
Subject Section

Block HSIC Lasso: model-free biomarker detection for ultra-high dimensional data

Héctor Climente-González^{1,2,3,6}, Chloé-Agathe Azencott^{3,1,2}, Samuel Kaski⁴ and Makoto Yamada^{5,6,7*}

¹Institut Curie, PSL Research University, F-75005 Paris, France

²INSERM, U900, F-75005 Paris, France

³MINES ParisTech, PSL Research University, CBIO-Centre for Computational Biology, F-75006 Paris, France

⁴Aalto University, Helsinki, Finland

⁵Kyoto University, Kyoto, 606-8501, Japan

⁶RIKEN AIP, Chuo-ku, Tokyo 103-0027, Japan.

⁷JST PRESTO, Saitama, 332-0012, Japan.

*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Abstract

Motivation: Finding nonlinear relationships between biomolecules and a biological outcome is computationally expensive and statistically challenging. Existing methods have crucial drawbacks, among others lack of parsimony, non-convexity, and computational overhead. Here we present the block HSIC Lasso, a nonlinear feature selector that does not present the previous drawbacks.

Results: We compare the block HSIC Lasso to other state-of-the-art feature selection techniques in synthetic data and real data, including experiments over three common types of genomic data: gene-expression microarrays, single-cell RNA-seq, and GWAS. In all the cases, we observe that features selected by block HSIC Lasso retain more information about the underlying biology than features of other techniques. As a proof of concept, we applied the block HSIC Lasso to a single-cell RNA-seq experiment on mouse hippocampus. We discovered that many genes linked in the past to brain development and function are involved in the biological differences between the types of neurons.

Availability: Block HSIC Lasso is implemented in the Python 2/3 package `pyHSICLasso`, available in [Github](https://github.com/riken-aip/pyHSICLasso) (<https://github.com/riken-aip/pyHSICLasso>) and [PyPi](https://pypi.org/project/pyHSICLasso) (<https://pypi.org/project/pyHSICLasso>).

Contact: myamada@i.kyoto-u.ac.jp

Supplementary information: Supplementary data are available at *Bioinformatics* online.

1 Introduction

Biomarker discovery, the goal of many bioinformatics experiments, aims at identifying a few key biomolecules that explain most of an observed phenotype. Without a strong prior hypothesis, these

molecular markers have to be identified from data generated by high-throughput technologies. Unfortunately, finding relevant molecules is a combinatorial problem: for d features, 2^d binary choices must be considered. As the number of features vastly outnumbers the number of samples, biomarker discovery is a high-dimensional problem. The statistical challenges posed by such high-dimensional spaces have been thoroughly reviewed elsewhere (Clarke *et al.*, 2008; Johnstone and

Titterton, 2009). In general, due to the *curse of dimensionality*, fitting models in many dimensions and on a small number of samples is extremely hard. Moreover, since biology is complex, a simple statistical model such as linear regression might not be able to find important biomarkers. Even when they are found, often they are hard to reproduce, suggesting overfitting. Exploring the solution space and finding true biomarkers is not only statistically challenging, but also computationally expensive.

In machine learning terms, biomarker discovery can be formulated as a problem of feature selection: identifying the best subset of features to separate between categories, or to predict a continuous response. In the past decades, many feature selection algorithms that deal with high-dimensional datasets have been proposed. Due to the difficulties posed by high-dimensionality, linear methods tend to be the feature selector of choice in bioinformatics. A widely used linear feature selector is the Least Absolute Shrinkage and Selection Operator, or Lasso (Tibshirani, 1996), a sparse learning algorithm. Lasso fits a linear model between the input features and phenotype by minimizing the sum of the least square loss and an ℓ_1 penalty term. The balance between the least square loss and the penalty ensures that the model explains the linear combination of features, while keeping the number of features in the model small. However, in many biology instances do not behave linearly. In such cases there is no guarantee that Lasso can capture those nonlinear relationships or an appropriate effect size to represent them.

In the past decade, several nonlinear feature selection algorithms for high-dimensional datasets have been proposed. Among them, the widely used sparse additive model, or SpAM (Ravikumar *et al.*, 2009), models the outcome as a sparse linear combination of nonlinear functions based on kernels, where the input of each function is a feature. However, since SpAM assumes an additive model over the selected features, it cannot select important features if the phenotype cannot be represented by the additive functions of input features (for example, if there exist a multiplicative relationship between features (Yamada *et al.*, 2014)).

Another family of nonlinear feature selectors are association-based: they compute the association scores between each input feature and the outcome and rank features according to the association scores. Since these approaches do not assume any model about the output, they can detect important features as long as there exist some association. When using a nonlinear association measure, such as the mutual information (Cover and Thomas, 2006) or the Hilbert-Schmidt Independence Criterion (HSIC) (Gretton *et al.*, 2005), they select the features with the strongest dependence with the phenotype. However, association-based methods do not account for the redundancy between the features, which is frequent in biological data sets, since they do not take into account relationships between features. Hence, many redundant features are typically selected, hindering interpretability. This is important in applications like drug target discovery, where only a small number of targets can be validated, and it is crucial to discriminate the most important target out of many other top-ranked targets.

To deal with the problem of redundant features, Peng *et al.* (2005) proposed the minimum redundancy maximum relevance (mRMR) algorithm. The goal of mRMR is to select a set of non-redundant features that have high association to the phenotype, while penalizing the selection of mutually dependent features by a mutual information term. Ding and Peng (2005) used mRMR to extract biomarkers from microarray data, finding that the selected genes captured better the variability in the phenotypes than other feature selectors. However, mRMR has three main drawbacks: the optimization problem is discrete,

it must be solved by a greedy approach, and the mutual information estimation is hard (Walters-Williams and Li, 2009). To be precise, it is unknown whether the objective function of mRMR has good theoretical properties such as submodularity (Fujishige, 2005), which would guarantee that the optimisation approach actually finds the mathematically optimal solution.

Recently, Yamada *et al.* (2014) proposed a kernel-based minimum redundancy maximum relevance algorithm called HSIC Lasso. Instead of mutual information, HSIC Lasso employs the Hilbert-Schmidt Independence Criterion (HSIC) (Gretton *et al.*, 2005) to measure dependency between variables. In addition, it uses an ℓ_1 penalty term to select a small number of features. This results in a convex optimization problem, for which one can therefore find a globally optimal solution. In practice, HSIC Lasso has been found to outperform mRMR in several experimental settings (Yamada *et al.*, 2014). Moreover, thanks to the sparsity assumption, it scales well in memory with the number of features. However, HSIC Lasso is memory-intensive: its memory complexity is $O(dn^2)$, where d is the number of features and n is the number of samples. Hence, HSIC Lasso cannot be applied to datasets with thousands of samples, nowadays widespread in biology. A MapReduce version of HSIC Lasso has been proposed to solve this drawback (Yamada *et al.*, 2018). MapReduce HSIC Lasso is able to select features in ultra-high dimensional settings (10^6 features, 10^4 samples) in a matter of hours (Yamada *et al.*, 2018). However, it requires a large number of computing nodes, inaccessible to common laboratories. And, since it relies on the Nyström approximation (Schölkopf and Smola, 2002) of Gram matrices, the final optimization problem is no longer convex, and hence finding a globally optimal solution cannot be easily guaranteed.

