
Composite Feature Selection Using Deep Ensembles

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In many real world problems, features do not act alone but in combination with
2 each other. For example, in genomics, diseases might not be caused by any single
3 mutation but require the presence of multiple mutations. Prior work on feature
4 selection either seeks to identify individual features or can only determine relevant
5 groups from a predefined set. We investigate the problem of discovering groups
6 of predictive features without predefined grouping. To do so, we define predictive
7 groups in terms of linear and non-linear interactions between features. We introduce
8 a novel deep learning architecture that uses an ensemble of feature selection models
9 to find predictive groups, without requiring candidate groups to be provided. The
10 selected groups are sparse and exhibit minimum overlap. Furthermore, we propose
11 a new metric to measure similarity between discovered groups and the ground truth.
12 We test our model on multiple synthetic tasks, semi-synthetic chemistry datasets
13 and image datasets to demonstrate its utility.

14 1 Introduction

15 Feature selection is a key problem permeating statistics, machine learning and broader science.
16 Typically in high-dimensional datasets, the majority of features will not be responsible for the target
17 response and thus an important goal is to identify which variables are truly predictive. For example,
18 in healthcare there may be many features (such as age, sex, medical history, etc.) that could be
19 considered, while only a small subset might in fact be relevant for predicting the likelihood of
20 developing a specific disease. By eliminating irrelevant variables, feature selection algorithms can be
21 used to drive discovery, improve model generalisation/robustness, and improve interpretability [16].

22 However, features often do not act alone but instead in *combination*. In genetics, for instance, it has
23 been noted that understanding the origins of many diseases may require methods able to identify more
24 complex genetic models than single variants [57]. While feature selection might be able to identify a
25 set of features associated with a particular response, the underlying structure of how features interact
26 is not captured. Further, the resulting predictive models can be complex, hard to interpret, and not
27 amenable to the generation of hypotheses that can be experimentally tested [41]. This limits the
28 impact such models can have in furthering scientific understanding across many domains where
29 variables are known to interact, such as genetics [57, 61, 54], medicine [79, 11], and economics [7].

30 Group feature selection is a generalisation of standard feature selection, where instead of selecting
31 individual features, groups of features are either entirely chosen or entirely excluded. A primary
32 application of group feature selection is when features are jointly measured, for example by different
33 instruments. In such scenarios, groups are readily defined as features measured by the same instrument.
34 A natural question is which instruments give the most meaningful measurements. Group feature
35 selection has also been applied in situations where there is extensive domain knowledge regarding the

36 group structure [64] or where groups are defined by the correlation structure between features (e.g.
37 neighbouring pixels in images are highly correlated). The pervasive issue with current group feature
38 selection methods is that a predetermined grouping *must* be provided, and the groups are selected
39 from the given candidates. In reality, we may not know how to group the variables.

40 In this paper we seek to solve a related but ultimately different and more challenging problem, which
41 we call *Composite Feature Selection*. We wish to find groups of variables *without* prior knowledge,
42 where each group acts as a separate predictive subset of the features and the overall predictive power
43 is greatest when all groups are used in unison. We call each group of features a composite feature.¹
44 By imposing this structure on the discovered features, we attempt to isolate pathways from features to
45 the response variable. [Discovering groups of features offers deeper insights](#) into *why* specific features
46 are important than standard feature selection.

47 **Contributions.** (1) We formalise *composite* feature selection as an extension of standard feature
48 selection, defining composite features in terms of linear and non-linear interactions between variables
49 (Sec. 3). (2) We propose a new deep learning architecture for composite feature selection using an
50 ensemble-based approach (Sec. 4). (3) To assess our solution, we introduce a metric for assessing
51 composite feature similarity based on Jaccard similarity. (4) We demonstrate the utility of our model
52 on a range of synthetic and semi-synthetic tasks where the ground truth group features are known
53 (Sec. 5). We see that our model not only frequently recovers the relevant features, but also often
54 discovers the underlying group structure.

55 2 Related Work

56 Significant attention has been placed on feature selection ([see Appendix B for further discussion of](#)
57 [standard feature selection](#)), and several approaches have been extended to select predefined groups
58 instead of individual features. For example, LASSO [67, 73] is a linear method that uses an L1 penalty
59 to impose sparsity among coefficients. Group LASSO [82] generalises this to allow predefined groups
60 to be selected or excluded jointly, rather than single features, by replacing the L1 penalty with L2
61 penalties on each group. Other feature selection methods, such as SLOPE [10], have been similarly
62 extended to group feature selection to give Group-SLOPE [12]. Further examples of group feature
63 selection using adapted loss functions are SCAD-L2 [84] and hierarchical LASSO [89]. [Similarly,](#)
64 [Bayesian approaches to feature selection \[25\] have also been generalised to the group setting \[31\].](#)
65 Finally, the Knockoff procedure [8, 15, 35, 51, 65, 70] is a [generative procedure that creates fake](#)
66 [covariates \(knockoffs\), obeying certain symmetries under permutations of real and knockoff features.](#)
67 [By subsequently carrying out Feature Selection on the combined real and knockoff data, there are](#)
68 [guarantees on the False Discovery Rate.](#) Generalisations of the Knockoff procedure to the group
69 setting also exist [21, 90], [where symmetries under permutations of entire groups must exist.](#)

70 The key commonality is that none of these methods *discover* groups, but instead can only *select*
71 groups from a set of predefined candidates. Therefore, while they may be applicable when we can
72 split inputs into groups, they are not able to find groups of predictors on their own. Our work differs
73 from these methods by considering the challenge of finding such groups in the absence of prior
74 knowledge. Additionally, unlike prior work, we do not make assumptions about correlations between
75 features or restrictions on groups, such as requiring that the candidate groups partition the features.

76 3 Problem Description

77 Let $\mathbf{X} \in \mathcal{X}^p$ be a p -dimensional signal (such as gene expressions or patient covariates) and $Y \in \mathcal{Y}$ be
78 a response (such as disease traits). Informally, we wish to group features into the maximum number
79 of subsets, $\mathcal{G}_i \subset [p]$, where [the predictive power of any single group significantly decreases when any](#)
80 [feature is removed](#), allowing us to separate the groups into different pathways from the signal to the
81 response. Note that we do not enforce assumptions on the groups such as non-overlapping groups or

¹We will often refer to composite features as groups for brevity; in this paper, they refer to the same thing.

82 every feature being in at least one group. In this section, we begin with a description of traditional
 83 feature selection before formalizing composite feature selection.

84 3.1 Feature Selection

85 The goal of traditional feature selection is to select a subset, $\mathcal{S} \subset [p]$, of features that are relevant for
 86 predicting the response variable. In particular, in the case of embedded feature selection [28], this is
 87 conducted jointly with the model selection process.

88 Let $*$ denote any point not in \mathcal{X} and define $\mathcal{X}_{\mathcal{S}} = (\mathcal{X} \cup \{*\})^p$. Then, given $\mathbf{X} \in \mathcal{X}^p$, the selected
 89 subset of features can be denoted as $\mathbf{X}_{\mathcal{S}} \in \mathcal{X}_{\mathcal{S}}$ where $x_{\mathcal{S},k} = x_k$ if $k \in \mathcal{S}$ and $x_{\mathcal{S},k} = *$ if $k \notin \mathcal{S}$.
 90 Let $f : \mathcal{X}_{\mathcal{S}} \rightarrow \mathcal{Y}$ be a function in **some space \mathcal{F} (such as the space of neural networks)** taking subset
 91 $\mathbf{X}_{\mathcal{S}}$ as input to yield Y . Then, selecting relevant features for predicting a response can be achieved
 92 by solving the following optimization problem:

$$\underset{f \in \mathcal{F}, \mathcal{S} \subset [p]}{\text{minimize}} \mathbb{E}_{\mathbf{x}, y \sim p_{XY}} \left[\ell_Y(y, f(\mathbf{x}_{\mathcal{S}})) \right] \text{ subject to } |\mathcal{S}| \leq \delta, \quad (1)$$

93 where δ constrains the number of selected features and $\ell_Y(y, y')$ is a task-specific loss function.

