Variable selection in transcriptomics data using knockoffs in a classification framework

Julie Cartier Supervisors : Florian Massip (CBIO)

> **Chloé-Agathe Azencott** (CBIO) **Adeline Fermanian** (Califrais)

Workshop on statistical inference for complex data: 12/09/2024

PhD: 14/11/2022 -









Large-scale biological data analysis comes with challenges



		gei	nes		
	-	ENSG0000131697	ENSG00000116661	ENSG00000157330	0
les	S1_84_NS_2				14.14.14.1
du	S1_85_NS_2		*		
Sa	S1_86_NS_2				
		ae	ne expression lev	el	

Large-scale biological data analysis comes with challenges





1

Large-scale biological data analysis comes with challenges

Biomarker identification (Variable selection)

- → **High-dimension**: Vast amount of variables (20 000 genes), small cohorts (hundreds of patients)
- → Correlations: Interactions between different functional units of the biological system



05223 9/24/20 (c) Kanehisa Laboratorie

Perform variable selection with the knockoff procedure⁽¹⁾

The Knockoff procedure⁽¹⁾:

- Designed for correlated settings
- Handle high dimension (Model-X knockoff⁽²⁾)

🔽 Can be used with no prior knowledge about the structure of the data

- Control the expected number of false discoveries
- □ Already used with biological data (GWAS)⁽³⁾

Study the applicability of the knockoff procedure to transcriptomic data in classification

(1): Barber and Candes : Controlling the false discovery rate via knockoffs (2015)

(2): Candes et al - Panning for Gold: Model-X Knockoffs for High-dimensional Controlled Variable Selection (2017)

^{(3):} Sesia et al - False discovery rate control in genome wide association studies with population structure

The knockoff variable selection procedure

Knockoff = False copy



X

$$X_{j} \qquad \qquad \widetilde{X}_{j} \qquad \qquad \widetilde{X}_{j}$$

$$X = \begin{pmatrix} X_{1,1} \cdots X_{1,j} \cdots X_{1,p} \\ X_{2,1} \cdots X_{2,j} \cdots X_{2,p} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ X_{n,1} \cdots & X_{n,j} \cdots & X_{n,p} \end{pmatrix} \qquad \widetilde{X} = \begin{pmatrix} \widetilde{X}_{1,1} \cdots & \widetilde{X}_{1,j} \cdots & \widetilde{X}_{1,p} \\ \widetilde{X}_{2,1} \cdots & \widetilde{X}_{2,j} & \cdots & \widetilde{X}_{2,p} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \widetilde{X}_{n,1} \cdots & \widetilde{X}_{n,j} & \cdots & \widetilde{X}_{n,p} \end{pmatrix}$$

$$\widetilde{X} \perp Y \mid X$$



The knockoff variable selection procedure

Knockoff = False copy



X

$$X = \begin{pmatrix} X_{1,1} \cdots X_{1,j} \cdots X_{1,p} \\ X_{2,1} \cdots X_{2,j} \cdots X_{2,p} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ X_{n,1} \cdots & X_{n,j} \cdots & X_{n,p} \end{pmatrix} \qquad \tilde{X} = \begin{pmatrix} \tilde{X}_{1,1} \cdots \tilde{X}_{1,j} \cdots \tilde{X}_{1,p} \\ \tilde{X}_{2,1} \cdots \tilde{X}_{2,j} \cdots \tilde{X}_{2,p} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ \tilde{X}_{n,1} \cdots \tilde{X}_{n,j} \cdots \tilde{X}_{n,p} \end{pmatrix} \qquad \tilde{X} = \begin{pmatrix} \tilde{X}_{1,1} \cdots \tilde{X}_{1,p} \cdots \tilde{X}_{1,p} \\ \tilde{X}_{1,1} \cdots \tilde{X}_{n,j} \cdots \tilde{X}_{n,p} \end{pmatrix} \\ \tilde{X} \perp Y | X \qquad \tilde{X} \perp Y | X \qquad \tilde{X}_{1,j} \cdots \tilde{X}_{1,p} \\ \vdots & \ddots & \vdots & \vdots & \ddots \\ X_{n,1} \cdots X_{n,j} \cdots X_{n,p} \quad \tilde{X}_{n,1} \cdots \begin{pmatrix} \tilde{X}_{1,j} \cdots \tilde{X}_{1,p} \\ \vdots & \ddots & \vdots \\ X_{n,j} \cdots & X_{n,p} \quad \tilde{X}_{n,1} \cdots \begin{pmatrix} \tilde{X}_{1,j} \cdots \tilde{X}_{1,p} \\ \vdots & \ddots & \vdots \\ \tilde{X}_{n,j} \cdots & \tilde{X}_{n,p} \end{pmatrix}$$
Effect J
Effect J
Effect J