In this paper, we propose the block HSIC Lasso: a simple yet effective nonlinear feature selection algorithm based on HSIC Lasso. The key idea is to estimate the HSIC terms with the recently proposed block HSIC estimator (Zhang *et al.*, 2018). Thanks to the block HSIC estimator, the memory complexity goes from $O(dn^2)$ down to $O(dnB)$, where $B \ll n$ is the block size. Moreover, as opposed to MapReduce HSIC Lasso, the optimization problem of the block HSIC Lasso remains convex. Through its application to synthetic data and biological datasets, we show that block HSIC Lasso can be applied to a variety of biomarker discovery settings and compares favorably with the vanilla HSIC Lasso algorithm and other feature selection approaches, linear and nonlinear, as it selects features more informative of the biological outcome. Both the block HSIC Lasso and HSIC Lasso are available in the Python 2/3 package pyHSICLasso (<https://pypi.org/project/pyHSICLasso>).

2 Materials and methods

2.1 Problem formulation

Assume a data set with n samples described by d features, each corresponding to a biomolecule (for example, the expression of one transcript, or the number of major alleles observed at a given SNP), and a label, continuous or binary, describing the outcome of interest (for example, the abundance of a target protein, or disease status). We denote the i -th sample by $\mathbf{x}_i = [x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(d)}]^\top \in \mathbb{R}^d$, where $^\top$ denotes transpose; and its label by $y_i \in \mathcal{Y}$, where $\mathcal{Y} = \{0, 1\}$ for a binary outcome, corresponding to a classification problem, and $\mathcal{Y} = \mathbb{R}$ for a continuous outcome, corresponding to a regression problem. In

addition, we denote by $\mathbf{f}_k = [x_1^{(k)}, x_2^{(k)}, \dots, x_n^{(k)}]^\top \in \mathbb{R}^n$ the k -th feature in the data.

The goal of supervised feature selection is to find m features ($m \ll d$) that are the most relevant for predicting the output y for a sample \mathbf{x} .

2.2 HSIC Lasso

Measuring the dependence between two random variables X and Y can be achieved by the Hilbert-Schmidt Independence Criterion, or HSIC (Gretton *et al.*, 2005):

$$\begin{aligned} \text{HSIC}(X, Y) &= \mathbb{E}_{x, x', y, y'} [K(x, x')L(y, y')] \\ &+ \mathbb{E}_{x, x'} [K(x, x')] \mathbb{E}_{y, y'} [L(y, y')] \\ &- 2\mathbb{E}_{x, y} [\mathbb{E}_{x'} [K(x, x')] \mathbb{E}_{y'} [L(y, y')]], \end{aligned} \quad (1)$$

where $K(x, x')$ and $L(y, y')$ are positive definite kernels, and $\mathbb{E}_{x, x', y, y'}$ denotes the expectation over independent pairs (x, y) and (x', y') drawn from $p(x, y)$. HSIC can be used as an independence measure if we use the characteristic kernels (Fukumizu *et al.*, 2004; Sriperumbudur *et al.*, 2011). That is, X and Y are independent if and only if $\text{HSIC}(X, Y) = 0$ and take positive value otherwise. Thus, we can select important features by ranking HSIC scores $\{\text{HSIC}(X_m, Y)\}_{m=1}^d$ in descending order, where X_m is the random variable of the m -th feature (Song *et al.*, 2012).

In practice, for a given Gram matrix $\mathbf{K}_k \in \mathbb{R}^{n \times n}$, computed from k -th feature, and a given output Gram matrix $\mathbf{L} \in \mathbb{R}^{n \times n}$, HSIC is computed as

$$\text{HSIC}_v(\mathbf{f}_k, \mathbf{y}) = \text{tr}(\overline{\mathbf{K}}_k \overline{\mathbf{L}}), \quad (2)$$

where for any Gram matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$, $\overline{\mathbf{K}}$ is defined as

$$\overline{\mathbf{K}} = \frac{\mathbf{H}\mathbf{K}\mathbf{H}}{\|\mathbf{H}\mathbf{K}\mathbf{H}\|_F}$$

with $\mathbf{H} \in \mathbb{R}^{n \times n}$ a centering matrix defined by $H_{ij} = \delta_{ij} - \frac{1}{n}$. Here δ_{ij} is equal to 1 if $i = j$ and 0 otherwise, and tr denotes the trace. $\text{HSIC}_v(\mathbf{f}_k, \mathbf{y})$ is equal to zero when \mathbf{f}_k and \mathbf{y} are independent, and has a non-negative value (the larger the more dependence between the inputs) when they are not. Note that we employ the normalized variant of the original empirical HSIC.

In this paper, we use the following kernels:

- The RBF Gaussian kernel for continuous variables (or outcomes):

$$K : x_i^{(k)}, x_j^{(k)} \mapsto \exp\left(-\frac{\|x_i^{(k)} - x_j^{(k)}\|_2^2}{2\sigma^2}\right),$$

where $\sigma^2 > 0$ is the bandwidth of the kernel;

- The normalized Dirac kernel for categorical variables (or outcomes):

$$L : y_i, y_j \mapsto \begin{cases} \frac{1}{n_c} & \text{if } y_i = y_j = c \\ 0 & \text{otherwise,} \end{cases}$$

where n_c is the number of samples in class c .

The goal of HSIC Lasso (Yamada *et al.*, 2014) is to select independent features that have a high dependence with the outcome.

For that purpose, it introduces a vector $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_d]^\top$ of feature weights and solves the following optimization problem:

$$\begin{aligned} \max_{\boldsymbol{\alpha} \geq 0} \sum_{k=1}^d \alpha_k \text{HSIC}_v(\mathbf{f}_k, \mathbf{y}) - \frac{1}{2} \sum_{k, k'=1}^d \alpha_k \alpha_{k'} \text{HSIC}_v(\mathbf{f}_k, \mathbf{f}_{k'}) \\ - \lambda \|\boldsymbol{\alpha}\|_1. \end{aligned} \quad (3)$$

The selected features are those that have a non-zero coefficient α_k . Here $\lambda > 0$ is a regularization parameter that controls the sparsity of the solution: the larger λ , the fewer features have a non-zero coefficient.