94 This can be solved by introducing a selection vector $\mathbf{M} = (M_1, \dots, M_p) \in \{0, 1\}^p$, consisting of
 95 binary random variables governed by distribution p_M , with realization \mathbf{m} indicating selection of the
 96 corresponding features. Then, the selected features given vector \mathbf{m} can be written as

$$\tilde{\mathbf{x}} \triangleq \mathbf{m} \odot \mathbf{x} + (1 - \mathbf{m}) \odot \hat{\mathbf{x}}, \quad (2)$$

97 where \odot indicates element-wise multiplication and $\hat{\mathbf{x}}$ are the values assigned to features that are not
 98 selected (typically $\hat{\mathbf{x}} \equiv 0$ or $\bar{\mathbf{x}}$). (1) can be (approximately) solved by jointly learning the model f
 99 and the selection vector distribution p_M based on the following optimization problem:

$$\underset{f, p_M}{\text{minimize}} \mathbb{E}_{\mathbf{x}, y \sim p_{XY}} \mathbb{E}_{\mathbf{m} \sim p_M} \left[\ell_Y(y, f(\tilde{\mathbf{x}})) + \beta \|\mathbf{m}\|_0 \right], \quad (3)$$

100 where β is a balancing coefficient that controls the number of features to be selected.

101 3.2 Composite Feature Selection

102 The goal of composite feature selection is to not only find the predictive features, but also to group
 103 them based on how they are predictive. For example, assume features x_1 and x_2 are only predictive
 104 when both are known by the model, but make the same prediction independent of x_3 . Then we
 105 wish to group x_1, x_2 from x_3 . In this section, we define the embedded composite feature selection
 106 problem; that is, we want to find a valid model f and groups $\{\mathcal{G}_1, \dots, \mathcal{G}_N\}$ in parallel. **A model is**
 107 **only valid when the group representations are combined in a way where we can view each group as**
 108 **contributing an independent piece of information for the final prediction. A valid model acts on a**
 109 **set of groups [83], thus when combining groups, we require order not to matter. Therefore, we must**
 110 **combine the representations using a permutation invariant aggregator.**

111 Let $A : (\prod_i \mathbb{R}^n) \rightarrow \mathbb{R}^N$ be a general permutation invariant aggregation function. It is well established
 112 that for a specific choice of $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $\rho : \mathbb{R}^m \rightarrow \mathbb{R}^N$, A can be decomposed as $\rho[\sum_i \phi(\cdot)]$
 113 (see [83] for examples). This gives $f(\mathbf{x}) = g(\rho[\sum_i \phi(f_i(\mathbf{x}_{\mathcal{G}_i}))])$, where f_i encodes group i , ρ
 114 and ϕ give the permutation invariant aggregation and g is any final non-linear function, softmax for
 115 instance. The function composition of ϕ and f_i can be relabelled as $\tilde{f}_i = \phi \circ f_i$, and the composition
 116 of g and ρ can be relabelled as $\tilde{\rho} = g \circ \rho$. This leads to $f(\mathbf{x}) = \tilde{\rho}[\sum_i \tilde{f}_i(\mathbf{x}_{\mathcal{G}_i})]$, this gives a definition
 117 for what a valid model structure can be in composite feature selection.

118 **Definition 3.1.** The most general valid model for acting on N composite features is given by:

$$f(\mathbf{x}) = \rho \left[\sum_{i=1}^N f_i(\mathbf{x}_{\mathcal{G}_i}) \right]. \quad (4)$$

119 That is, the groups must interact only once, all groups must be included and the interaction is a
 120 summation; all other interactions can (and often should) be non-linear.

121 Depending on the task, a specific permutation invariant aggregation may be chosen (e.g. $\text{Max}()$).
 122 However, any permutation invariant aggregator can be (approximately) expressed in the form of Def.
 123 3.1; thus, when learning from data, the general structure of Def. 3.1 means that this is not necessary.

124 The embedded composite feature selection problem can now be phrased in an analogous way to
 125 traditional feature selection. Let $*$ denote some point not in \mathcal{X} and define $\mathcal{X}_{\mathcal{G}_i} = (\mathcal{X} \cup \{*\})^p$. Then,
 126 given $\mathbf{X} \in \mathcal{X}^p$, the selected group of features is denoted as $\mathbf{X}_{\mathcal{G}_i} \in \mathcal{X}_{\mathcal{G}_i}$ where $x_{\mathcal{G}_i,k} = x_k$ if $k \in \mathcal{G}_i$
 127 and $x_k = *$ if $k \notin \mathcal{G}_i$. Let $f_i : \mathcal{X}_{\mathcal{G}_i} \rightarrow \mathcal{Z}$ be a function in \mathcal{F} that takes as input the subset $\mathbf{X}_{\mathcal{G}_i}$ and
 128 outputs a latent representation \mathbf{z}_i . Then, finding the groups of features can be achieved by solving the
 129 optimization problem:

$$\underset{\rho, f_i \in \mathcal{F}, \mathcal{G}_i \subset [p]}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, y \sim p_{XY}} \left[\ell_Y \left(y, \rho \left[\sum_{i=1}^N f_i(\mathbf{x}_{\mathcal{G}_i}) \right] \right) \right] \quad \text{subject to} \quad \begin{cases} |\mathcal{G}_i| \leq \delta_i \quad \forall i, \\ N \geq \Delta, \end{cases} \quad (5)$$

130 where δ constrains the number of selected features in each group and Δ gives the minimum number
 131 of groups. This objective leads to multiple smaller groups, rather than one group containing all
 132 features, which is consistent with our motivation of the problem.

133 Continuing to expand from traditional feature selection, we can also extend the solution to the
 134 composite setting. For N groups we can introduce a selection *matrix* $M \in \{0, 1\}^{N \times p}$, governed by
 135 distribution p_M . For a realization \mathbf{m} , the selected features from group i are given by

$$\tilde{\mathbf{x}}_i \triangleq \mathbf{m}_i \odot \mathbf{x} + (1 - \mathbf{m}_i) \odot \hat{\mathbf{x}}, \quad (6)$$

136 where \mathbf{m}_i is the i^{th} row of M . We can approximately solve (5) by solving the optimization problem:

$$\underset{f, p_M}{\text{minimize}} \quad \mathbb{E}_{\mathbf{x}, y \sim p_{XY}} \mathbb{E}_{\mathbf{m} \sim p_M} \left[\ell_Y(y, f(\mathbf{x})) + R_e(\mathbf{m}) \right], \quad (7)$$

137 where $f(\mathbf{x})$ obeys Def. (4) and R_e is a regularisation term which controls how features are selected
 138 in each group. R_e should capture both group size (i.e. encourage as few features as possible to be
 139 selected) but also the relationships between groups (i.e. groups should be distinct and not redundant).

140 3.3 Challenges

141 There are various challenges in solving the composite feature selection problem. While the ultimate
 142 task is to find predictive groups of features, there remains the necessity simply to identify predictive
 143 features, which is already an NP-hard problem [2]. Composite feature selection not only inherits
 144 this property but introduces additional complexity since we can think of each group as solving a
 145 separate feature selection problem. Consider the number of potential solutions: in traditional feature
 146 selection (assuming not all features are selected), there are $2^n - 2$ ways of selecting a subset from n
 147 features; even restricting to at most $m \ll n$ quickly becomes unfeasible for even modest values
 148 of m . In composite feature selection, *every group* has the same number of solutions as traditional
 149 feature selection, drastically increasing the total number of possible solutions. A challenge specific
 150 to composite feature selection arises when the ground truth group structure contains groups with
 151 overlapping features (e.g. feature x_1 interacts independently with both x_2 and x_3). In this scenario, it
 152 is difficult to separate these two effects while penalizing the inclusion of additional features.

153 4 Method: CompFS

154 In this section, we propose a novel architecture for finding predictive groups of features, which we
 155 refer to as **Composite Feature Selection (CompFS)**. In order to discover groups of features, our model
 156 is composed of a set of group selection models and an aggregate predictor. Our approach resembles
 157 an ensemble of “weak” feature selection models, where each learner attempts to solve the task using
 158 a sparse set of features (Figure 1). These models are then trained in such a way as to discover distinct
 159 predictive groups. We first consider the group selection models in more detail before describing how
 160 the group selection models are combined and the training procedure.

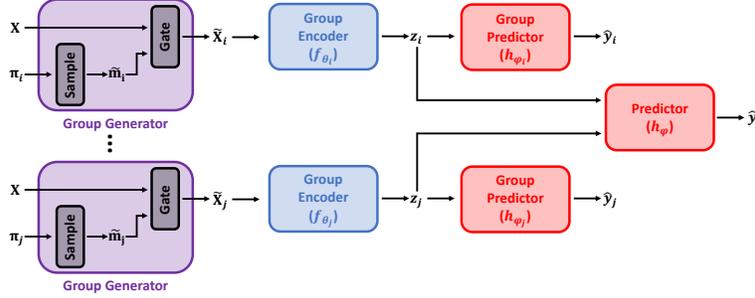


Figure 1: An illustration of CompFS. We use an ensemble of group selection models to discover composite features and an aggregate predictor to combine these features when issuing predictions.

