

Identifying relevant covariates in RNA-seq analysis by pseudo-variable augmentation

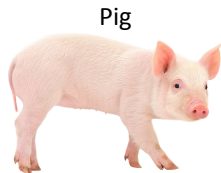
Yet Nguyen, Ph.D.¹ and Dan Nettleton, Ph.D.²

2024-11-02

¹Department of Mathematics and Statistics, Old Dominion University

²Department of Statistics, Iowa State University

Experimental Design



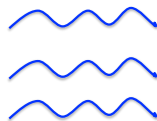
Diet
Line
RFI

Blood Sample



Basophil
Eosinophil
Lymphocyte
Monocyte
Neutrophil

RNA sample



Block
Order
RIN before globin depletion (GD)
RIN after GD
RNA Concentration before GD
RNA Concentration after GD

Prototypical RNA-seq Dataset

	Treatment 1				Treatment 2			
	u_{11}	u_{12}	\dots	u_{1n_1}	u_{21}	u_{22}	\dots	u_{2n_2}
\mathbf{x}_1	0	0	\dots	0	1	1	\dots	1
\mathbf{x}_2	0.5	0.95	\dots	-1.42	-0.45	.89	\dots	1.2
\vdots	\vdots	\vdots		\vdots	\vdots	\vdots		\vdots
\mathbf{x}_k	0	1	\dots	0	1	0	\dots	0
gene 1	10	13	\dots	2017	31	975	\dots	3289
gene 2	0	2	\dots	1	0	0	\dots	1
gene 3	1	3	\dots	0	0	0	\dots	0
\vdots	\vdots	\vdots		\vdots	\vdots	\vdots		\vdots
gene G	17301	2464	\dots	7345	3214	534	\dots	934

How to Handle the Available Covariates?

- Including all available covariates (Full)
- Excluding all available covariates (OnlyLine)
- Backward selection that maximizes the number of DE genes with respect to the main factor of interest (BS15, Nguyen et al. 2015)

Our Proposed Method (Nguyen and Nettleton 2024+)

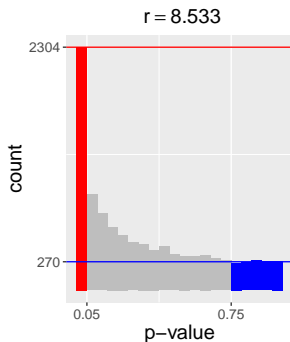
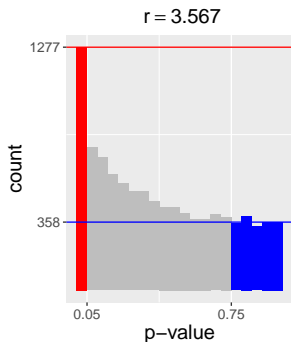
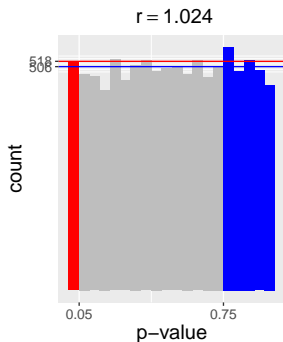
- Using `limma-voom` (Law et al. (2014)) for differential expression analysis to obtain vectors of p -values of tests for significance of regression coefficients w.r.t each of the covariates
- Selecting the most relevant covariates by a **backward selection** strategy intending to control the **false selection rate (FSR)** using **pseudo-variables** (Wu et al. (2007), 'Controlling Variable Selection by the Addition of Pseudovariables', JASA)
- Wu et al. (2007) method published for one response variable
- We extend Wu et al. (2007)'s method to thousands of response variables

Measure of Covariate Relevance

Definition

With $\mathbb{1}$ representing an indicator function, a relevance measure for covariate j is defined as

$$r(\mathbf{p}_j) = \frac{\sum_{g=1}^G \mathbb{1}(p_{gj} \leq 0.05)}{\max\{\sum_{g=1}^G \mathbb{1}(p_{gj} \geq 0.75)/5, 1\}}. \quad (1)$$