The HSIC Lasso optimization problem can be rewritten as

$$\min_{\boldsymbol{\alpha} \geq 0} \|\text{vec}(\overline{\mathbf{L}}) - [\text{vec}(\overline{\mathbf{K}}_1), \dots, \text{vec}(\overline{\mathbf{K}}_d)]\boldsymbol{\alpha}\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1,$$

where $\text{vec} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n^2}$, $\mathbf{K} \mapsto [K_{11}, \dots, K_{1n}, K_{21}, \dots, K_{nn}]$ is the vectorization operator. Using this formulation, we can solve the problem using an off-the-shelf non-negative Lasso solver.

HSIC Lasso performs well in particular for high-dimensional data. However, it requires a large memory space ($O(dn^2)$), since it stores d Gram matrices. To handle this issue, two approximation methods have been proposed. The first approach uses a memory lookup to dramatically reduce the memory space (Yamada *et al.*, 2014). However, since this method needs to perform a large number of memory lookups, it is computationally expensive. Another approach (Yamada *et al.*, 2018) is to rewrite the problem using the Nyström approximation (Schölkopf and Smola, 2002) and solve the problem using a cluster. However using the Nyström approximation makes the problem non-convex.

2.3 Block HSIC Lasso

In this paper, we propose an alternative HSIC Lasso method for large-scale problems, the *block HSIC Lasso*, which is convex and can be efficiently solved on a reasonably-sized server.

Block HSIC Lasso employs the block HSIC estimator (Zhang *et al.*, 2018) instead of the V-statistics estimator of Equation (2). More specifically, to compute the block HSIC, we first partition the training dataset into $\frac{n}{B}$ partitions $\{(x_i^\ell, y_i^\ell)\}_{i=1}^B\}_{\ell=1}^{n/B}$, where B is the number of samples in each block. Note that the block size B is set to a relatively small number such as 10 or 20 ($B \ll n$). Then, the block HSIC estimator can be written as

$$\text{HSIC}_b(\mathbf{f}_k, \mathbf{y}) = \frac{B}{n} \sum_{\ell=1}^{n/B} \text{HSIC}_v(\mathbf{f}_k^{(\ell)}, \mathbf{y}^{(\ell)}),$$

where $\mathbf{f}_k^{(\ell)} \in \mathbb{R}^B$ represents the k -th feature vector of the ℓ -th partition. Note that the computation of $\text{HSIC}_v(\mathbf{f}_k^{(\ell)}, \mathbf{y}^{(\ell)})$ requires $O(B^2)$ memory. That is, the required memory for the block HSIC is $O(nB)$, where $nB \ll n^2$.

If we denote by $\overline{\mathbf{K}}_k^{(\ell)} \in \mathbb{R}^{B \times B}$ the restriction of $\overline{\mathbf{K}}_k$ to the ℓ -th partition, and by $\overline{\mathbf{L}}^{(\ell)} \in \mathbb{R}^{B \times B}$ the restriction of $\overline{\mathbf{L}}$ to the ℓ -th partition, then

$$\text{HSIC}_v(\mathbf{f}_k^{(\ell)}, \mathbf{y}^{(\ell)}) = \text{tr}(\overline{\mathbf{K}}_k^{(\ell)} \overline{\mathbf{L}}^{(\ell)}) = \text{vec}(\overline{\mathbf{K}}_k^{(\ell)})^\top \text{vec}(\overline{\mathbf{L}}^{(\ell)}).$$

The block HSIC Lasso is obtained by replacing the HSIC estimator HSIC_v with the block HSIC estimator HSIC_b in Equation (3):

$$\max_{\alpha \geq 0} \sum_{k=1}^d \alpha_k \text{HSIC}_b(\mathbf{f}_k, \mathbf{y}) - \frac{1}{2} \sum_{k,k'=1}^d \alpha_k \alpha_{k'} \text{HSIC}_b(\mathbf{f}_k, \mathbf{f}_{k'}) - \lambda \|\alpha\|_1. \quad (4)$$

Using the vectorization operator, the block estimator can be written as

$$\text{HSIC}_b(\mathbf{f}_k, \mathbf{f}_{k'}) = \mathbf{u}_k^\top \mathbf{u}_{k'}, \quad \text{HSIC}_b(\mathbf{f}_k, \mathbf{y}) = \mathbf{u}_k^\top \mathbf{v},$$

where

$$\mathbf{u}_k = \sqrt{\frac{B}{n}} \left[\text{vec} \left(\overline{\mathbf{K}}_k^{(1)} \right)^\top, \dots, \text{vec} \left(\overline{\mathbf{K}}_k^{(n/B)} \right)^\top \right]^\top \in \mathbb{R}^{nB},$$

$$\mathbf{v} = \sqrt{\frac{B}{n}} \left[\text{vec} \left(\overline{\mathbf{L}}^{(1)} \right)^\top, \dots, \text{vec} \left(\overline{\mathbf{L}}^{(n/B)} \right)^\top \right]^\top \in \mathbb{R}^{nB}.$$

Hence, the block HSIC Lasso can also be written as

$$\min_{\alpha \geq 0} \|\mathbf{v} - \mathbf{U}^\top \alpha\|_2^2 + \lambda \|\alpha\|_1,$$

where $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_d] \in \mathbb{R}^{nB \times d}$.

Since the objective function of the block HSIC Lasso is convex, we can obtain a globally optimal solution. And, as with HSIC Lasso, we can solve the block HSIC Lasso using an off-the-shelf Lasso solver. Here, we use the non-negative least angle regression-LASSO, or LARS-LASSO (Efron *et al.*, 2004), to solve the problem in a greedy manner.

The required memory space for block HSIC Lasso is $O(dnB)$, which compares favorably to vanilla HSIC Lasso's $O(dn^2)$; as the block size $B \ll n$, the memory space is dramatically reduced. However, the computational cost of the proposed method is still large when both d and n are large. Thus, we implemented the proposed algorithm using multiprocessing by parallelizing the computation of $\overline{\mathbf{K}}_k^{(\ell)}$. Thanks to the combination of the block HSIC Lasso and the multiprocessing implementation we can efficiently find solutions on large datasets with a reasonably-sized server.