161 4.1 Group Selection Models

162 CompFS is composed of a set of group selection models, each of which primarily aims to solve the
 163 traditional feature selection problem specified by (1). We achieve this by solving (3) using a neural
 164 network-based approach with stochastic gating of the input features. Each group selection model
 165 consists of the following three components (Figure 1):

- 166 • *Group Selection Probability*, $\pi_i = (\pi_{1,i}, \dots, \pi_{p,i}) \in [0, 1]^p$, which is a trainable vector that
 167 governs the Bernoulli distribution used to generate the gate vector \mathbf{m}_i . Each element of the
 168 selection probability $\pi_{k,i}$ indicates the importance of the corresponding feature to the target.
- 169 • *Group Encoder*, $f_{\theta_i} : \mathcal{X}^p \rightarrow \mathcal{Z}$, that takes as input the selected subset of features $\tilde{\mathbf{x}}_i$ and outputs
 170 latent representations $\mathbf{z}_i \in \mathcal{Z}$.
- 171 • *Group Predictor*, $h_{\phi_i} : \mathcal{Z} \rightarrow \mathcal{Y}$, that takes as input the latent representations of the selected subset
 172 of features, $\mathbf{z}_i = f_{\theta_i}(\tilde{\mathbf{x}}_i)$, and outputs predictions on the target outcome.

173 Solving (3) directly is not possible since the sampling step has no differentiable inverse. Instead, we
 174 use the relaxed Bernoulli distribution [53, 34] and apply the reparameterization trick as follows.

175 Formally, given selection probability $\pi = (\pi_1, \dots, \pi_p)$ and independent Uniform(0, 1) random
 176 variables (U_1, \dots, U_p) , we can generate a relaxed gate vector $\tilde{\mathbf{m}} = (\tilde{m}_1, \dots, \tilde{m}_p) \in (0, 1)^p$ based
 177 on the following reparameterization trick [53]:

$$\tilde{m}_k = \sigma\left(\frac{1}{\tau} (\log \pi_k - \log(1 - \pi_k) + \log U_k - \log(1 - U_k))\right), \quad (8)$$

178 where $\sigma(x) = (1 + \exp(-x))^{-1}$ is the sigmoid function. This relaxation is parameterized by π and
 179 temperature $\tau \in (0, \infty)$. Further, as $\tau \rightarrow 0$, the gate vectors \tilde{m}_k converge to Bernoulli(π_k) random
 180 variables. Crucially this is differentiable with respect to π .

181 Given group selection probability π_i , we first sample relaxed Bernoulli random variable $\tilde{\mathbf{m}}_i$ according
 182 to (8) and then use $\tilde{\mathbf{m}}_i$ in a gating procedure to select the group of features. The output of the gate is:

$$\tilde{\mathbf{x}}_i = \text{gate}_i(\mathbf{x}) = \tilde{\mathbf{m}}_i \odot \mathbf{x} + (1 - \tilde{\mathbf{m}}_i) \odot \bar{\mathbf{x}}, \quad (9)$$

183 where we replace the variables that were not selected by their mean value $\bar{\mathbf{x}}$. The mean is used because
 184 in particular tasks a feature having a value of 0 may be particularly meaningful. The gate output $\tilde{\mathbf{x}}_i$
 185 is then fed into the group encoder f_{θ_i} to yield representation $\mathbf{z}_i = f_{\theta_i}(\tilde{\mathbf{x}}_i)$. This representation is finally
 186 passed to the group predictor h_{ϕ_i} to produce the prediction for an individual learner, $\hat{y}_i = h_{\phi_i}(\mathbf{z}_i)$.

187 4.2 Group Aggregation

188 The final component necessary for CompFS is a way to aggregate the individual group selection
 189 models. This is achieved via an overall *predictor*, $h_{\phi} : \mathcal{Z} \rightarrow \mathcal{Y}$, that takes as input the set of latent
 190 representations $\{\mathbf{z}_1, \dots, \mathbf{z}_N\}$ produced by the individual learners and outputs predictions on the
 191 target outcome. For simplicity, we apply a linear prediction head to the latent representations and use
 192 element-wise summation to aggregate. Thus, the prediction of the ensemble is given by:

$$\hat{y} = h_{\phi}(\{\mathbf{z}_1, \dots, \mathbf{z}_N\}) = \rho \left[\sum_{i=1}^N \mathbf{W}_i \mathbf{z}_i + \mathbf{b}_i \right], \quad (10)$$

193 where N is the number of members of the ensemble and ρ is a suitable transformation (e.g. softmax).
 194 Note that by using element-wise summation, our model satisfies (4) for acting on composite features.

195 4.3 Loss Functions

196 **Group Selection Models.** The individual learners can be trained to perform (traditional) feature
 197 selection (1) by minimizing the following loss function:

$$\mathcal{L}_{\mathcal{G}_i} = \mathbb{E}_{\mathbf{x}, y \sim p_{XY}} \left[\ell_Y(y, h_{\phi_i}(f_{\theta_i}(\text{gate}_i(\mathbf{x})))) + \beta \langle \boldsymbol{\pi}_i \rangle^2 \right], \quad (11)$$

198 where ℓ_Y is a suitable loss function for the prediction task (e.g. cross-entropy for classification tasks
 199 and MSE for regression tasks) and $\beta \geq 0$ balances the two terms. Note the selections probabilities $\boldsymbol{\pi}_i$
 200 are not regularized with the typical L1 penalty. Instead, we apply an L2 penalty to the mean selection
 201 probability $\langle \boldsymbol{\pi}_i \rangle$ for each individual learner. This is justified as follows. Recall the optimization
 202 problem given by (5). We desire a solution with the maximal number of predictive groups N , while
 203 minimizing the number of selected features per group $\sum_{i=1}^N |\mathcal{G}_i|$. The standard L1 penalty term does
 204 not achieve this goal since adding an additional feature to either group \mathcal{G}_i or \mathcal{G}_j incurs the same
 205 penalty. In contrast, the L2 penalty imposed on $\langle \boldsymbol{\pi}_i \rangle$ penalizes adding extra features to already large
 206 groups, favoring the construction of smaller groups over larger ones.

207 **Aggregate Predictor.** The aggregate predictor can be trained jointly with the group feature selection
 208 models by minimizing a standard prediction loss (where ℓ_Y is the same as in (11)):

$$\mathcal{L}_E = \mathbb{E}_{\mathbf{x}, y \sim p_{XY}} \left[\ell_Y(y, h_{\phi}(\{\mathbf{z}_1, \dots, \mathbf{z}_n\})) \right]. \quad (12)$$

209 **Additional Regularization.** If we simply apply the losses given by Eqs. (11), (12), there will be
 210 limited (or even no) differentiation among the individual learners and the optimal solution would
 211 be for each learner to simply solve the traditional feature selection problem (1). This results in
 212 all learners selecting the same features, which does not achieve our aim of discovering groups of
 213 predictive features. In order to encourage differentiation between the models, we introduce an
 214 additional loss that penalizes the selection of the same features in multiple groups:

$$\mathcal{L}_R = \mathbb{E}_{\mathbf{x}, y \sim p_{XY}} \left[\sum_{i=1}^N \sum_{j>i} \boldsymbol{\pi}_i \cdot \boldsymbol{\pi}_j \right]. \quad (13)$$

215 **Overall Loss.** Combining the above, our overall loss function therefore can be written as follows:

$$\mathcal{L} = \sum_{i=1}^N \mathcal{L}_{\mathcal{G}_i} + \beta_E \mathcal{L}_E + \beta_R \mathcal{L}_R, \quad (14)$$

216 where $\beta_E, \beta_R \geq 0$ are hyperparameters to balance the losses.

217 Training CompFS with the loss given by (14) is designed to achieve the following: (1) The overall
 218 ensemble network should be a good predictor (\mathcal{L}_E). (2) Each individual learner should to solve the
 219 traditional feature selection problem ($\mathcal{L}_{\mathcal{G}_i}$), which requires the group predictor to be accurate while
 220 selecting minimal features. (3) Finally, we want the groups to be distinct and thus discourage highly
 221 similar groups (\mathcal{L}_R). However, note that we do not exclude the possibility of some overlap of features
 222 between groups. The model is end-to-end differentiable, so we train with gradient descent.

223 **Evaluation.** During evaluation, only the gating procedure changes. The way features can be selected
 224 is chosen by the user. A standard solution which we adopt in this paper is using a threshold λ and
 225 calculating the gate vector \mathbf{m}_i as follows: $m_{k,i} = 1$, if $\pi_{k,i} > \lambda$ and 0 otherwise.

226 5 Experiments

227 We evaluate CompFS using several synthetic and semi-synthetic datasets where ground truth feature
 228 importances and group structure are known. Specific architectural details are given in App. C.

