcal{S}' we have that $I(\mathbf{X}_{\mathcal{S}'}; Y) < I(\mathbf{X}; Y)$. Indeed, if we let $i \in \mathcal{S}^* \cap \mathcal{S}'^c$ then we have

$$\begin{aligned} I(\mathbf{X}_{\mathcal{S}'}; Y) &\leq I(\mathbf{X}_{\setminus\{i\}}; Y) \\ &= I(\mathbf{X}; Y) - I(\mathbf{X}_i; Y | \mathbf{X}_{\setminus\{i\}}) \\ &< I(\mathbf{X}; Y), \end{aligned}$$

where the final inequality follows by Assumption 1. On the other hand, for any subset \mathcal{S}' such that $\mathcal{S}^* \subset \mathcal{S}'$, based on Assumption 2 we get that $I(\mathbf{X}_{\mathcal{S}'}; Y) = I(\mathbf{X}; Y)$. Now, when we consider the Bernoulli optimization problem we have

$$\max_{\boldsymbol{\pi}} I(\mathbf{X} \odot \tilde{\mathbf{S}}; Y) \quad \text{s.t.} \quad \sum_l \pi_l \leq k \text{ and } 0 \leq \pi_l \leq 1.$$

The mutual information can be expanded as

$$I(\mathbf{X} \odot \tilde{\mathbf{S}}; Y) = \sum_{\tilde{\mathbf{s}}} I(\mathbf{X} \odot \tilde{\mathbf{s}}; Y) p_{\boldsymbol{\pi}}(\tilde{\mathbf{S}} = \tilde{\mathbf{s}}),$$

where we have used the fact that $\tilde{\mathbf{S}}$ is independent of everything else. Our goal is to understand the form of the distributing of the random vector $\tilde{\mathbf{S}}$ that maximizes the objective subject to its constrains. Recall that in optimization (7) the coordinates of $\tilde{\mathbf{S}}$ are sampled at random, which is motivated by practical needs. This means that the distribution that is being optimized over $p_{\boldsymbol{\pi}}$ is a product distribution. Here, we will show that we can remove independence assumption. If we can show that the distribution found by solving the more general optimization problem is still a product distribution, then we obtain a solution to the original optimization (7).

Now, from above we know that the optimal value of the optimization is $I(\mathbf{X}_{\mathcal{S}'}; Y)$ for any set $\mathcal{S}^* \subset \mathcal{S}'$. Hence, any unconstrained distribution should place all of its mass on such subsets in order to maximize the mutual information. As a result $\sum_{l \in \mathcal{S}^*} p_{\boldsymbol{\pi}}(\tilde{\mathbf{S}}_l = 1) = k$. However, there is an optimization constraint that $\mathbb{E}[\sum_l \tilde{\mathbf{S}}_l] \leq k$. Therefore, $\mathbb{E}[\tilde{\mathbf{S}}_l] = 0$ for any $l \notin \mathcal{S}^*$. Hence, the optimal solution is to select the distribution so that all of the mass is placed on the subset \mathcal{S}^* and no mass elsewhere. As this is also a product distribution, this complete the proof of the claim.

S2 Bridging the Two Perspectives

To motivate the introduction of randomness into the risk, we have looked at the feature selection problem from a MI perspective. Based on the MI objective, we have observed that introducing randomness into the constrained maximization procedure does not change the objective (See Proposition 1 in Section 4) Here we provide a relation between the MI objective (6) and the empirical risk (1), which supports our proposed procedure.