2.4 Improving selection stability using bagging

Since we need to compute the block HSIC of the paired data $\{(\mathbf{x}_i^\ell, \mathbf{y}_i^\ell)\}_{i=1}^B\}_{\ell=1}^{n/B}$ with a fixed partition, the performance can be highly affected by the partition. Thus, we propose to use a bagging version of the block HSIC estimator. Given M random permutations of the n samples, we define *bagging block HSIC* as

$$\text{HSIC}_{bb}(\mathbf{f}_k, \mathbf{y}) = \frac{1}{M} \sum_{m=1}^M \frac{B}{n} \sum_{\ell=1}^{n/B} \text{HSIC}_v(\mathbf{f}_k^{(\ell,m)}, \mathbf{y}^{(\ell,m)}) = \overline{\mathbf{u}}_k^\top \overline{\mathbf{v}},$$

where $\mathbf{f}_k^{(\ell,m)}$ is the k -th feature vector restricted to the ℓ -th block as defined by the m -th permutation,

$$\overline{\mathbf{u}}_k = \sqrt{\frac{1}{M}} \left[\mathbf{u}_k^{(1)\top}, \dots, \mathbf{u}_k^{(M)\top} \right]^\top \in \mathbb{R}^{nBM},$$

$$\overline{\mathbf{v}} = \sqrt{\frac{1}{M}} \left[\mathbf{v}^{(1)\top}, \dots, \mathbf{v}^{(M)\top} \right]^\top \in \mathbb{R}^{nBM},$$

and $\mathbf{u}_k^{(m)} \in \mathbb{R}^{nB}$ and $\mathbf{v}_k^{(m)} \in \mathbb{R}^{nB}$ are the vectors of the m -th block HSIC Lasso, respectively.

Hence, the bagging block HSIC Lasso can be written as

$$\min_{\alpha \geq 0} \|\overline{\mathbf{v}} - \overline{\mathbf{U}}^\top \alpha\|_2^2 + \lambda \|\alpha\|_1,$$

where $\overline{\mathbf{U}} = [\overline{\mathbf{u}}_1, \dots, \overline{\mathbf{u}}_d] \in \mathbb{R}^{nBM \times d}$.

Note that the memory space $O(dnBM)$ required for $B = 60$ and $M = 1$ is equivalent to $B = 30$ and $M = 2$. It is not clear which choice of parameters provides better feature selection accuracy.

3 Experimental setup

3.1 Feature selection methods

HSIC Lasso and block HSIC Lasso: We used HSIC Lasso and block HSIC Lasso implemented in Python 2/3 package *pyHSICLasso*. In block HSIC Lasso, M was set to 3 in all experimental settings; the block size B was set on an experiment-dependent fashion. In all the experiments, when we wanted to select k features, HSIC Lasso versions were required to first retrieve 50 features, and then the top k features were selected as the solution.

mRMR: Minimum Redundancy Maximum Relevance (mRMR) selects features that have high relevance with respect to the outcome and are non-redundant (Peng *et al.*, 2005). To that end, it uses mutual information between different variables and between the outcome and the variables. The mRMR score of a set of features \mathbf{V} is defined as

$$\text{mRMR}(\mathbf{V}) = \frac{1}{m} \sum_{k=1}^d \overline{\text{MI}}(\mathbf{f}_k, \mathbf{y}) - \frac{1}{m^2} \sum_{k,k'=1}^d \overline{\text{MI}}(\mathbf{f}_k, \mathbf{f}_{k'})$$

where $\overline{\text{MI}}(\mathbf{f}_k, \mathbf{y})$ is an empirical estimate of mutual information (Peng *et al.*, 2005). As the problem of finding the set of features \mathbf{V} is nonconvex, mRMR implementations rely on a greedy search.

We used a C++ implementation of mRMR (Peng, 2005). The maximum number of samples and the maximum number of features were set to the actual number of samples and features in the data. In regression problems, discretization was set to binarization.

LARS: Least angle regression (LARS) is a forward stagewise feature selector (Efron *et al.*, 2004). It is an efficient way of solving the same problem as Lasso. We used the SPAMS implementation of LARS (Mairal *et al.*, 2010), with the default parameters.

3.2 Evaluation of the selected features

Selection accuracy on simulated data: We simulated high-dimensional data where a few variables were truly related to the outcome. We used these datasets to find out the ability of the tested algorithms to find the true causal variables, instead of other, likely spuriously correlated. To that end, we requested each algorithm to retrieve the known number of causal features. Then, we studied how many of them were actually causal.

Classification with a random forest: To evaluate the information retained in the features selected by a method, we trained a random forest classifier using only those features to recover the original categories. We used the random forest because it is able to handle nonlinearity. We selected the features on the whole training set. We estimated the best parameters by cross-validation on the training set: the number of trees (200, 500), the maximum depth of the trees (4, 6, 8), the number

Table 1. Summary of benchmark datasets.

Type	Dataset	Features (d)	Samples (n)	Classes
Image	ARI10P	2,400	130	10
	PIE10P	2,400	210	10
	PIX10P	10,000	100	10
	ORL10P	10,000	100	10
Microarray	CLL-SUB-111	11,340	111	3
	GLIOMA	4,434	50	4
	SMK-CAN-187	19,993	187	2
	TOX-171	5,748	171	4
scRNA-seq	Haber <i>et al.</i> (2017)	15,972	7,216	19
	Habib <i>et al.</i> (2016)	25,393	13,302	8
	Villani <i>et al.</i> (2017)	23,395	1,140	10
GWA data	RA vs. controls	362,577	3,479	2
	T1D vs. controls	393,676	3,443	2
	T2D vs. controls	393,722	3,479	2

of features to consider (\sqrt{d} , $\log_2 d$), and the criterion to measure the quality of the chosen features (Gini impurity, information gain). Then, we trained a model with those parameters on the whole training set and made predictions on a separate testing set to find out the accuracy.

3.3 Datasets

We evaluated the performance of the different algorithms on synthetic data and four types of real-world high dimensional datasets (Table 1). All real-world datasets are on classification problems here, though HSIC Lasso can handle regression problems (continuous-valued outcomes) as well, as we show on synthetic data.

Synthetic data: We simulated random matrices of features $X \sim \mathcal{N}(0, 1)$. A number of variables were selected as related to the phenotype, and functions that are nonlinear in the data range were selected (cosine, sine and square) and combined additively to create the outcome vector y .

Images: Facial recognition is a classical supervised nonlinear feature selection problem. We used four face image datasets from Arizona State University feature selection repository (Li *et al.*, 2016): pixraw10P, warpAR10P, orlraws10P, and warpPIE10P.

Gene expression microarrays: We looked for biomarkers in four gene expression microarray datasets from Arizona State University feature selection repository (Li *et al.*, 2016). The phenotypes were subtypes of B-cell chronic lymphocytic leukemia (CLL-SUB-111), hepatocyte phenotypes under different diets (TOX-171), glioma (GLIOMA) and smoking-driven carcinogenesis (SMK-CAN-187).

Single-cell RNA-seq: Single-cell RNA-seq (scRNA-seq) measures gene expression at cell resolution, allowing to characterize the diversity in a tissue. We looked for biomarkers on the three most popular datasets in Broad Institute’s Single Cell Portal, related to mouse small intestinal epithelium (Haber *et al.*, 2017), mouse hippocampus (Habib *et al.*, 2016), and human blood cells (Villani *et al.*, 2017). Gene expression was imputed with MAGIC (van Dijk *et al.*, 2018).