229 Additional information regarding experiments, benchmarks, and datasets can be found in App. D.
 230 Additional ablations and sensitivity analysis are in App. A.

231 **Benchmarks.** The primary goal of our experiments is to demonstrate the utility of discovering
 232 composite features over traditional feature selection. Our main benchmark is an oracle feature
 233 selection method (“Oracle”) that perfectly selects the ground truth features but provides no structure,
 234 giving all features as one group. We also include comparisons to a linear feature selection method
 235 (LASSO) [73] and two non-linear, state of the art approaches, Stochastic Gates (STG) [81] and
 236 Supervised Concrete Autoencoder (Sup-CAE) [6]. Finally, we compare with Group LASSO [82],
 237 where we enumerate all groups with 1 or 2 features as predefined groups. Note this represents a
 238 significant simplification of the task for Group Lasso (see App. G for additional baselines).

239 **Metrics.** Since the ground truth feature groups $\mathcal{G}_1, \dots, \mathcal{G}_N$ are known, we use True Positive Rate
 240 (TPR) and False Discovery Rate (FDR) to assess the discovered features against the ground truth. To
 241 assess composite features, i.e. grouping, we define the Group Similarity (G_{sim}) as the normalized
 242 Jaccard similarity between ground truth feature groups and the most similar proposed group:

$$G_{\text{sim}} = \frac{1}{\max(N, K)} \sum_{i=1}^N \max_{j \in [K]} \mathcal{J}(\mathcal{G}_i, \hat{\mathcal{G}}_j), \quad (15)$$

243 where \mathcal{J} is the Jaccard index [33] and $\hat{\mathcal{G}}_1, \dots, \hat{\mathcal{G}}_K$ are the discovered groups. $G_{\text{sim}} \in [0, 1]$, where
 244 $G_{\text{sim}} = 1$ corresponds to perfect recovery of the ground truth groups, while $G_{\text{sim}} = 0$ when none of the
 245 correct features are discovered. See App. E for additional details together with examples. We assess
 246 the models by seeing if the ground truth features have been correctly discovered, using TPR and FDR.
 247 We then see if the underlying grouping has been uncovered (and correct features) using G_{sim} .

248 5.1 Synthetic Experiments.

249 **Dataset Description.** We begin by evaluating our method on a range of synthetic datasets where
 250 the ground truth feature importance is known (Table 1). We generate synthetic datasets by sampling
 251 from the Gaussian distribution with initially no correlations across the features ($X \sim \mathcal{N}(0, I)$). We
 252 construct binary classification tasks, where the class y is determined by the following decision rules:

- 253 • **(Syn1)** $y = 1$ if $x_1 > 0.55$ or $x_2 > 0.55$, 0 otherwise. The ground truth groups are $\{\{1\}, \{2\}\}$.
 254 This task assesses whether the model can separate two features rather than group them together.
- 255 • **(Syn2)** $y = 1$ if $x_1x_2 > 0.30$ or $x_3x_4 > 0.30$, 0 otherwise. The ground truth groups are
 256 $\{\{1, 2\}, \{3, 4\}\}$. This task requires identifying groups consisting of more than one variable.
- 257 • **(Syn3)** $y = 1$ if $x_1x_2 > 0.30$ or $x_1x_3 > 0.30$, 0 otherwise. The ground truth groups are
 258 $\{\{1, 2\}, \{1, 3\}\}$. This task investigates whether a model can split the features into two *overlapping*
 259 groups of two, rather than one group with all three features.
- 260 • **(Syn4)** $y = 1$ if $x_1x_4 > 0.30$ or $x_7x_{10} > 0.30$, 0 otherwise. The ground truth groups are
 261 $\{\{1, 4\}, \{7, 10\}\}$. This task is equivalent to **Syn2**, however, here the features exhibit strong
 262 correlation in collections of 3. This task demonstrates the difficulty of carrying out group feature
 263 selection (and indeed standard feature selection) when the features are highly correlated.

264 The decision rules are created such that there is minimal class imbalance. We use signals with 500
 265 dimensions to demonstrate the utility in the high dimensional regime. We use 20000 samples to train
 266 and 200 to test. Each experiment is repeated 10 times.

267 **Analysis.** On both Syn1 and Syn2, CompFS achieves high TPR with no false discoveries (0%
 268 FDR) and significantly higher G_{sim} than the Oracle. Despite allowing CompFS to discover up to
 269 5 groups, CompFS typically finds the correct number of groups (2), demonstrating that it is not
 270 necessary for the number of allowed composite features to match the ground truth, which is vital
 271 in real-world use cases where this is unknown. Syn3 is significantly more challenging due to the
 272 overlapping structure and we observe essentially the same performance as Oracle. Despite finding all
 273 the correct features and no false discoveries, CompFS typically finds the union $\{1, 2, 3\}$ rather than
 274 the underlying group structure $\{\{1, 2\}, \{1, 3\}\}$. Finally, for Syn4, while CompFS has a relatively
 275 high FDR, it frequently finds the ground truth relevant features and groups with similar G_{sim} to Oracle.

Table 1: Performance on Synthetic Datasets, values are recorded with their standard deviations.

DATASET	MODEL	TPR	FDR	G _{sim}	NO. GROUPS	ACCURACY (%)
SYN1	COMPFS(5)	100.0 ± 0.0	0.0 ± 0.0	0.91 ± 0.14	2.2 ± 0.4	98.9 ± 0.5
	ORACLE	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	LASSO	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	81.8 ± 2.0
	GROUP LASSO	100.0 ± 0.0	0.0 ± 0.0	0.67 ± 0.00	3.0 ± 0.0	83.8 ± 1.4
	STG	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	97.8 ± 1.4
	SUP-CAE	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	97.8 ± 1.4
SYN2	COMPFS(5)	95.0 ± 15.0	0.0 ± 0.0	0.90 ± 0.20	1.8 ± 0.4	95.5 ± 5.4
	ORACLE	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	LASSO	0.0 ± 0.0	0.0 ± 0.0	0.00 ± 0.00	0.0 ± 0.0	52.6 ± 2.9
	GROUP LASSO	0.0 ± 0.0	0.0 ± 0.0	0.00 ± 0.00	0.0 ± 0.0	52.2 ± 0.9
	STG	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	93.9 ± 2.2
	SUP-CAE	37.5 ± 31.7	42.5 ± 44.2	0.24 ± 0.20	1.0 ± 0.0	61.9 ± 12.8
SYN3	COMPFS(5)	100.0 ± 0.0	0.0 ± 0.0	0.68 ± 0.05	1.3 ± 0.5	97.4 ± 1.1
	ORACLE	100.0 ± 0.0	0.0 ± 0.0	0.67 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	LASSO	0.0 ± 0.0	0.0 ± 0.0	0.00 ± 0.00	0.0 ± 0.0	56.5 ± 4.0
	GROUP LASSO	0.0 ± 0.0	0.0 ± 0.0	0.00 ± 0.00	0.0 ± 0.0	54.6 ± 1.3
	STG	100.0 ± 0.0	0.0 ± 0.0	0.67 ± 0.00	1.0 ± 0.0	95.3 ± 1.7
	SUP-CAE	23.3 ± 31.6	66.7 ± 47.1	0.23 ± 0.31	1.0 ± 0.0	62.6 ± 12.6
SYN4	COMPFS(5)	90.0 ± 12.2	51.9 ± 13.8	0.47 ± 0.20	2.5 ± 0.7	95.8 ± 1.8
	ORACLE	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	LASSO	0.0 ± 0.0	0.0 ± 0.0	0.00 ± 0.00	0.0 ± 0.0	51.8 ± 3.2
	GROUP LASSO	0.0 ± 0.0	10.0 ± 31.6	0.00 ± 0.00	0.1 ± 0.3	53.0 ± 1.1
	STG	100.0 ± 0.0	66.7 ± 0.0	0.17 ± 0.00	1.0 ± 0.0	94.2 ± 2.1
	SUP-CAE	72.5 ± 14.2	16.7 ± 14.7	0.39 ± 0.08	1.0 ± 0.0	72.2 ± 13.2

276 This is a challenging task with significant correlation between features. Despite this, CompFS is
 277 able to uncover the underlying group structure, providing additional insight over traditional feature
 278 selection. **STG typically performs well in terms of traditional feature selection, but scores poorly in**
 279 **terms of G_{sim} due to not providing any group information.**

280 5.2 Semi-Synthetic Experiments.