Effect $\mathbf{J} > \text{Effect } \mathbf{\tilde{J}}$:

the jth variable is a

relevant feature

Effect $J \sim$ Effect \tilde{J} : the jth variable has "spurious" effect

Variable selection with Model-X knockoffs⁽²⁾ procedure



- \rightarrow Adjustable to the context
- → There are different ways to use the KO procedure

(2) : Candes et al - Panning for Gold: Model-X Knockoffs for High-dimensional Controlled Variable Selection

Study the applicability of the knockoff procedure to transcriptomic data in classification

-> Performance of the method

How to build relevant knockoff features ?

Which type of statistic performs well?

Stability

Study the reliability of selected sets of features



Data : Nasal transcriptomic data used to assess lung cancer risk



Good classification performances⁽⁴⁾:

- AUC-PR : 0.82
- AUC-ROC : 0.83

• Features are not **stably** selected

Study the applicability of the knockoff procedure to transcriptomic data in classification

-> Performance of the method

How to build relevant knockoff features ?

Which type of statistic performs well?

Compute the statistics



Compute the statistics

- W_j is a function of Y and $[X, \widetilde{X}], \forall j \in \{1, ..., p\}$
- Large values for features that must be selected
- Some other specific properties



Effect $\mathbf{J} > \text{Effect } \mathbf{\tilde{J}}$ or Effect $\mathbf{J} \sim \text{Effect } \mathbf{\tilde{J}}$?

Test statistics benchmark

$$[X, \widetilde{X}] = \begin{pmatrix} X_{1,1} & \cdots & X_{1,j} \\ \vdots & \ddots & \vdots \\ X_{n,1} & \cdots & X_{n,j} \end{pmatrix} \cdots X_{1,p} \quad \widetilde{X}_{1,1} & \cdots & \widetilde{X}_{1,j} \\ \vdots & \ddots & \vdots \\ X_{n,j} & \cdots & X_{n,p} \quad \widetilde{X}_{n,1} & \cdots & \widetilde{X}_{n,j} \end{pmatrix} \cdots \qquad \widetilde{X}_{n,p} \end{pmatrix}$$

Variable Importance (VI) - based methods:

• Fit a ML model to
$$[X, \widetilde{X}], Y$$
.
• Get $VI(j)$ and $VI(j+p), \forall j \in \{1, ..., p\}$
• $W_j = \underline{VI}(j) - \underline{VI}(j+p)$

Test statistics benchmark

$$[X, \widetilde{X}] = \begin{pmatrix} X_{1,1} & \cdots & X_{1,j} \\ \vdots & \ddots & \vdots \\ X_{n,1} & \cdots & X_{n,j} \end{pmatrix} \cdots \begin{pmatrix} \widetilde{X}_{1,1} & \cdots & \widetilde{X}_{1,j} \\ \vdots & \ddots & \vdots \\ X_{n,j} & \cdots & X_{n,p} \end{pmatrix} \widetilde{X}_{n,1} \cdots \begin{pmatrix} \widetilde{X}_{1,j} & \cdots & \widetilde{X}_{1,p} \\ \vdots & \ddots & \vdots \\ \widetilde{X}_{n,j} & \cdots & \widetilde{X}_{n,p} \end{pmatrix}$$

Variable Importance (VI) - based methods:

- Fit a ML model to $[X, \widetilde{X}], Y$.
- Get VI(j) and $VI(j\!+\!p), \forall j \in \{1,...,p\}$

•
$$W_j = \underline{VI}(j) - \underline{VI}(j+p)$$

- LASSO
- Elastic-net (EN)
- Random Forest (RF)
- Boosted Tree
- Deep learning

Data simulation

Simulated outcomes from a given matrix $X \in \mathbb{R}^{369 \times 749}$:

Classification:



-

$$Y = \mathbb{1}_{\left(\frac{1}{1+e^{-X_n\beta}} > 0.5\right)} = \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

Linear setting:

$$X_n\beta = X_1\beta_1 + X_2\beta_2 + \dots + X_p\beta_p$$

Test statistics benchmark

$$[X, \widetilde{X}] = \begin{pmatrix} X_{1,1} & \cdots & X_{1,j} \\ \vdots & \ddots & \vdots \\ X_{n,1} & \cdots & X_{n,j} \end{pmatrix} \overset{\widetilde{X}_{1,1}}{\cdots} & \cdots & \widetilde{X}_{1,j} \\ \vdots & \ddots & \vdots & \vdots & \ddots \\ X_{n,j} & \cdots & X_{n,p} & \widetilde{X}_{n,1} & \cdots & \widetilde{X}_{n,j} \end{pmatrix} \overset{\widetilde{X}_{1,j}}{\cdots} \overset{\widetilde{X}_{1,j}}{\cdots} \overset{\widetilde{X}_{1,p}}{\widetilde{X}_{n,j}}$$

Variable Importance (VI) - based methods:

- Fit a ML model to $[X, \widetilde{X}], Y$.
- Get VI(j) and $VI(j+p), \forall j \in \{1, ..., p\}$
- $W_j = \underline{VI}(j) \underline{VI}(j+p)$

LASSO

- Elastic-net (EN)
- Random Forest (RF)
- Boosted Tree
- Deep learning

LASSO Coefficient-Difference⁽²⁾ (LCD)

+ Other methods

LCD based knockoff improves variable selection



→ The LCD knockoff is more powerful than the lasso penalized logistic regression for reasonable values of FDR/FDP.

 The LCD based knockoff is more powerful than the lasso penalized logistic regression for small numbers of false discoveries Study the applicability of the knockoff procedure to transcriptomic data in classification

Stability

Study the reliability of selected sets of features

(+ assess method's reliability with no ground truth)

How to get selection frequency



Lasso systematically selects both important and noisy features



X	real matrix

Linear simulation Y

regression

The KO framework removes false positive genes from the set of most selected genes...



The KO framework removes false positive genes from the set of most selected genes...



The KO framework removes false positive genes from the set of most selected genes...



Additional source of instability with the knockoff framework



Knockoff aggregation

Run the knockoff procedure multiple times (in parallel):



Knockoff aggregation

Run the knockoff procedure multiple times (in parallel):



pierre.neuvial@math.univ-toulouse.fr

Knockoffs		
Alexandre Blain	Bertrand Thirion	
Université Paris-Saclay	CEA	
alexandre.blain@inria.fr	bertrand.thirion@inria.fr	
Olivier Grisel	Pierre Neuvial	
INRIA	Institut de Mathématiques de Toulous	
olivier.grisel@inria.fr	Université de Toulouse	

Evaluate aggregation knockoff efficiency



The aggregation of knockoffs tends to remove alternatively selected genes



📃 False discoveries 📕 True genes

The aggregation of knockoffs tends to remove alternatively selected genes



False discoveries **F**rue genes

- The LCD based knockoff is more powerful than the lasso penalized logistic regression for small numbers of false discoveries
- The KO procedure does not improve the stability of features selection in our setting
 - > Aggregating knockoffs eliminates some of the effect of knockoff stochasticity

- The LCD based knockoff is more powerful than the lasso penalized logistic regression for small numbers of false discoveries
- The KO procedure does not improve the stability of features selection in our setting
 - > Aggregating knockoffs eliminates some of the effect of knockoff stochasticity

→ Identified limitations of the knockoff framework

- > The knockoff procedure fails with interacting variables
- \succ The power decreases with the number of non null features

- The LCD based knockoff is more powerful than the lasso penalized logistic regression for small numbers of false discoveries
- The KO procedure does not improve the stability of features selection in our setting
 - > Aggregating knockoffs eliminates some of the effect of knockoff stochasticity
- Identified limitations of the knockoff procedure
 - Powerless in the interaction setting
 - > Decrease in power as the number of non-null features increases
- The knockoff procedure does not select any genes with real outcomes

Group knockoffs

Why?