GWA datasets: We studied the WTCCC1 datasets (Burton *et al.*, 2007) for rheumatoid arthritis (RA), type 1 diabetes (T1D) and type 2 diabetes (T2D) (2,000 samples each), using the 1958BC cohort as control (1,504 samples). Affymetrix 500K was used for genotyping. We removed the samples and the SNPs that did not pass WTCCC’s quality controls, as well as SNPs in sex chromosomes and those that

were not genotyped in both cases and controls. Missing genotypes were imputed with CHIAMO. When they could not be imputed, we replaced the missing value by the major allele in homozygosis.

Preprocessing: Images, microarrays and scRNA-seq data were normalized feature-wise by subtracting the mean and dividing by the standard deviation. GWAS data did not undergo any normalization.

3.4 Computational resources

We ran the experiments on synthetic data, images, microarrays and scRNA-seq on CentOS 7 machines with Intel® Xeon® 2.6GHz and 50 GB RAM memory. For the GWA datasets experiments, we used the CentOS 7 server with 96 core Intel® Xeon® 2.2GHz and 1TB RAM memory.

3.5 Software availability

Block HSIC Lasso was implemented in the Python 2/3 package *pyHSICLasso*. The source code is available at Github (<https://github.com/riken-aip/pyHSICLasso>), and the package can be installed from PyPi (<https://pypi.org/project/pyHSICLasso>). All the analyses and the scripts needed to reproduce them are available in Github (<https://github.com/hclimente/nori>).

4 Results

4.1 Block HSIC Lasso performance is comparable to state of the art

At first, we worked on synthetic, nonlinear data (section 3.2). We generated synthetic data with combinations of the following experimental parameters: $n = \{100, 1000, 10000\}$ samples; $d = \{100, 2500, 5000, 10000\}$ features; and 5, 10 and 20 causal features i.e. features truly related to the outcome. We evaluated the performance of different feature selectors at retrieving the causal features. The algorithms compared were the block HSIC Lasso, two other nonlinear methods (vanilla HSIC Lasso (Yamada *et al.*, 2014) and mRMR (Peng *et al.*, 2005)) and a linear method (LARS (Efron *et al.*, 2004)).

Each of the methods was required to retrieve the number of true causal features. In Figure 1 we show the proportion of the causal features retrieved by each method. The different versions of HSIC Lasso outperform the other approaches in virtually all settings. Block HSIC Lasso with decreasing block sizes results in worse performances. As expected, vanilla HSIC Lasso outperforms the block versions in accuracy, but increases memory use. Crucially, block HSIC Lasso on a larger number of samples performs better than vanilla HSIC Lasso on fewer samples. Hence, when the number of samples is in the thousands, it is better to apply block HSIC Lasso on the whole dataset, than to apply vanilla HSIC Lasso on a subsample.

We wanted to test these conclusions using a nonlinear, real-world dataset. We selected four image-based face recognition tasks (section 3.3). In this case, we selected different numbers of features (10, 20, 30, 40 and 50). Then, we trained random forest classifiers on these subsets of the features, and compared the accuracy of the different classifiers on a test set (Figure S1). Block HSIC Lasso displayed a performance comparable to vanilla HSIC Lasso, and comparable or superior to the other methods. This is remarkable, since it shows that, in many practical

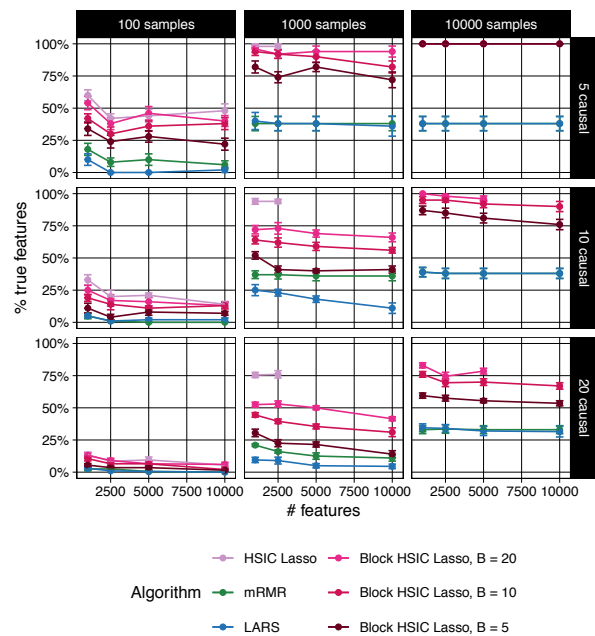


Fig. 1. Percentage of true causal features extracted by different feature selectors. Each data point represents the mean over 10 replicates, and the error bars represent the standard error of the mean. Lines are discontinued when the algorithm required more memory than the provided (50 GB). Note that in some conditions mRMR's line cannot be seen due to the overlap with LARS.

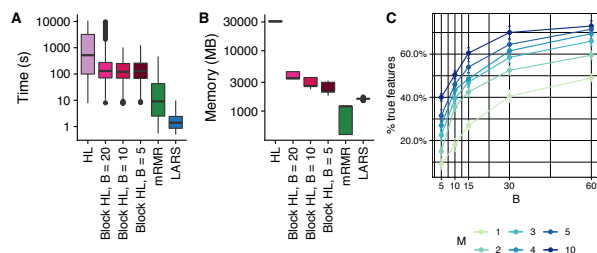


Fig. 2. Computational resources used by the different methods. A. Time elapsed in a multiprocess setting. B. Memory usage in a single-core setting. C. Number of correct features retrieved on synthetic data ($n = 1000$, $d = 2500$, 20 causal features) by block HSIC Lasso at different block sizes B and number of permutations M .

cases, block HSIC Lasso does not need more samples to achieve vanilla HSIC Lasso performance.

4.2 Block HSIC Lasso is computationally efficient

In our experiments on synthetic data, vanilla HSIC Lasso runs into memory issues already with 1,000 samples (Figure 1). This experiment shows how block HSIC Lasso keeps the good properties of HSIC Lasso, while extending it to more experimental settings. Block HSIC Lasso with $B = 20$ reaches the memory limit only at 10,000 samples, which is already sufficient for most common bioinformatics applications. If larger datasets need to be handled, that can be done by using smaller block sizes or a larger computer cluster.