281 **Dataset Description.** Next, we assess our ability to identify composite features using semi-synthetic
 282 molecular datasets. These tasks are analogs of real-world problems, such as identifying biologically
 283 active chemical groups; however, the labels are determined by a synthetic “binding logic” so that the
 284 ground truth feature relevance is known. We use several of the datasets constructed by [56], some of
 285 which were also used by [66].² The synthetic “binding logics” are expressed as a combination of
 286 molecular fragments that must either be present or absent for binding to occur and are used to label
 287 molecules from the ZINC database [32]. Each logic includes up to four functional groups (Table 6).
 288 Molecules are featurized using a set of 84 functional groups, where feature $x_i = 1$ if the molecule
 289 contains functional group i and 0 otherwise. The specific binding logics are given in App. F.

Table 2: Performance on Chemistry Datasets, values are recorded with their standard deviations.

DATASET	MODEL	TPR	FDR	G _{sim}	NO. GROUPS	ACCURACY (%)
CHEM1	COMPFS(5)	100.0 ± 0.0	0.0 ± 0.0	0.82 ± 0.20	1.9 ± 0.5	100.0 ± 0.0
	ORACLE	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	LASSO	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	75.8 ± 0.0
	GROUP LASSO	100.0 ± 0.0	0.0 ± 0.0	0.67 ± 0.00	3.0 ± 0.0	100.0 ± 0.0
	STG	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	SUP-CAE	62.5 ± 13.2	23.3 ± 17.5	0.37 ± 0.07	1.0 ± 0.0	77.8 ± 11.0
CHEM2	COMPFS(5)	100.0 ± 0.0	0.0 ± 0.0	0.72 ± 0.24	2.2 ± 0.6	100.0 ± 0.0
	ORACLE	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	LASSO	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	81.6 ± 0.0
	GROUP LASSO	100.0 ± 0.0	0.0 ± 0.0	0.40 ± 0.00	5.0 ± 0.0	81.6 ± 0.0
	STG	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	SUP-CAE	66.7 ± 0.0	0.0 ± 0.0	0.42 ± 0.00	1.0 ± 0.0	80.9 ± 9.5
CHEM3	COMPFS(5)	100.0 ± 0.0	7.3 ± 11.7	0.62 ± 0.17	2.4 ± 0.5	100.0 ± 0.0
	ORACLE	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	LASSO	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	87.4 ± 5.2
	GROUP LASSO	100.0 ± 0.0	20.0 ± 0.0	0.20 ± 0.00	10.0 ± 0.0	91.5 ± 0.0
	STG	100.0 ± 0.0	0.0 ± 0.0	0.50 ± 0.00	1.0 ± 0.0	100.0 ± 0.0
	SUP-CAE	62.5 ± 13.2	23.3 ± 17.5	0.37 ± 0.07	1.0 ± 0.0	77.8 ± 11.0

²Data from https://github.com/google-research/graph-attribution/raw/main/data/all_16_logics_train_and_test.zip.

290 **Analysis.** All methods are able to identify the ground truth relevant features; however, only CompFS
291 provides deeper insights. Unlike for Syn1-4, LASSO correctly selects the ground truth features
292 since the dataset consists of binary variables and thus it is possible to find performant linear models.
293 However, while discovering the correct features, Group LASSO selects all possible combinations of
294 these features, adding no benefit over standard feature selection.

295 For Chem1-2, CompFS perfectly recovers the group structure in the majority of experiments, leading
296 to high G_{sim} far exceeding traditional feature selection. On Chem3, we occasionally discover
297 additional features that are not part of the binding logic. However, a number of molecular fragments
298 are strongly correlated with the binding logic, even though they are not themselves included. In fact,
299 some features contain information about *multiple* functional groups. For example, esters contain
300 a carbonyl and an ether; both are in the binding logic for Chem3, while ester is not, despite being
301 highly informative, and thus occasionally CompFS incorrectly selects this feature. In spite of this,
302 CompFS achieves significantly higher G_{sim} than even Oracle. This demonstrates the benefit of
303 the grouping discovered by CompFS, even with a modest number of false discoveries. As before,
304 CompFS typically finds the correct number of groups (2), despite being able to discover up to 5
305 groups, further demonstrating that the number of composite features need not be known *a priori*,
306 which is the case in real-world applications.

307 5.3 Real-World Data: METABRIC

308 **Dataset Description.** Finally, we assess CompFS on a real-
309 world dataset, METABRIC [19, 60], where the ground truth
310 group structure is *unknown*. METABRIC contains gene ex-
311 pression, mutation, and clinical data for 1,980 primary breast
312 cancer samples. We evaluated the ability to predict the proges-
313 terone receptor (PR) status of the tissue based on the gene ex-
314 pression data, which consists of measurements for 489 genes.

315 **Analysis.** CompFS suffers limited performance degradation compared to using all features, despite
316 only using 5% of the features (Table 3). Despite imposing a more rigid structural form on how
317 features can interact in the predictive model, STG only had marginally greater predictive power than
318 CompFS. However, CompFS provides greater insight into how the features interact than STG.

319 We found supporting evidence in the scientific literature for all but 1 of the genes discovered by
320 CompFS (Table 10). In addition, within each group, we found further evidence of the interactions
321 between genes, demonstrating the ability for CompFS to learn informative groups of features. For
322 example, in Group 1, CXCR1 and PEN-2 (the protein encoded by PSENEN) are known to interact
323 [5]. In Group 2, BMP6 encodes a member of the TGF- β superfamily of proteins, and TGF- β
324 triggers activation of SMAD3 [17]. In the same group, MAPK1 activity is dependent on the activity
325 of PRKCQ in breast cancer cells [13], while MAPK1 is also known to interact with MAPT [45],
326 SMAD3 [23], and BMP6 [85]. Additional supporting evidence can be found in Appendix H.

327 6 Conclusion

328 In this paper, we introduced CompFS, an ensemble-based approach that tackles the newly proposed
329 challenge of composite feature selection. Using synthetic and semi-synthetic data, we assess our
330 ability to go beyond traditional feature selection and recover deeper underlying connections between
331 variables. CompFS is not without limitations: as with other methods, points of difficulty arise when
332 features are highly correlated, or if predictive composites contain overlapping features. Future work
333 may overcome this by using correlated gates. Further, as with many traditional feature selection
334 methods, there are no guarantees on false discovery rate. This could be tackled by first proposing
335 candidate composite features, and then using the Group Knockoff procedure. Additionally, to discover
336 groups, CompFS requires the introduction of additional hyperparameters which could be challenging
337 to tune in practice. More broadly, as with standard feature selection, groups found under composite
338 feature selection must be verified by domain experts (both features but additionally interactions).
339 However, we believe the additional structure provided by composite feature selection could be of
340 significant benefit to a wide variety of practitioners.

Table 3: METABRIC performance. We compare CompFS and STG using 25 features to an MLP using all 489 features.

Model	AUC ROC
MLP (All features)	0.869
CompFS(5)	0.830
STG	0.843

References

- [1] Wail Al Sarakbi, Sara Reefy, Wen G. Jiang, Terry Roberts, Robert F. Newbold, and Kefah Mokbel. Evidence of a tumour suppressor function for DLEC1 in human breast cancer. *Anticancer Research*, 30(4):1079–1082, 2010.
- [2] Edoardo Amaldi and Viggo Kann. On the approximability of minimizing nonzero variables or unsatisfied relations in linear systems. *Theoretical Computer Science*, 209(1):237–260, 1998.
- [3] M. Ampuja, E.L. Alarmo, P. Owens, R. Havunen, A.E. Gorska, H.L. Moses, and A. Kallioniemi. The impact of bone morphogenetic protein 4 (BMP4) on breast cancer metastasis in a mouse xenograft model. *Cancer Letters*, 375(2):238–244, 2016.
- [4] Sharon Arcuri, Georgia Pennarossa, Fulvio Gandolfi, and Tiziana A. L. Brevini. Generation of trophoblast-like cells from hypomethylated porcine adult dermal fibroblasts. *Frontiers in Veterinary Science*, 8, 2021.
- [5] Martina Bakele, Amelie S. Lotz-Havla, Anja Jakowetz, Melanie Carevic, Veronica Marcos, Ania C. Muntau, and Dominik Gersting, Soeren W.and Hartl. An interactive network of elastase, secretases, and PAR-2 protein regulates CXCR1 receptor surface expression on neutrophils. *Journal of Biological Chemistry*, 289(30):20516–20525, 2014.
- [6] Muhammed Fatih Balin, Abubakar Abid, and James Zou. Concrete autoencoders: Differentiable feature selection and reconstruction. In *International Conference on Machine Learning (ICML)*, 2019.
- [7] Hatice Ozer Balli and Bent E. Sørensen. Interaction effects in econometrics. *Empirical Economics*, 45(1):583–603, 2013.
- [8] Rina Foygel Barber and Emmanuel J Candès. Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5):2055–2085, 2015.
- [9] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [10] Malgorzata Bogdan, Ewout van den Berg, Weijie Su, and Emmanuel Candes. Statistical estimation and testing via the sorted L1 norm. *arXiv preprint arXiv:1310.1969*, 2013.
- [11] Terry Brown. Silica exposure, smoking, silicosis and lung cancer—complex interactions. *Occupational Medicine*, 59(2):89–95, 03 2009.
- [12] Damian Brzyski, Alexej Gossmann, Weijie Su, and Małgorzata Bogdan. Group SLOPE - adaptive selection of groups of predictors. *Journal of the American Statistical Association*, 114(525):419–433, 2019.
- [13] Jessica Byerly, Gwyneth Halstead-Nussloch, Koichi Ito, Igor Katsyv, and Hanna Y. Irie. PRKCQ promotes oncogenic growth and anoikis resistance of a subset of triple-negative breast cancer cells. *Breast Cancer Research*, 18(1):95, 2016.
- [14] Jessica H. Byerly, Elisa R. Port, and Hanna Y. Irie. PRKCQ inhibition enhances chemosensitivity of triple-negative breast cancer by regulating Bim. *Breast Cancer Research*, 22(1):72, 2020.
- [15] Emmanuel Candes, Yingying Fan, Lucas Janson, and Jinchi Lv. Panning for gold: ‘model-X’ knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(3):551–577, 2018.
- [16] Girish Chandrashekar and Ferat Sahin. A survey on feature selection methods. *Computers & Electrical Engineering*, 40(1):16–28, 2014.