We first note that the MI maximization over the set \mathcal{S} can be reformulated as the minimization of the conditional entropy $H(Y | \mathbf{X}_{\mathcal{S}})$ since $H(Y)$ does not depend on \mathcal{S} :

$$\max_{\mathcal{S}} I(\mathbf{X}_{\mathcal{S}}; Y) \iff \min_{\mathcal{S}} H(Y | \mathbf{X}_{\mathcal{S}}).$$

Recall that $\mathbf{X}_{\mathcal{S}} = \mathbf{X} \odot \tilde{\mathbf{S}}$. By Proposition 1, we rewrite the deterministic search over the set \mathcal{S} by a search over the Bernoulli parameters $\boldsymbol{\pi}$:

$$\begin{aligned} \min_{\boldsymbol{\pi}} H(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) &= \min_{\boldsymbol{\pi}} \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[-\log P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) \right] \\ &= \min_{\boldsymbol{\pi}, \boldsymbol{\theta}} \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[-\log P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) \right], \end{aligned}$$

where the expectation is over $\mathbf{X}, \mathbf{Y} \sim P_{\boldsymbol{\theta}^*}$, which is the true data distribution, and $\tilde{\mathbf{S}} \sim \text{Bern}(\tilde{\mathbf{S}}|\boldsymbol{\pi})$ and we put our model distribution as $P_{\boldsymbol{\theta}}$. Then we can rewrite the right hand side as:

$$\begin{aligned} \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \log P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) &= \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[\log P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) \frac{P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})}{P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})} \right] \\ &= \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[\log \frac{P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})}{P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})} \right] + \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[\log P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) \right]. \end{aligned}$$

Since $\text{KL}(P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})||P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}))$ is non-negative, $\mathbb{E}_{\tilde{\mathbf{S}}} \text{KL}(P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})||P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}))$ is also non-negative because it is a weighted sum of non-negative terms. Noting that

$$\mathbb{E}_{\tilde{\mathbf{S}}} \text{KL}(P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})||P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})) = \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[\log \frac{P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})}{P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})} \right]$$

we conclude that

$$\mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[-\log P_{\boldsymbol{\theta}^*}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) \right] \leq \mathbb{E}_{\mathbf{X}, \mathbf{Y}, \tilde{\mathbf{S}}} \left[-\log P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}}) \right].$$

If we consider the negative log likelihood of the target given the observations (i.e. $-\log P_{\boldsymbol{\theta}}(\mathbf{Y}|\mathbf{X} \odot \tilde{\mathbf{S}})$) as a loss function L (in Eq. (1)), then we see that the minimizing the risk approximately maximizes the MI objective in (6).

S3 Details of the Regularization Term

Here we provide a detailed description of our regularization term. For the vector of stochastic gates $\mathbf{z} \in \mathbb{R}^D$, the regularization term is expressed as follows:

$$\begin{aligned} \mathbb{E}_{\mathbf{Z}} \|\mathbf{Z}\|_0 &= \sum_{d=1}^D \mathbb{P}[z_d > 0] = \sum_{d=1}^D \mathbb{P}[\mu_d + \sigma \epsilon_d > 0] \\ &= \sum_{d=1}^D \{1 - \mathbb{P}[\mu_d + \sigma \epsilon_d \leq 0]\} \\ &= \sum_{d=1}^D \{1 - \Phi\left(\frac{-\mu_d}{\sigma}\right)\} \\ &= \sum_{d=1}^D \Phi\left(\frac{\mu_d}{\sigma}\right). \end{aligned}$$

Note that the derivative of the regularization term with respect to the distribution parameter μ_d is simply the Gaussian PDF:

$$\frac{\partial}{\partial \mu_d} \mathbb{E}_Z \|\mathbf{Z}\|_0 = \frac{\partial}{\partial \mu_d} \Phi\left(\frac{\mu_d}{\sigma}\right) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\mu_d^2}{2\sigma^2}}.$$

We now turn our attention towards providing a scheme for selecting σ . The effect of σ can be understood by looking at the value of $\frac{\partial}{\partial \mu_d} \mathbb{E}_Z \|\mathbf{Z}\|_0$. In the first training step, μ_d is 0.5. Therefore, at initial training phase, $\frac{\partial}{\partial \mu_d} \mathbb{E}_Z \|\mathbf{Z}\|_0$ is close to $\frac{1}{\sqrt{2\pi\sigma_d^2}} e^{-\frac{1}{8\sigma_d^2}}$. In order to remove irrelevant features, this term (multiplied by the regularization parameter λ) has to be greater than the derivative of the loss with respect to μ_d because otherwise μ_d is updated in the incorrect direction. To encourage such behavior, we set $\sigma = 0.5$, which is around the maximum of the gradient during the initial phase as shown in Fig. S1. Although the point that attains the maximum moves as μ changes, we empirically observe that setting $\sigma = 0.5$ performs well in our experiments when the regularization parameter λ is appropriately set. Note that in our implementation we divide the regularization term by the size of the features D . This rescaling normalizes the hyper-parameter search into a similar range.