We next quantified the computational efficiency improvement the block HSIC estimator brings. We compared the runtime and the peak

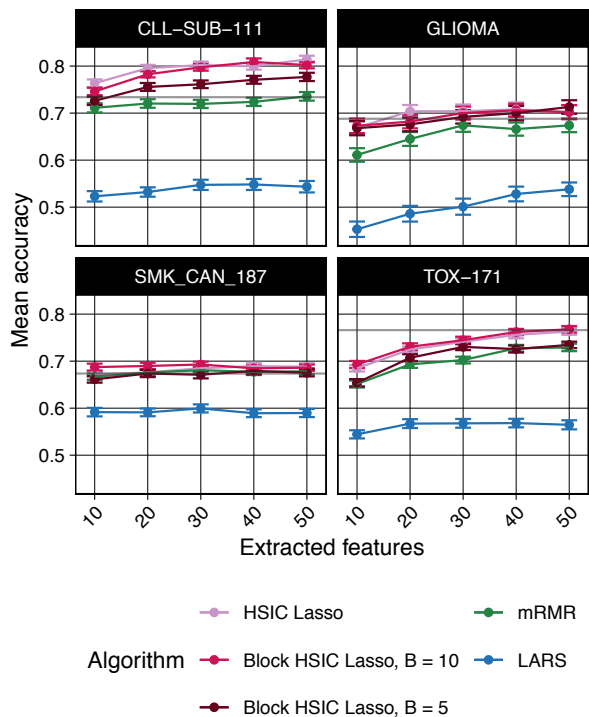


Fig. 3. Random forest classification accuracy of microarray gene expression samples after feature extraction by the different methods. The gray line represents the mean accuracy of 10 classifiers trained on all the dataset.

memory usage in the highest-dimensional setting where all methods could run ($n = 1000$, $d = 2500$, 20 causal features) (Figure 2). We observe how, as expected, block HSIC Lasso requires an order of magnitude less memory than vanilla HSIC Lasso. Block versions also run notoriously faster, thanks to the lower number of operations and the parallelization. mRMR is ten times faster than block HSIC Lasso, at the expense of a clearly lower accuracy. However, a fraction of this gap is likely due to mRMR having been implemented in C++, while HSIC Lasso is written in Python. In this regard, there is potential for other faster implementations of (block) HSIC Lasso.

4.3 HSIC estimator improves with more permutations

We were interested in the trade-off between the block size and the number of permutations, which affect both the computation time and accuracy of the result. We tested the performance of block HSIC Lasso with $B = \{5, 10, 15, 30, 60\}$ and $M = \{1, 2, 3, 5\}$ in datasets of $n = 1000$, $d = 2500$ and 20 causal features. As expected, causal feature recovery increases with M and B (Figure 2C), as the HSIC estimator approaches its true value.

The memory usage $O(dnBM)$ of several of the conditions was the same e.g. $B = 10$, $M = 3$ and $B = 30$, $M = 1$. Such conditions are indistinct from the points of view of both accuracy, and memory requirements. In practice, we found no major differences in runtime between different combinations of B and M . Hence, a reasonable strategy is to fix B to a given size, and tune the M to the available memory/desired amount of information.

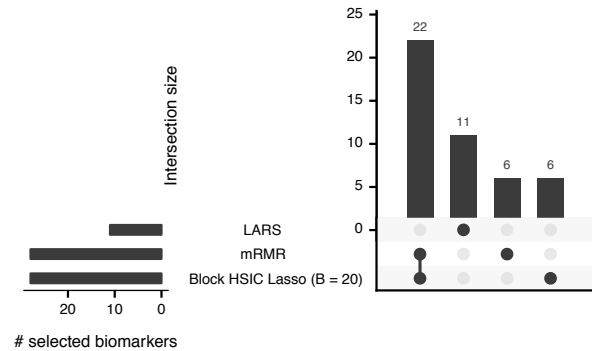


Fig. 4. Overlap between the pooled SNPs selected by three feature selectors (block HSIC Lasso, mRMR and LARS) on three GWAS datasets. We asked each method to select 10 features in each dataset. The horizontal histogram represents the total number of SNPs selected by each algorithm. The vertical histogram represents the size of the overlap between SNP sets (22 SNPs are selected by both block HSIC Lasso and mRMR; 11 were exclusively selected by LARS and so on).

4.4 Block HSIC Lasso finds relevant features efficiently

We tested the dimensionality reduction potential of different feature selectors. We selected a variable number of features from different multi-class biological datasets, then used a random forest classifier to retrieve the original classes (section 3.2). The underlying assumption is that only selected features which are biologically relevant will be useful to classify unseen data. To that end, we evaluated the classification ability of the biomarkers selected in four gene expression microarrays (Figure 3) and three scRNA-seq experiments (Figure S2). Unsurprisingly, we observe that nonlinear feature selectors perform notably better than linear selectors. Of the nonlinear methods, in virtually all cases block HSIC Lasso showed similar or superior performance to mRMR. Interestingly, as little as 20 selected genes retain enough information to achieve a plateau accuracy in most experiments.

Surveying $10^5 - 10^6$ SNPs in $10^3 - 10^4$ patients, genome-wide association (GWA) datasets are among the most high-dimensional in biology, an unbalance which worsens the statistical and computational challenges. We performed the same evaluation on three WTCCC1 phenotypes (section 3.3). As a baseline, we also computed the accuracy of a classifier trained on all the SNPs (Table S2). We observe that a feature selection prior step is not always favourable: LARS worsens the classification accuracy by 10%. On top of that, LARS could not select any SNP in 3 out of the 15 experimental settings. On the other hand, nonlinear methods improve the classification accuracy by 15-20%, with mRMR and block HSIC Lasso achieving similar accuracies. In fact, mRMR and block HSIC Lasso selected the same 22 out of 30 SNPs when we selected 10 SNPs in each the three datasets with each method (Figure 4).

4.5 Block HSIC Lasso is robust to ill-conditioned problems

Single-cell RNA-seq datasets differ from microarray datasets in two ways. First, the number of features is larger, equaling the number of genes in the annotation ($>20,000$). Second, the expression matrices are very sparse, due to biological variability (genes actually not expressed in a particular cell) and dropouts (genes whose expression levels have not been measured, usually because they are low, i.e. technical zeroes).

In summary, the problem is severely ill-conditioned, and the feature selectors need to deal with this issue. We observed that block HSIC Lasso runs reliably when faced with variations in the data, even on ill-conditioned problems like scRNA-seq. In the different scRNA-seq datasets, LARS was unable to select the requested number of biomarkers in any of the cases, returning always a lower number (Figure 5). mRMR did in all cases. However, the implementation of mRMR that we used crashed while selecting features on the full Villani *et al.* (2017) dataset.

4.6 Block HSIC Lasso for biomarker discovery

4.6.1 New biomarkers in mouse hippocampus scRNA-seq

To study the potential of block HSIC lasso for biomarker discovery in scRNA-seq data, we focused on the mouse hippocampus dataset from Habib *et al.* (2016), as a list of 1669 known biomarkers for the different cell types is also provided by the authors. We requested block HSIC Lasso, mRMR and LARS to select the best 20 genes for classification of 8 cell types (Table 1). The cell types were four different hippocampal anatomical subregions (DG, CA1, CA2 and CA3), glial cells, ependymal, GABAergic and unidentified cells.