- 383 [17] Bijun Chen, Ruoshui Li, Silvia C. Hernandez, Anis Hanna, Kai Su, Arti V. Shinde, and
384 Nikolaos G. Frangogiannis. Differential effects of smad2 and smad3 in regulation of macrophage
385 phenotype and function in the infarcted myocardium. *Journal of Molecular and Cellular*
386 *Cardiology*, 171:1–15, 2022.
- 387 [18] Tianyi Cheng, Peiyong Chen, Jingyi Chen, Yingtong Deng, and Chen Huang. Landscape
388 analysis of matrix metalloproteinases unveils key prognostic markers for patients with breast
389 cancer. *Frontiers in Genetics*, 12, 2022.
- 390 [19] Christina Curtis, Sohrab P. Shah, Suet-Feung Chin, Gulisa Turashvili, Oscar M. Rueda, Mark J.
391 Dunning, Doug Speed, Andy G. Lynch, Shamith Samarajiwa, Yinyin Yuan, Stefan Gräf, Gavin
392 Ha, Gholamreza Haffari, Ali Bashashati, Roslin Russell, Steven McKinney, METABRIC
393 Group, Anita Langerød, Andrew Green, Elena Provenzano, Gordon Wishart, Sarah Pinder,
394 Peter Watson, Florian Markowetz, Leigh Murphy, Ian Ellis, Arnie Purushotham, Anne-Lise
395 Børresen-Dale, James D. Brenton, Simon Tavaré, Carlos Caldas, and Samuel Aparicio. The
396 genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups.
397 *Nature*, 486(7403):346–352, 2012.
- 398 [20] Kun Dai, Hong-Yi Yu, and Qing Li. A semisupervised feature selection with support vector
399 machine. *Journal of Applied Mathematics*, 64:141–158, 2013.
- 400 [21] Ran Dai and Rina Barber. The knockoff filter for FDR control in group-sparse and multitask
401 regression. In *International Conference on Machine Learning (ICML)*, 2016.
- 402 [22] Mei Dong, Tam How, Kellye C. Kirkbride, Kelly J. Gordon, Jason D. Lee, Nadine Hempel,
403 Patrick Kelly, Benjamin J. Moeller, Jeffrey R. Marks, and Gerard C. Blobe. The type III
404 TGF- β receptor suppresses breast cancer progression. *The Journal of Clinical Investigation*,
405 117(1):206–217, 2007.
- 406 [23] Wei Bin Fang, Iman Jokar, An Zou, Diana Lambert, Prasanthi Dendukuri, and Nikki Cheng.
407 CCL2/CCR2 chemokine signaling coordinates survival and motility of breast cancer cells
408 through smad3 protein- and p42/44 mitogen-activated protein kinase (MAPK)-dependent mech-
409 anisms. *Journal of Biological Chemistry*, 287(43):36593–36608, 2012.
- 410 [24] Michael Y. Fessing, Ruzanna Atoyan, Ben Shander, Andrei N. Mardaryev, Vladimir
411 V. Botchkarev Jr., Krzysztof Poterlowicz, Yonghong Peng, Tatiana Efimova, and Vladimir A.
412 Botchkarev. BMP signaling induces cell-type-specific changes in gene expression programs of
413 human keratinocytes and fibroblasts. *Journal of Investigative Dermatology*, 130(2):398–404,
414 2010.
- 415 [25] Edward I. George and Robert E. McCulloch. Approaches for bayesian variable selection.
416 *Statistica Sinica*, pages 339–373, 1997.
- 417 [26] Christophe Ginestier, Suling Liu, Mark E. Diebel, Hasan Korkaya, Ming Luo, Marty Brown,
418 Julien Wicinski, Olivier Cabaud, Emmanuelle Charafe-Jauffret, Daniel Birnbaum, Jun-Lin Guan,
419 Gabriela Dontu, and Max S. Wicha. CXCR1 blockade selectively targets human breast cancer
420 stem cells in vitro and in xenografts. *The Journal of Clinical Investigation*, 120(2):485–497,
421 2010.
- 422 [27] Santiago M. Gómez Bergna, Abril Marchesini, Leslie C. Amorós Morales, Paula N. Arrías,
423 Hernán G. Farina, Víctor Romanowski, M. Florencia Gottardo, and Matias L. Pidre. Exploring
424 the metastatic role of the inhibitor of apoptosis BIRC6 in breast cancer. *bioRxiv*, 2021.
- 425 [28] Isabelle Guyon and André Elisseeff. An introduction to variable and feature selection. *Journal*
426 *of Machine Learning Research*, 3(1):1157–1182, 2003.
- 427 [29] Kai Han, Yunhe Wang, Chao Zhang, Chao Li, and Chao Xu. Autoencoder inspired unsupervised
428 feature selection. In *IEEE International Conference on Acoustics, Speech and Signal Processing*
429 *(ICASSP)*, 2018.

- 430 [30] Xiaofei He, Deng Cai, and Partha Niyogi. Laplacian score for feature selection. In *Advances in*
431 *Neural Information Processing Systems (NeurIPS)*, 2005.
- 432 [31] Daniel Hernández-Lobato, José Miguel Hernández-Lobato, and Pierre Dupont. Generalized
433 spike-and-slab priors for bayesian group feature selection using expectation propagation. *Journal*
434 *of Machine Learning Research*, 14(7), 2013.
- 435 [32] John J. Irwin, Teague Sterling, Michael M. Mysinger, Erin S. Bolstad, and Ryan G. Coleman.
436 ZINC: A free tool to discover chemistry for biology. *Journal of Chemical Information and*
437 *Modeling*, 52(7):1757–1768, 2012.
- 438 [33] Paul Jaccard. The distribution of the flora in the alpine zone. *The New Phytologist*, 11(2):37–50,
439 1912.
- 440 [34] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with Gumbel-Softmax.
441 In *International Conference on Learning Representations (ICLR)*, 2017.
- 442 [35] James Jordon, Jinsung Yoon, and Mihaela van der Schaar. KnockoffGAN: Generating knockoffs
443 for feature selection using generative adversarial networks. In *International Conference on*
444 *Learning Representations (ICLR)*, 2018.
- 445 [36] Päivi Järvensivu, Taija Heinosalu, Janne Hakkarainen, Pauliina Kronqvist, Niina Saarinen, and
446 Matti Poutanen. HSD17B1 expression induces inflammation-aided rupture of mammary gland
447 myoepithelium. *Endocrine-Related Cancer*, 25(4):393 – 406, 2018.
- 448 [37] Tadashi Kato, Atsushi Yamada, Mikiko Ikehata, Yuko Yoshida, Kiyohito Sasa, Naoko Morimura,
449 Akiko Sakashita, Takehiko Iijima, Daichi Chikazu, Hiroaki Ogata, and Ryutaro Kamijo. FGF-2
450 suppresses expression of nephronectin via JNK and PI3K pathways. *FEBS Open Bio*, 8(5):836–
451 842, 2018.
- 452 [38] Pora Kim, Feixiong Cheng, Junfei Zhao, and Zhongming Zhao. ccmGDB: a database for cancer
453 cell metabolism genes. *Nucleic Acids Research*, 44(D1):D959–D968, 2015.
- 454 [39] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
455 *arXiv:1412.6980*, 2014.
- 456 [40] Kenji Kira and Larry A. Rendell. A practical approach to feature selection. In *Machine Learning*
457 *Proceedings*, pages 249–256. 1992.
- 458 [41] Theo A. Knijnenburg, Gunnar W. Klau, Francesco Iorio, Mathew J. Garnett, Ultan McDermott,
459 Ilya Shmulevich, and Lodewyk F. A. Wessels. Logic models to predict continuous outputs
460 based on binary inputs with an application to personalized cancer therapy. *Scientific Reports*,
461 6(1):36812, 2016.
- 462 [42] Ron Kohavi and George H. John. Wrappers for feature subset selection. *Artificial Intelligence*,
463 97(1):273–324, 1997.
- 464 [43] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
465 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 466 [44] Ismael Lemhadri, Feng Ruan, and Rob Tibshirani. LassoNet: Neural networks with feature
467 sparsity. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2021.
- 468 [45] Chad Leugers, Ju Yong Koh, Willis Hong, and Gloria Lee. Tau in MAPK activation. *Frontiers*
469 *in Neurology*, 4, 2013.
- 470 [46] Xinghua Li, Weijiang Liang, Junling Liu, Chuyong Lin, Shu Wu, Libing Song, and Zhongyu
471 Yuan. Transducin (β)-like 1 X-linked receptor 1 promotes proliferation and tumorigenicity in
472 human breast cancer via activation of beta-catenin signaling. *Breast Cancer Research*, 16(5):465,
473 2014.