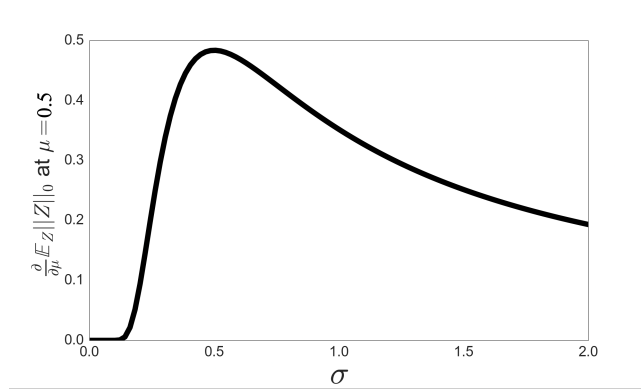


Figure S1: The value of $\frac{\partial}{\partial \mu} \mathbb{E}_Z \|\mathbf{Z}\|_0|_{\mu=0.5} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{8\sigma^2}}$ for $\sigma = [0.001, 2]$.

S4 Issues in Gradient Estimation of Discrete Random Variables

In Section 3.1, we have introduced Bernoulli random variables $\tilde{s}_d, d = 1, \dots, D$ with corresponding parameters π_d into the risk objective (4). Taking the expectation over the ℓ_0 norm of $\tilde{\mathbf{S}}$ boils down to the sum of the Bernoulli parameters π_d . However, a gradient-based optimization over the resulting objective suffers from high variance due to the discrete nature of $\tilde{\mathbf{S}}$. Here, we attempt to convey this problem by analyzing the risk term in the objective (4).

Using the Bernoulli parameterization the empirical risk $\hat{R}(\boldsymbol{\theta}, \boldsymbol{\pi})$ is expressed as

$$\sum_{\mathbf{z}:\{0,1\}^D} \left[\sum_{n=1}^N [L(f_{\boldsymbol{\theta}}(\mathbf{z} \odot \mathbf{x}_n), \mathbf{y}_n)] \prod_{d=1}^D \pi_d^{z_d} (1 - \pi_d)^{1-z_d} \right].$$

In practice, as the outer sum involves enumerating 2^D possibilities of the indicator variables, one can replace the outer sum with Monte Carlo samples from the product of Bernoulli distributions $B(\mathbf{z}|\boldsymbol{\pi})$. However, a Monte Carlo estimate of $\frac{\partial}{\partial \pi_d} \hat{R}(\boldsymbol{\theta}, \boldsymbol{\pi})$ suffers from high variance. To see this, consider the following exact gradient of the empirical risk with respect to π_d , which is

$$\sum_{\mathbf{z}:\{0,1\}^D, z_d=1} [L(\mathbf{z})p_{z_{i \neq d}}] - \sum_{\mathbf{z}:\{0,1\}^D, z_d=0} [L(\mathbf{z})p_{z_{i \neq d}}],$$

where $p(z_{i \neq d}) = \prod_{i \neq d}^D \pi_i^{z_i} (1 - \pi_i)^{1-z_i}$, by absorbing the model $f_\theta(\cdot)$ and the data into $L(\cdot)$. Due to the discrete nature of \mathbf{z} , we see that even the sign of the gradient estimate becomes inaccurate if we can only access a small number of Monte Carlo samples. While a score-function estimator such as REINFORCE [60] can be used, it is known that the reparametrization trick is more effective for variance reduction [26, 25, 27].