Backward Selection to Control FSR

- Run backward selection procedure using $r(\cdot)$ on k_T covariates of \mathbf{X}
- Let $BS(\mathbf{X}, \lambda)$ denote the subset of \mathbf{X} selected by this backward selection, i.e., the largest subset of \mathbf{X} for which each variable has r -value at least λ
- Define $S(\lambda) = \text{Card}\{BS(\mathbf{X}, \lambda)\}$. Then $S(\lambda) = R(\lambda) + I(\lambda)$, where $R(\lambda)$ and $I(\lambda)$ denote the number of selected relevant and irrelevant covariates, respectively
- False selection rate (FSR) is calculated as $\alpha(\lambda) = \frac{E(I(\lambda))}{E(S(\lambda)+1)}$
- Calculate the tuning parameter λ_* to control FSR at level α_0

$$\lambda_* = \inf\{\lambda : \alpha(\lambda) \leq \alpha_0\}.$$

Estimating FSR I

Generate B sets of k_P pseudo-variables \mathbf{Z}_b

Define $\alpha_P(\lambda) = \frac{E(I_{P,b}^*(\lambda))}{E(1+S_{P,b}(\lambda))}$ where

- $R_{P,b}(\lambda)$: number of truly relevant covariates selected from \mathbf{X}, \mathbf{Z}_p
- $I_{P,b}(\lambda)$: number of truly irrelevant covariates selected from \mathbf{X}, \mathbf{Z}_p
- $I_{P,b}^*(\lambda)$: number of pseudo-covariates selected from \mathbf{X}, \mathbf{Z}_p
- $S_{P,b} = R_{P,b}(\lambda) + I_{P,b}(\lambda) + I_{P,b}^*(\lambda)$

Estimating FSR II

Assumptions

(A1) $E(I(\lambda)) = E(I_{P,b}(\lambda)) = k_U E(I_{P,b}^*(\lambda))/k_P$, where k_U is the unknown number of truly irrelevant covariates

(A2) $E(R_{P,b}(\lambda)) = E(R(\lambda))$

- (A1) & (A2) imply: $\alpha_P(\lambda) = \frac{k_P \alpha(\lambda)}{k_P \alpha(\lambda) + k_U}$
- Let $\bar{I}_P^*(\lambda) = B^{-1} \sum_{b=1}^B I_{P,b}^*(\lambda)$, $\bar{S}_P(\lambda) = B^{-1} \sum_{b=1}^B S_{P,b}(\lambda)$
- Estimate $\alpha_P(\lambda)$ by $\hat{\alpha}_P(\lambda) = \frac{\bar{I}_P^*(\lambda)}{1 + \bar{S}_P(\lambda)}$
- If k_U is known, estimate $\alpha(\lambda)$ by solving

$$\hat{\alpha}_P(\lambda) = \frac{k_P \alpha(\lambda)}{k_P \alpha(\lambda) + k_U}$$

Generating Pseudo-covariates $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_{k_p})$

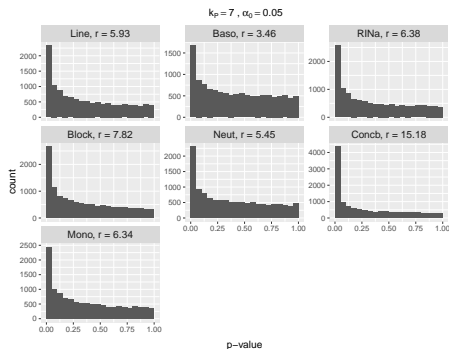
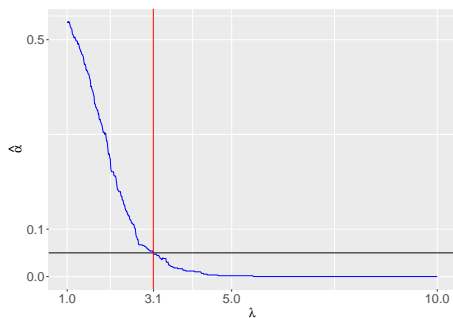
- Option 1 (WN): $\mathbf{z}_1, \dots, \mathbf{z}_{k_p}$ i.i.d. $\sim \mathbf{N}(\mathbf{0}, \mathbf{1})$
- Option 2 (RX): The n rows of \mathbf{Z} are obtained by randomly permuting the rows and the columns of \mathbf{X}
- Options 3 & 4 (OWN & ORX): $(\mathbf{I} - \mathbf{H}_\mathbf{X})\mathbf{Z}$, where $\mathbf{H}_\mathbf{X} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, where \mathbf{Z} is generated either by option 1 or 2, respectively

RFI RNA-seq Data Analysis

Table 1: Covariates removed from the full model and their r values at each iteration of the backward selection algorithm applied to the RFI RNA-seq dataset.