- 474 [47] Yifeng Li Li, Chih-Yu Chen, and Wyeth W. Wasserman. Deep feature selection: theory and
475 application to identify enhancers and promoters. *Journal of Computational Biology*, 23(5):322–
476 336, 2016.
- 477 [48] Faming Liang, Qizhai Li, and Lei Zhou. Bayesian neural networks for selection of drug sensitive
478 genes. *Journal of the American Statistical Association*, 113(523):955–972, 2018.
- 479 [49] Ofir Lindenbaum, Uri Shaham, Erez Peterfreund, Jonathan Svirsky, Nicolas Casey, and Yuval
480 Kluger. Differentiable unsupervised feature selection based on a gated laplacian. In *Advances
481 in Neural Information Processing Systems (NeurIPS)*, 2021.
- 482 [50] Huan Liu and Rudy Setiono. A probabilistic approach to feature selection - a filter solution. In
483 *International Conference on Machine Learning (ICML)*, 1996.
- 484 [51] Ying Liu and Cheng Zheng. Deep latent variable models for generating knockoffs. *Stat*,
485 8(1):e260, 2019.
- 486 [52] Huanyu Lu, Yue Guo, Gaurav Gupta, and Xingsong Tian. Mitogen-activated protein kinase
487 (MAPK): New insights in breast cancer. *Journal of Environmental Pathology, Toxicology and
488 Oncology*, 38(1):51–59, 2019.
- 489 [53] Chris J. Maddison, Andriy Mnih, and Yee Whye Teh. The Concrete distribution: A con-
490 tinuous relaxation of discrete random variables. In *International Conference on Learning
491 Representations (ICLR)*, 2017.
- 492 [54] Ramamurthy Mani, Robert P. St.Onge, John L. Hartman, Guri Giaever, and Frederick P. Roth.
493 Defining genetic interaction. *Proceedings of the National Academy of Sciences*, 105(9):3461–
494 3466, 2008.
- 495 [55] Pulak R. Manna, Ahsen U. Ahmed, Shengping Yang, Madhusudhanan Narasimhan, Joëlle
496 Cohen-Tannoudji, Andrzej T. Slominski, and Kevin Pruitt. Genomic profiling of the steroido-
497 genic acute regulatory protein in breast cancer: In silico assessments and a mechanistic perspec-
498 tive. *Cancers*, 11(5), 2019.
- 499 [56] Kevin McCloskey, Ankur Taly, Federico Monti, Michael P. Brenner, and Lucy J. Colwell. Using
500 attribution to decode binding mechanism in neural network models for chemistry. *Proceedings
501 of the National Academy of Sciences*, 116(24):11624–11629, 2019.
- 502 [57] Sofia Papadimitriou, Andrea Gazzo, Nassim Versbraegen, Charlotte Nachtegael, Jan Aerts,
503 Yves Moreau, Sonia Van Dooren, Ann Nowé, Guillaume Smits, and Tom Lenaerts. Predicting
504 disease-causing variant combinations. *Proceedings of the National Academy of Sciences*,
505 116(24):11878–11887, 2019.
- 506 [58] Ui-Hyun Park, Mi Ran Kang, Eun-Joo Kim, Young-Soo Kwon, Wooyoung Hur, Seung Kew
507 Yoon, Byoung-Joon Song, Jin Hwan Park, Jin-Taek Hwang, Ji-Cheon Jeong, and Soo-Jong Um.
508 ASXL2 promotes proliferation of breast cancer cells by linking ER α to histone methylation.
509 *Oncogene*, 35(28):3742–3752, 2016.
- 510 [59] Hanna M. Peltonen, Annakaisa Haapasalo, Mikko Hiltunen, Vesa Kataja, Veli-Matti Kosma,
511 and Arto Mannermaa. Γ -secretase components as predictors of breast cancer outcome. *PLOS
512 ONE*, 8(11), 2013.
- 513 [60] Bernard Pereira, Suet-Feung Chin, Oscar M. Rueda, Hans-Kristian Moen Vollan, Elena Proven-
514 zano, Helen A. Bardwell, Michelle Pugh, Linda Jones, Roslin Russell, Stephen-John Sammut,
515 Dana W. Y. Tsui, Bin Liu, Sarah-Jane Dawson, Jean Abraham, Helen Northen, John F. Peden,
516 Abhik Mukherjee, Gulisa Turashvili, Andrew R. Green, Steve McKinney, Arusha Oloumi,
517 Sohrab Shah, Nitzan Rosenfeld, Leigh Murphy, David R. Bentley, Ian O. Ellis, Arnie Pu-
518 rushotham, Sarah E. Pinder, Anne-Lise Børresen-Dale, Helena M. Earl, Paul D. Pharoah,
519 Mark T. Ross, Samuel Aparicio, and Carlos Caldas. The somatic mutation profiles of 2,433

- 520 breast cancers refine their genomic and transcriptomic landscapes. *Nature Communications*,
521 7(1):11479, 2016.
- 522 [61] Patrick C. Phillips. Epistasis — the essential role of gene interactions in the structure and
523 evolution of genetic systems. *Nature Reviews Genetics*, 9(11):855–867, 2008.
- 524 [62] Michael W. Pickup, Laura D. Hover, Eleanor R. Polikowsky, Anna Chytil, Agnieszka E. Gorska,
525 Sergey V. Novitskiy, Harold L. Moses, and Philip Owens. BMPR2 loss in fibroblasts promotes
526 mammary carcinoma metastasis via increased inflammation. *Molecular Oncology*, 9(1):179–
527 191, 2015.
- 528 [63] Barbara Maria Piskór, Andrzej Przylipiak, Emilia Dąbrowska, Iwona Sidorkiewicz, Marek
529 Niczyporuk, Maciej Szmikowski, and Sławomir Ławicki. Plasma level of MMP-10 may be a
530 prognostic marker in early stages of breast cancer. *Journal of Clinical Medicine*, 9(12), 2020.
- 531 [64] Franck Rapaport, Emmanuel Barillot, and Jean-Philippe Vert. Classification of arrayCGH data
532 using fused SVM. *Bioinformatics*, 24(13):i375–i382, 2008.
- 533 [65] Yaniv Romano, Matteo Sesia, and Emmanuel Candès. Deep knockoffs. *Journal of the American
534 Statistical Association*, 115(532):1861–1872, 2020.
- 535 [66] Benjamin Sanchez-Lengeling, Jennifer Wei, Brian Lee, Emily Reif, Peter Wang, Wesley Qian,
536 Kevin McCloskey, Lucy Colwell, and Alexander Wiltchko. Evaluating attribution for graph
537 neural networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 538 [67] Fadil Santosa and William W. Symes. Linear inversion of band-limited reflection seismograms.
539 *SIAM Journal on Scientific and Statistical Computing*, 7(4):1307–1330, 1986.
- 540 [68] Razieh Sheikhpour, Mehdi Agha Sarram, Sajjad Gharaghani, and Mohammad Ali Zare Cha-
541 hooki. A survey on semi-supervised feature selection methods. *Pattern Recognition*, 64:141–158,
542 2017.
- 543 [69] Prajwal K. Singha, Srilakshmi Pandeswara, Hui Geng, Rongpei Lan, Manjeri A. Venkatachalam,
544 Albert Dobi, Shiv Srivastava, and Pothana Saikumar. Increased smad3 and reduced smad2
545 levels mediate the functional switch of TGF- β from growth suppressor to growth and metastasis
546 promoter through TMEPAI/PMEPA1 in triple negative breast cancer. *Genes & cancer*, 10(5-
547 6):134, 2019.
- 548 [70] Mukund Sudarshan, Wesley Tansey, and Rajesh Ranganath. Deep direct likelihood knockoffs.
549 In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- 550 [71] Mina Takahashi, Fumio Otsuka, Tomoko Miyoshi, Hiroyuki Otani, Junko Goto, Misuzu Ya-
551 mashita, Toshio Ogura, Hirofumi Makino, and Hiroyoshi Doihara. Bone morphogenetic protein
552 6 (BMP6) and BMP7 inhibit estrogen-induced proliferation of breast cancer cells by suppressing
553 p38 mitogen-activated protein kinase activation. *Journal of Endocrinology*, 199(3):445 – 455,
554 2008.
- 555 [72] Jiliang Tang, Salem Alelyani, and Huan Liu. Feature selection for classification: A review.
556 *Data classification: Algorithms and applications*, page 37, 2014.
- 557 [73] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal
558 Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- 559 [74] Nicholas Turner, Alex Pearson, Rachel Sharpe, Maryou Lambros, Felipe Geyer, Maria A.
560 Lopez-Garcia, Rachael Natrajan, Caterina Marchio, Elizabeth Iorns, Alan Mackay, Cheryl
561 Gillett, Anita Grigoriadis, Andrew Tutt, Jorge S. Reis-Filho, and Alan Ashworth. FGFR1
562 amplification drives endocrine therapy resistance and is a therapeutic target in breast cancer.
563 *Cancer Research*, 70(5):2085–2094, 2010.