S5 Hard-Concrete Distribution (HC)

In the main text, we have compared the proposed embedded method using our STG distribution and an alternative based on the Hard-Concrete (HC). Here, we provide a full description for the HC distribution. The HC was introduced in [27] as a modification of Binary Concrete [26, 25], whose sampling procedure is as follows:

$$\begin{aligned} u &\sim U(0, 1), L = \log(U) - \log(1 - U), \\ s &= \frac{1}{1 + \exp\left(\frac{-(\log \alpha + L)}{\beta}\right)}, \\ \bar{s} &= s(\zeta - \tau) + \tau, \\ z &= \min(1, \max(0, \bar{s})), \end{aligned}$$

where (τ, ζ) is an interval, with $\tau < 0$ and $\zeta > 1$. This induces a new distribution, whose support is $[0, 1]$ instead of $(0, 1)$. With $0 < \beta < 1$, the probability density concentrates its mass near the end points, since values larger than $\frac{1-\tau}{\zeta-\tau}$ are rounded to one, whereas values smaller than $\frac{-\tau}{\zeta-\tau}$ are rounded to zero.

The CDF of s is

$$Q_s(s|\beta, \log \alpha) = \text{Sigmoid}((\log s - \log(1 - s))\beta - \log \alpha), \quad (10)$$

and so the CDF of \bar{s} is

$$Q_{\bar{s}}(\bar{s}|\phi) = \text{Sigmoid}\left(\left(\log\left(\frac{\bar{s} - \tau}{\zeta - \tau}\right) - \log\left(1 - \frac{\bar{s} - \tau}{\zeta - \tau}\right)\right)\beta - \log \alpha\right), \quad (11)$$

where $\phi = (\beta, \log \alpha, \zeta, \tau)$. Using this distribution to model the gates, the probability of a gate z being active is $1 - Q_{\bar{s}}(0|\phi)$ and can be written as

$$1 - Q_{\bar{s}}(0|\phi) = \text{Sigmoid}\left(\log \alpha - \beta \log \frac{-\tau}{\zeta}\right). \quad (12)$$

The hyperparameters β, ξ and τ are set as in [27] and fixed throughout training, while α is learned to determine whether the gate is active or not. Note that L is distributed as a Logistic distribution, which causes high variance in \bar{s} due to the heavy-tailness of the distribution. By replacing the logistic distribution with the Gaussian distribution, we reduce variance in gradient estimates, leading to stability of feature selection (See Section 8).

S6 Additional Experiments

S6.1 Two Moons classification with nuisance features

In this experiment, we construct a dataset based on "two moons" shape classes, concatenated with noisy features. The first two coordinates x_1, x_2 are generated by adding a Gaussian noise with zero mean and the variance of $\sigma_r^2 = 0.1$ onto two nested half circles, as presented in Fig. 2(a). Nuisance features $x_i, i = 3, \dots, D$, are drawn from a Gaussian distribution with zero mean and variance of $\sigma_n^2 = 1$. We reserve the 70% as a test set, and use 10% of the remaining training set as a validation set. We follow the same hyperparameter tuning procedure as in the XOR experiment. The classification accuracy is in Fig. 2(b). Based on the classification accuracies, it is evident that for a small number of nuisance dimensions all methods correctly identify the most relevant features. The proposed method (STG) and Random Forest (RF) are the only methods that achieve near perfect classification accuracy for a wide range of nuisance dimensions. The other NN based method (DFS) seem to converge to sub-optimal solutions. We note that all the methods achieve the median rank 1.5, which is the optimal median rank in this example.