The overlap between the genes selected by different algorithms was empty. We compared the selected genes to the known biomarkers. Out of the 20 genes selected by mRMR, 14 are known biomarkers, a number that goes down to 0 in the case of block HSIC Lasso (Figure 5A). Hence, these 20 genes, which are sufficient for accurately separating the cell types, are potential novel biomarkers. However, we have no reason to believe that HSIC Lasso generally has a higher tendency to return novel biomarkers than other approaches; we merely emphasize that it suggests alternative, statistically plausible biological hypotheses that can be worth investigating.

We therefore evaluated whether the novel biomarkers found by block HSIC Lasso participate in biological functions known to be different between the cell classes. To obtain the biological processes responsible for the differences between classes, we mapped the known biomarkers to GO Biological process categories using the GO2MSIG database (Powell, 2014). Then we repeated the process using the genes selected by the different feature selectors, and compared the overlap between them. The overlap between the different techniques increases when we consider the biological process instead of specific genes (Figure 5B). Specifically, one biological process term that is shared between mRMR and block HSIC Lasso, “Adult behaviour” (associated to *Sez6* and *Klhl1*, respectively), is clearly related to hippocampus function. This reinforces the notion that the selected biomarkers are relevant for the studied phenotypes.

Then we focused on potential biomarkers and biologically interesting molecules among those genes selected by block HSIC Lasso. As it is designed specifically to select non-redundant features, often-used GO enrichment analyses are not meaningful: we expect biomarkers belonging to the same GO annotation to be correlated, and HSIC lasso should not accumulate them. Among the top 5 genes we find 2 genes mapped to a biological processes known to be involved: the aforementioned *Klhl1* and *Pou3f1* (related to Schwann cell development). *Klhl1* is gene expressed in 7 of the studied cell types and which has been related to neuron development in the past (He *et al.*, 2006). *Pou3f1* is a transcription factor which in the past has been linked to myelination, and neurological damage in its absence (Jaegle *et al.*, 1996). The only gene among the top 5 that was expressed exclusively in one of the clusters is the micro RNA *Mir670*, expressed exclusively

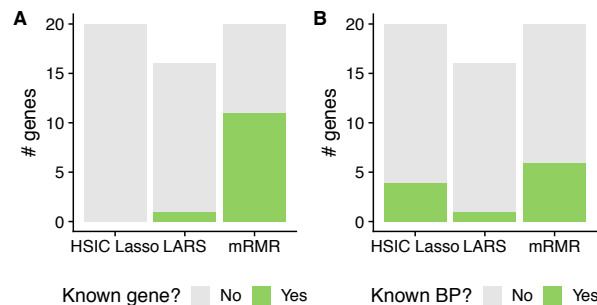


Fig. 5. Genes selected by different feature selection algorithms on single-cell RNA-seq data from mouse hippocampus (Habib et al., 2016). A. Count of genes that are (not) known biomarkers. B. Count of genes that share (not) a GO Biological process with a known biomarker.

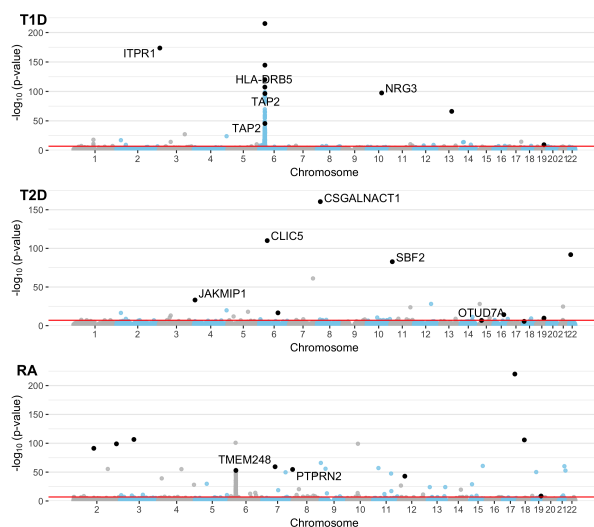


Fig. 6. Manhattan plot of the GWA datasets using p-values from the genotypic test. A constant of 10^{-220} was added to all p-values to allow plotting p-values of 0. SNPs in black are the SNPs selected by block HSIC Lasso ($B = 20$), 10 per phenotype. When SNPs are located within the boundaries of a gene (± 50 kb), the gene name is indicated. The red line represents the Bonferroni threshold with $\alpha = 0.05$.

in CA1. According to miRDB (Wong and Wang, 2015), *Mir670* top predicted target of its 3' arm is *Pcnt*, which is involved in neocortex development.

4.6.2 GWAS without assumptions on genetic architecture

We applied block HSIC Lasso ($B = 20$) to all the three GWA datasets. Typically the GWAS practitioner assumes a genetic model before the biomarker discovery. Two common, well-known models are the additive model – the minor allele in homozygosity has twice the effect as the minor allele in heterozygosity – and the dominant model – any number of copies of the minor allele have a phenotypic outcome. Using nonlinear models such as block HSIC Lasso to explore the relationship between SNPs and outcome is attractive since no assumptions are needed on the genetic architecture of the trait. On top of that, by penalizing the selection of redundant features, block HSIC Lasso avoids selecting multiple SNPs in high linkage disequilibrium.

In our experiments, we selected 10 SNPs with block HSIC Lasso for each of the three phenotypes. These are the SNPs that best balance high relatedness to the phenotype and not giving redundant information, be it through linkage disequilibrium or through an underlying shared biological mechanism. We compared these SNPs to those selected by the univariate statistical tests implemented in PLINK 1.9 (Chang et al., 2015). Some of them explicitly account for nonlinearity by considering dominant and recessive models of inheritance. We found that the genotypic test is the statistical procedure that displays the largest overlap with block HSIC Lasso selected biomarkers (28 out of 30). It consists of a 2df χ^2 test on a 3-by-2 contingency table representing the genotypes and the outcome. The genotypic test also displays the largest overlap with the other methods (Figure S3). Hence, in this case relaxing the assumptions on the genetic model of inheritance does not bring advantages over making stronger assumptions on the data, but on the other hand, roughly the same results can be achieved without the assumptions.

We compared the genome-wide genotypic p-values to the SNPs selected by block HSIC Lasso (Figure 6). Block HSIC Lasso selects SNPs among those with the most extreme p-values. However, not being constrained by a conservative p-value threshold, block HSIC Lasso selects two SNPs in type 2 diabetes with low, albeit non-Bonferroni significant, p-values when they improve classification accuracy (rs597414 , $p = 1.781 \times 10^{-7}$; rs543759 , $p = 2.904 \times 10^{-6}$). Moreover, the selected SNPs are scattered all across the genome, displaying the lack of redundancy between them. This strategy gives a more representative set of SNPs than other approaches common in bioinformatics, like selecting the smallest 10 p-values.