- 564 [75] Dongfeng Wang, Jian Li, Fengling Cai, Zhi Xu, Li Li, Huanfeng Zhu, Wei Liu, Qingyu Xu, Jian
565 Cao, Jingfeng Sun, and Jinhai Tang. Overexpression of MAPT-AS1 is associated with better
566 patient survival in breast cancer. *Biochemistry and Cell Biology*, 97(2):158–164, 2019.
- 567 [76] Dongsheng Wang, Chenglong Zhao, Liangliang Gao, Yao Wang, Xin Gao, Liang Tang, Kun
568 Zhang, Zhenxi Li, Jing Han, and Jianru Xiao. NPNT promotes early-stage bone metastases in
569 breast cancer by regulation of the osteogenic niche. *Journal of Bone Oncology*, 13:91–96, 2018.
- 570 [77] Chang-Yuan Wei, Qi-Xing Tan, Xiao Zhu, Qing-Hong Qin, Fei-Bai Zhu, Qin-Guo Mo, and
571 Wei-Ping Yang. Expression of CDKN1A/p21 and TGFBR2 in breast cancer and their prognostic
572 significance. *International Journal of Clinical and Experimental Pathology*, 8(11):14619, 2015.
- 573 [78] Michael K. Wendt, Molly A. Taylor, Barbara J. Schiemann, Khalid Sossey-Alaoui, and
574 William P. Schiemann. Fibroblast growth factor receptor splice variants are stable mark-
575 ers of oncogenic transforming growth factor β 1 signaling in metastatic breast cancers. *Breast
576 Cancer Research*, 16(2):R24, 2014.
- 577 [79] Tim Wilson, Tim Holt, and Trisha Greenhalgh. Complexity and clinical care. *BMJ*,
578 323(7314):685–688, 2001.
- 579 [80] Xinxing Wu and Qiang Cheng. Algorithmic stability and generalization of an unsupervised
580 feature selection algorithm. In *Advances in Neural Information Processing Systems (NeurIPS)*,
581 2021.
- 582 [81] Yutaro Yamada, Ofir Lindenbaum, Sahand Negahban, and Yuval Kluger. Feature selection
583 using stochastic gates. In *International Conference on Machine Learning (ICML)*, 2020.
- 584 [82] Ming Yuan and Yi Lin. Model selection and estimation in regression with grouped variables.
585 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(1):49–67, 2006.
- 586 [83] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov,
587 and Alexander Smola. Deep sets. In *Advances in Neural Information Processing Systems
588 (NeurIPS)*, 2017.
- 589 [84] Lingmin Zeng and Jun Xie. Group variable selection via SCAD-L2. *Statistics*, 48(1):49–66,
590 2014.
- 591 [85] Xin-Yue Zhang, Hsun-Ming Chang, Elizabeth L. Taylor, Rui-Zhi Liu, and Peter C. K. Leung.
592 BMP6 downregulates GDNF expression through SMAD1/5 and ERK1/2 signaling pathways in
593 human granulosa-lutein cells. *Endocrinology*, 159(8):2926–2938, 2018.
- 594 [86] Yong-ping Zhang, Wen-ting Na, Xiao-qiang Dai, Ruo-fei Li, Jian-xiong Wang, Ting Gao,
595 Wei-bo Zhang, and Cheng Xiang. Over-expression of SRD5A3 and its prognostic significance
596 in breast cancer. *World Journal of Surgical Oncology*, 19(1):260, 2021.
- 597 [87] Jidong Zhao, Ke Lu, and Xiaofei He. Locality sensitive semi-supervised feature selection.
598 *Neurocomputing*, 71:1842—1849, 2008.
- 599 [88] Ting Zhong, Feifei Xu, Jinhui Xu, Liang Liu, and Yun Chen. Aldo-keto reductase 1C3
600 (AKR1C3) is associated with the doxorubicin resistance in human breast cancer via PTEN loss.
601 *Biomedicine & Pharmacotherapy*, 69:317–325, 2015.
- 602 [89] Nengfeng Zhou and Ji Zhu. Group variable selection via a hierarchical lasso and its oracle
603 property. *arXiv preprint arXiv:1006.2871*, 2010.
- 604 [90] Guangyu Zhu and Tingting Zhao. Deep-gKnock: Nonlinear group-feature selection with deep
605 neural networks. *Neural Networks*, 135:139–147, 2021.

606 **Checklist**

- 607 1. For all authors...
- 608 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
609 contributions and scope? [\[Yes\]](#)
- 610 (b) Did you describe the limitations of your work? [\[Yes\]](#) See Conclusion for details.
- 611 (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) In the
612 Conclusion, we caution that as with any feature selection method, discovered features
613 must be verified or evaluated by domain experts. This verification or evaluation might
614 be costly, and should the method perform poorly, could result in wasted resources. In
615 addition, without additional oversight (primarily in dataset construction but also when
616 validating features), features that contain bias could remain and be identified by feature
617 selection algorithms.
- 618 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
619 them? [\[Yes\]](#)
- 620 2. If you are including theoretical results...
- 621 (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
- 622 (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
- 623 3. If you ran experiments...
- 624 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
625 mental results (either in the supplemental material or as a URL)? [\[Yes\]](#) We provide all
626 code needed to reproduce all results in the supplemental material.
- 627 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
628 were chosen)? [\[Yes\]](#) Hyperparameters for each experiment are provided in Table 4.
629 Architecture details are provided in Appendix C and further experimental details are
630 provided in Appendix D.
- 631 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
632 iments multiple times)? [\[Yes\]](#) All experiments are repeated 10 times and results are
633 reported along with standard deviations.
- 634 (d) Did you include the total amount of compute and the type of resources used (e.g.,
635 type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) All experiments can be run
636 easily on a commercially-available laptop. We provide further details of the compute
637 resources used in Appendix D.
- 638 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 639 (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#) We used several
640 existing methods and datasets (see Experiments). All benchmark methods and datasets
641 are clearly cited.
- 642 (b) Did you mention the license of the assets? [\[Yes\]](#) Licenses of assets (benchmark methods
643 and datasets) is provided in Appendix D and F.
- 644 (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#)
645 The code to run our experiments is included in the supplementary material. The code
646 will also be released publicly after the review period.
- 647 (d) Did you discuss whether and how consent was obtained from people whose data you’re
648 using/curating? [\[N/A\]](#)
- 649 (e) Did you discuss whether the data you are using/curating contains personally identifiable
650 information or offensive content? [\[N/A\]](#)
- 651 5. If you used crowdsourcing or conducted research with human subjects...
- 652 (a) Did you include the full text of instructions given to participants and screenshots, if
653 applicable? [\[N/A\]](#)

- 654 (b) Did you describe any potential participant risks, with links to Institutional Review
655 Board (IRB) approvals, if applicable? [N/A]
- 656 (c) Did you include the estimated hourly wage paid to participants and the total amount
657 spent on participant compensation? [N/A]