- Improve the interaction setting
- Perform selection at a higher functional level (more robust)
- Integrate more features

Group knockoffs

Why?

- Improve the interaction setting
- Perform selection at a higher functional level (more robust)
- Integrate more features

- How to make groups ?
 - ➢ e.g.: transcription factors, biological pathways,
- How to build knockoffs ?
 - Either use previous methods or use group's structure instead of correlations between features
- At which level does the statistic do the selection (feature or group) ?
 - > group level

How?

Acknowledgements

- Florian MASSIP
- Adeline FERMANIAN
- Chloé-Agathe AZENCOTT
- Johanna LAGOAS

Thank you for your attention !

CBIO members









749 risk genes



-115 shared DE genes + same gene expression dynamic/behavior

(3) : de Biase et al. - Smoking-dependent expression alterations in the general population reveal immune impairment linked to germline variation and lung cancer risk

749 risk genes

genes altered by smoke in the clinic 513 genes affected by smoke in the healthy group volunteer group 351 genes altered by smoke only in the clinic group

-115 shared DE genes + same gene expression dynamic/behavior

749 risk genes

(3) : de Biase et al. - Smoking-dependent expression alterations in the general population reveal immune impairment linked to germline variation and lung cancer risk

Generation of the knockoff (KO) : LSCIP algorithm



- Does not use the covariance matrix
- no assumption on features distribution

LSCIP pseudo-algorithm :

▶ j = 1, for all $j \in \{1, ..., p\}$:

(4): Blain et al - False discovery proportion control for aggregated knockoffs

- ≻ Fit a lasso model on (X_{-j}) with X_j as outcomes.
- ➤ Compute residuals $\epsilon = (\epsilon_1, ..., \epsilon_p), \epsilon = X_j \hat{X}_j$ where \hat{X}_j is the predicted value of X_j with the regression model
- > Permute the residuals vectors randomly. Let's denote ρ_p the permutation of $\{1, ..., p\}$. compute $\tilde{X}_j = \hat{X}_j + \epsilon_{\rho_p(j)}$
- \blacktriangleright Return \widetilde{X}

computationally more expensive

Variable selection with Model-X knockoffs⁽²⁾ procedure



Variable selection with Model-X knockoffs⁽²⁾ procedure



🖷 CI 📥 LSCIP 💻 MVR 📥 SDP



LSCIP + LCD/MLR based knockoff are the most powerful



Power and FDP obtained with the knockoff selection procedure







Stability selection with simulated data



KOPI is empirically at least as powerful as vanilla knockoff



Aggregation simulations



Impact of the aggregation scheme on the selection frequency



The KO framework does not improve stability (real data with subsampling)







Alpha effect



Study the applicability of the knockoff procedure to transcriptomic data in classification

→ Identified limitations of the knockoff framework

The knockoff procedure fails with interacting variables



The knockoff procedure fails with interacting variables



Power obtained with the knockoff selection procedure

The power decreases with the number of non null features





Power and FDP obtained with the LSCIP + LCD knockoff selection procedure for different numbers of non-null features (k)

Performances also depend on the features used to simulate Y...

- Setting (linear)
- KO framework (LSCIP + LCD)
- Features matrix





Power and FDP obtained with the LSCIP + LCD knockoff selection procedure

Application to the real data



Selection frequency

Application to the real data



Selection frequency

Application to the real data



Application to the real data vs simulated data





(3) : de Biase et al. - Smoking-dependent expression alterations in the general population reveal immune impairment linked to germline variation and lung cancer risk