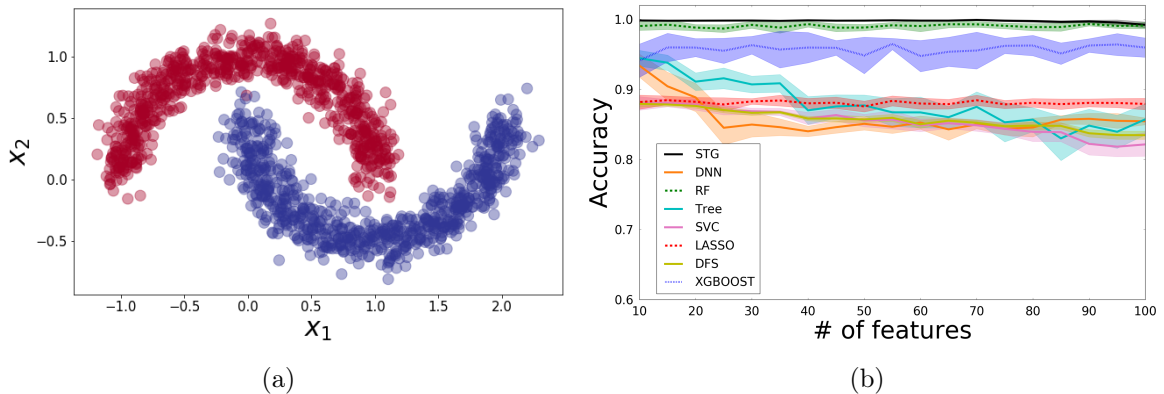


Figure S2: (a) Realizations from the "Two moons" shaped binary classification class. x_1 and x_2 are the relevant features, $x_i, i = 3, \dots, D$ are noisy features drawn from a Gaussian with zero mean and variance of 1. (b) Classification accuracy (mean and standard deviation based on 20 runs) vs. the number of irrelevant noisy dimension.

S6.2 Convergence Comparison to Hard-Concrete distribution

In this section, we show additional experiments to compare STG with HC in terms of convergence speed.

The main difference between our proposed distribution (STG) and the Hard-Concrete [27] distribution is that the latter is based on the logistic distribution, which has a heavier tail than the Gaussian distribution we have employed. As shown in Fig S3, the heavy-tailness results in instability during training. Furthermore, the STG converges much faster and more reliably than the feature selection method using the HC distribution on a the two-moons, XOR and MADELON datasets (see Subsection 6.1 and S6.1).