5 Discussion

In this work, we presented block HSIC Lasso, a nonlinear feature selector. Block HSIC Lasso retains the properties of HSIC Lasso while extending its applicability to larger datasets. Among the attractive properties of block HSIC Lasso we find, first, its ability to handle both linear and nonlinear relationships between the variables and the outcome. Second, block HSIC Lasso has a convex formulation, ensuring that a global solution exists, and that it is accessible. Third, the HSIC score can be accurately estimated, as opposed to other measures of nonlinearity like mutual information. Fourth, block HSIC Lasso's memory consumption scales linearly with respect to both the number of features and the number of samples. Lastly, block HSIC Lasso can be easily adapted to different problems via different kernel functions that better capture similarities in new datasets. Due to all these properties, we show how block HSIC Lasso outperforms all other algorithms in the tested conditions.

Block HSIC Lasso can be applied to different kinds of datasets. As other nonlinear methods, block HSIC Lasso is particularly useful when we do not want to make assumptions about how the causal variables and the outcome are related. Owing to the advantages mentioned above, HSIC Lasso and block HSIC Lasso tend to outperform other state-of-the-art approaches when comparing the classification accuracy.

Regarding its potential in bioinformatics, we applied block HSIC Lasso to images, microarrays, single-cell RNA-seq and GWAS. The two latter involve thousands of samples and variables, making it unfeasible to find a server with enough memory for vanilla HSIC Lasso. The selected biomarkers are biologically plausible, agree with the outcome of other

methods, and provide a good classification accuracy when used to train a classifier. Such a ranking is useful, for instance, when selecting SNPs or genes to assay in *in vitro* experiments.

Block HSIC Lasso's main drawback is the memory complexity, markedly lower than in vanilla HSIC Lasso but still $O(dnB)$. Memory issues might appear in low-memory servers in cases with a large number of samples n , of features d , or both. However, through our work on GWA datasets, the largest type of dataset in bioinformatics, we show that working on these datasets is feasible. Another drawback, which block HSIC Lasso shares with the other nonlinear methods, is their black box nature. Block HSIC Lasso looks for biomarkers which, after an unknown, nonlinear transformation, would allow a linear separation between the samples. Unfortunately, we cannot access this transformed space and explore it, which makes the results hard to interpret.

Funding

Computational resources and support were provided by RIKEN AIP. HC is funded by the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No 666003. SK is supported by the Academy of Finland (grants 292334, 319264). MY is supported by the JST PRESTO program JPMJPR165A and partly supported by MEXT KAKENHI 16H06299 and the RIKEN engineering network funding.

References

- Burton, P. R. *et al.* (2007). Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*, **447**(7145), 661–678.
- Chang, C. C. *et al.* (2015). Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience*, **4**, 7.
- Clarke, R. *et al.* (2008). The properties of high-dimensional data spaces: implications for exploring gene and protein expression data. *Nature Reviews Cancer*, **8**(1), 37–49.
- Cover, T. M. and Thomas, J. A. (2006). *Elements of Information Theory*. John Wiley & Sons, Inc., Hoboken, NJ, USA, 2nd edition.
- Ding, C. and Peng, H. (2005). Minimum redundancy feature selection from microarray gene expression data. *Journal of Bioinformatics and Computational Biology*, **03**(02), 185–205.
- Efron, B. *et al.* (2004). Least angle regression. *The Annals of statistics*, **32**(2), 407–499.
- Fujishige, S. (2005). *Submodular functions and optimization*, volume 58. Elsevier.
- Fukumizu, K. *et al.* (2004). Dimensionality reduction for supervised learning with reproducing kernel hilbert spaces. *Journal of Machine Learning Research*, **5**(Jan), 73–99.
- Gretton, A. *et al.* (2005). Measuring statistical dependence with Hilbert-Schmidt norms. In *ALT*.
- Haber, A. L. *et al.* (2017). A single-cell survey of the small intestinal epithelium. *Nature*, **551**(7680), 333–339.
- Habib, N. *et al.* (2016). Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn neurons. *Science (New York, N.Y.)*, **353**(6302), 925–8.
- He, Y. *et al.* (2006). Targeted deletion of a single Sca8 ataxia locus allele in mice causes abnormal gait, progressive loss of motor coordination, and Purkinje cell dendritic deficits. *The Journal of neuroscience : the official journal of the Society for Neuroscience*, **26**(39), 9975–82.
- Jaegle, M. *et al.* (1996). The POU factor Oct-6 and Schwann cell differentiation. *Science (New York, N.Y.)*, **273**(5274), 507–10.
- Johnstone, I. M. and Titterton, D. M. (2009). Statistical challenges of high-dimensional data. *Philosophical transactions. Series A, Mathematical, physical, and engineering sciences*, **367**(1906), 4237–53.
- Li, J. *et al.* (2016). Feature selection: A data perspective. *arXiv:1601.07996*.
- Mairal, J. *et al.* (2010). Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, **11**, 19–60.
- Peng, H. (2005). mrmr. <http://home.penglab.com/proj/mRMR/>.
- Peng, H. *et al.* (2005). Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **27**, 1226–1237.
- Powell, J. A. C. (2014). GO2MSIG, an automated GO based multi-species gene set generator for gene set enrichment analysis. *BMC bioinformatics*, **15**, 146.
- Ravikumar, P. *et al.* (2009). Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **71**(5), 1009–1030.
- Schölkopf, B. and Smola, A. J. (2002). *Learning with Kernels*. MIT Press, Cambridge, MA.
- Song, L. *et al.* (2012). Feature selection via dependence maximization. *JMLR*, **13**, 1393–1434.
- Sriperumbudur, B. K. *et al.* (2011). Universality, characteristic kernels and RKHS embedding of measures. *Journal of Machine Learning Research*, **12**(Jul), 2389–2410.
- Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, **58**(1), 267–288.
- van Dijk, D. *et al.* (2018). Recovering Gene Interactions from Single-Cell Data Using Data Diffusion.
- Villani, A.-C. *et al.* (2017). Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors. *Science (New York, N.Y.)*, **356**(6335), 925–8.
- Walters-Williams, J. and Li, Y. (2009). Estimation of mutual information: A survey. In P. Wen, Y. Li, L. Polkowski, Y. Yao, S. Tsumoto, and G. Wang, editors, *Rough Sets and Knowledge Technology*, pages 389–396, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Wong, N. and Wang, X. (2015). miRDB: an online resource for microRNA target prediction and functional annotations. *Nucleic acids research*, **43**(Database issue), D146–52.
- Yamada, M. *et al.* (2014). High-dimensional feature selection by feature-wise kernelized lasso. *Neural computation*, **26**(1), 185–207.
- Yamada, M. *et al.* (2018). Ultra high-dimensional nonlinear feature selection for big biological data. *IEEE Transactions on Knowledge and Data Engineering*, **30**(7), 1352–1365.
- Zhang, Q. *et al.* (2018). Large-scale kernel methods for independence testing. *Statistics and Computing*, **28**(1), 113–130.