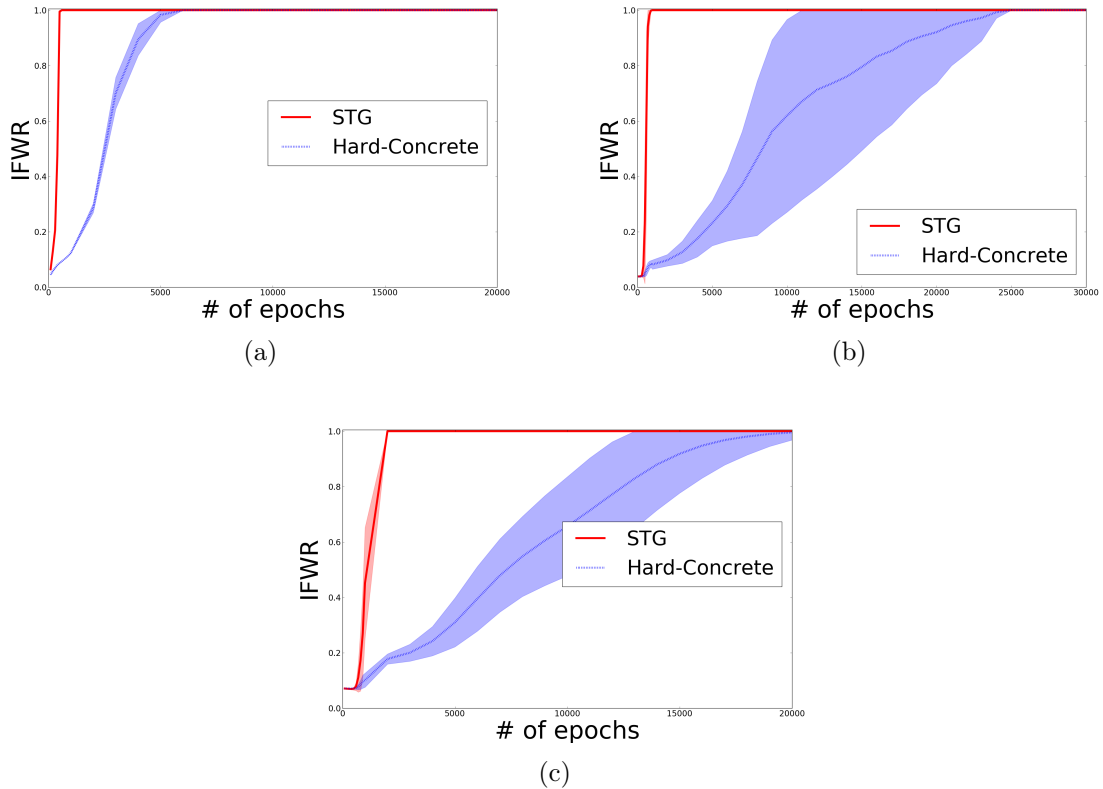


Figure S3: Comparison between STG and Hard-Concrete on the Two-Moon dataset (a), XOR (b) and MADELON (c). The shaded area represents the standard deviation, calculated by running the same experiment 10 times with different random initializations.

S7 Reproducibility

Here we provide a full description of the procedures and datasets we have used in the experimental parts of the paper.

S7.1 Datasets

S7.2 Reuters Corpus Volume I

The Reuters Corpus Volume I (RCV1) consists of 800,000 news stories manually labeled by 103 categories. This is a multilabel regime, where each story is assigned to multiple categories. Here we focus on a binary subset of this corpus, with 23,203 stories, where we use 10%, 8.5%, and 81.5% for the train, validation and test set, respectively. The total number of feature is 47,236. This example demonstrates that our method is also effective in an extremely high dimensional regime of nonlinear function estimation.

S7.3 Purified populations of peripheral blood monocytes (PBMCs)

Single-cell RNA sequencing (scRNA-seq) is a novel technology that measures gene expression levels of hundreds of thousands of individual cells. [55], have subjected more than 90,000

purified populations of peripheral blood monocytes (PBMCs) to scRNA-seq analysis. Here we focus on classifying two subpopulations of T-cells, namely the Naive and regulatory T-cells.

The data consists of $D = 11,382$ genes and $N = 20742$ cells, of which we only use 10% of the data for training. We apply the proposed method for different values of λ and report the number of selected features and classification accuracy on the test set. In this example, the STG identifies 20 – 30 genes which are sufficient for the classification task.

S7.4 Details of Cox Proportional Hazard Model

Survival times are assumed to follow a distribution, which is characterized by the survival function $S(t) = P(T > t)$. A hazard function, which measures the instantaneous rate of death, is defined by $h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t < T \leq t + \Delta t | T > t)}{\Delta t} = \frac{p(t)}{S(t)}$. We can relate the two functions in the following way: $S(t) = e^{-\int_0^t h(t) dt}$.

Proportional hazard models assume a multiplicative effect of the covariates x on the hazard function such that $h(t|x) = h_0(t)e^{\boldsymbol{\theta}^T \mathbf{x}}$, where $h_0(t)$ is a baseline hazard function, which is often the exponential or Weibull distribution, and $\boldsymbol{\theta}$ is the parameter of interests.

One of the difficulties in estimating $\boldsymbol{\theta}$ in survival analysis is that a large portion of the available data is censored. However, in order to obtain estimates, Cox observed that it is sufficient to maximize the partial-likelihood, which is defined as follows:

$$L(\boldsymbol{\theta}) = \prod_{T_i \text{ uncensored}} \frac{e^{\boldsymbol{\theta}^T \mathbf{x}_i}}{\sum_{T_j \geq T_i} e^{\boldsymbol{\theta}^T \mathbf{x}_j}}.$$

In [56], the authors propose DeepSurv, which uses a deep neural network model to replace the linear relations between the covariate \mathbf{x} and $\boldsymbol{\theta}$, demonstrating improvements of survival time prediction over existing models such as CPH and the random survival forest [61], [58].

The Molecular Taxonomy of Breast Cancer International Consortium (METABRIC) dataset consists of gene expression data and clinical features for 1,980 patients, and 57.72% have an observed death due to breast cancer with a median survival time of 116 months.

The METABRIC dataset involves 24,368 features (genes). Most genes are irrelevant for outcome prediction. To demonstrate the advantage of STG in the context of survival analysis, we selected 16 well-known genes relevant for survival (out of the 24,368 genes) that correspond to the Oncotype DX test, a gene panel used for treatment decision making. We also include five additional clinical features (hormone treatment indicator, radiotherapy indicator, chemotherapy indicator, ER-positive indicator, and age at diagnosis). We then added 200 additional irrelevant gene variables that we selected randomly from the remaining list of genes.

After we omit the null values, the number of samples is 1969. We reserve the 20% for test, and use the 20% of the remaining training set as validation (that is, train: 1260, valid: 315, test: 394 samples).

Experimental Detail For DeepSurv, we manually select the architecture using the validation set so that we obtain a similar performance reported in [56]. The learning rate decay is set to 1. The learning rate and the regularization parameter are optimized via Optuna [51] using the validation set, where the search range is set $[1e - 3, 1]$ for the learning rate and $[1e - 3, 1]$ for λ . The hyperparameters used in the experiment are the following: architecture : [60, 20, 3], activation: Selu (as suggested by [56]), learning rate : 0.152, learning rate decay: 1.0, σ : 0.5, λ : 0.023, training epoch: 2000. Note that to see the effect of feature selection, we used the hyperparameters optimized for DeepSurv to test our method (Cox-STG).

S7.5 Implementation details

Datasets are first split into train, validation and test. Validation is always 10% of the train, while the exact ratios between train and test is detailed for each experiment separately (see Table 1). All the neural network weights are initialized by drawing from $\mathcal{N}(0, 0.1)$ and bias terms are set to 0 (following Xavier initialization [62]). We use SGD for all the experiments, except for the Cox model where we use Adam. All the experiments are conducted using Intel(R) Xeon(R) CPU E5-2620 v3 @2.4Ghz x2 (12 cores total). We set the number of Monte Carlo samples $K = 1$, which worked well in our experiments. Hyper-parameters for all method are tuned using Optuna [51]. Optuna is a hyper-parameter optimization software that supports pruning and parallel computing across GPUs. We run $n = 1000$ trails with parameters search ranges as described in Table S1.

Table S1: List of the search range for the hyperparameters used in our experiments

Param	Search range
# dense layers	[1,5]
# hidden units	[10, 500]
activation	[Tanh, Relu, Sigmoid]
LR	[1e-4, 1e-1]
n-epoch (DFS, SG-NN)	[50, 20000]
α (SG-NN)	[1e-3, 10]
λ (SG-NN)	[1e-7, 10]
λ (STG, DFS)	[1e-3, 10]
λ (LASSO)	[1e-5, 1]
n-est (RF, XGB, ERT)	[5,100]
n-boost-round (XGB)	[1,100]
Thresh (RF, XGB, ERT)	[1e-7,1]
max-depth (XGB)	[1,5]
c (SVC)	[1e-7, 1]

For linear regression experiment, we use 0.1 as a learning rate. For the XOR problem, the exact architectures used for the NN based methods are: (STG/HC): [476, 490, 14] with Tanh, (DFS): [100, 10] with Tanh, (SG-NN): [100, 10, 5] with Tanh. For the two moons we use (STG): [490, 406, 18] with Tanh, (DFS): [158, 27, 224] with Tanh, (SG-NN): [88, 28, 27] with Tanh. For the XOR problem, we attempted to use Optuna to optimize parameters of DFS and SG-L1-NN, but we ended up using the architecture suggested by the authors [22, 23] as they outperform the values suggested by Optuna. The number of epochs used for the XOR problem is $20k$, $14k$, 800 for STG/HC, DFS and SG-L1-NN respectively. Regularization parameters are 0.17 , $3.3e - 5$ and $3e - 5$ respectively. The number of epochs used for the two-moons problem is $20K$, 1570 , 708 for STG/HC, DFS and SG-NN respectively. Regularization parameters are 0.48 , $9e - 3$ and $1e - 3$ respectively. We note that the regularization parameters and learning procedure is different in nature, as we use an ℓ_0 type penalty. For the PBMC experiment, the architecture use us [27, 10, 383] with Tanh activations, a learning rate was 0.0036, batch size is 1000 and the number of epochs 8000. The hyperparameter λ varies in the range [0.001, 0.11] to achieve different levels of sparsity. For MADELON, GISETTE, ISOLET, we use the architecture optimized for the binary XOR classification. The number of epochs used are $20K$, $20K$, $4k$, the learning rate are 0.06, 0.2, 0.1, batch size are 200, 1000, 40. The parameters used are the

following: (SE1) architecture [600, 200, 100, 50] with ReLu activations, number of epochs is 5, $\lambda : 5$, learning rate 0.0001 (SE2) architecture [600, 300, 150, 60, 20] with ReLu activations, number of epochs is 2000, $\lambda : 5$, learning rate is 0.001 (SE3) architecture [600, 300, 150, 60, 20] with ReLu activations, num epochs : 1000, $\lambda : 1$, learning rate 0.005. (RCP) architecture [1000, 300, 150, 60, 20] with ReLu activation, num-epochs: 2000, $\lambda : 5.0$, learning rate 0.001. (REL) architecture [26, 91, 63] with ReLu activation, number of epochs: 1600, $\lambda : 0.031$, learning rate 0.007. (RAI) architecture [10, 177] with ReLu activation, number of epochs: 1800, $\lambda : 0.019$, learning rate is 0.002. Architectures for SE1-SE3 and RCP were tuned manually. The ratio of train/test/valid split is 1:1:1 for synthetic regression data. For the real regression data (RCP and REL), the train size is 6000, the test and valid size is 1000 samples. For RAI the train size is 5000, the test and valid size is 1000 samples. For the Friedman data we use Tanh with an architecture of [500, 200, 100, 20] (HC/STG/DFS/SG-NN) with Tanh activations, a learning rate of 0.2 and a batch size of 200. For BASEHOCK (STG/HC) the architecture used is [500, 105, 25] with Tanh activations, a batch size of 50, learning rate of 0.5 and 2000 epochs. For RELATHE (STG/HC) the architecture used is [500, 200, 10] with Tanh activations, a batch size of 40, learning rate of 0.008 and $10k$ epochs. For COIL20 (STG/HC) the architecture used is [32, 438, 158, 20] with Tanh activations, a batch size of 150, learning rate of 0.12 and $10k$ epochs. For PCMAC (STG/HC) the architecture used is [109, 25, 455] with Tanh activations, a batch size of 450, learning rate of 0.5 and 6000 epochs. For RCV1 (STG/HC) the architecture used is [500, 100] with Tanh activations, a batch size of 500, learning rate of 0.5 and 650 epochs. For MNIST (STG/HC) the architecture used is [300, 100] with Relu activations (as in [53]), a batch size of 200, learning rate of 0.1 and 250 epochs.

In any experiment where we count the number of active features, we evaluate the set of indices such that: $\{d : \min(1, \max(0, \mu_d + 0.5)) > 0\}$ after training. In order to define the IFWR, for STG, the d^{th} feature weight is set to $\max(0, \min(1, \mu_d + 0.5))$. For other neural net based methods, it is given by $\sum_j W_{dj}$, where W is the weight matrix of the first layer. For other methods we just used the feature relevance returned by the trained model. Finally, the LASSO’s IFWR in the XOR experiment was omitted from the manuscript as it suffered from high variance.

Regarding the comparison performed in the two-moons and XOR problem, we believe that adding IFWR along with classification accuracy versus number of feature selected provides a complementary perspective in demonstrating the efficacy of feature selection techniques. We emphasize that our goal is not to just rank features but select features by assigning the weight of 0 to irrelevant features while simultaneously obtaining good predictive accuracy.