Iteration	1	2	3	4	5	6	7	8	9	10	11	12	13
Covariate	RINb	Eosi	Order	Conca	Diet	RFI	Lymp	Baso	RINa	Block	Neut	Concb	Mono
r	0.26	0.49	0.62	0.65	0.53	2.07	2.87	3.46	6.3	7.71	7.85	9.42	11.45

$k_p = 7, \alpha_0 = 0.05$



Simulation Study - Setting

Table 2: Six simulation scenarios corresponding to six sets of truly relevant covariates.

Number of relevant covariates k_R	Relevant covariates
0	
1	Mono
2	Concb, Mono
6	Baso, RINa, Block, Neut, Concb, Mono
7	Lymp, Baso, RINa, Block, Neut, Concb, Mono
8	RFI, Lymp, Baso, RINa, Block, Neut, Concb, Mono

- Number of genes: 2000
- Number of replications: 100

Simulation Study - FSR Results

FSR Variable Selection Method Results

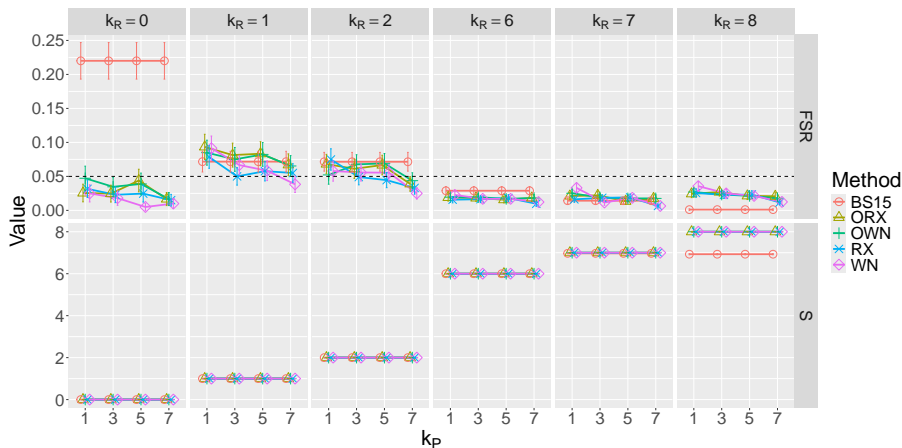


Figure 1: The figure displays the variable selection performance of four variants of the proposed method and BS15.

Simulation Study - Differential Expression Analysis Results

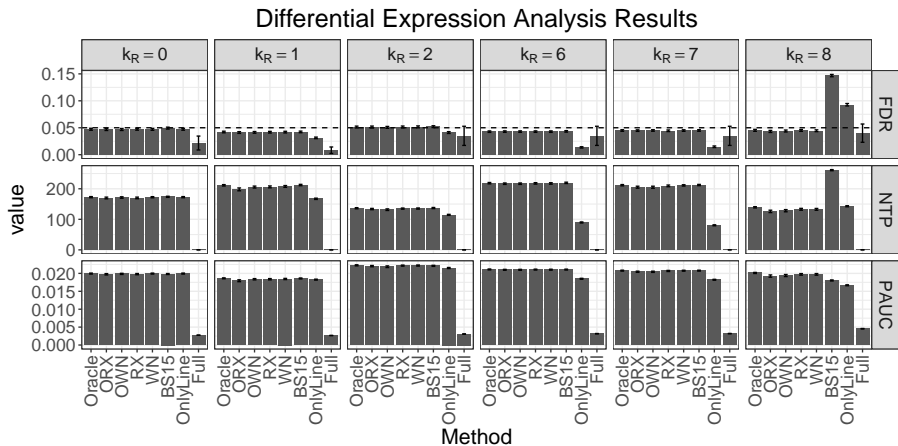


Figure 2: The figure presents the performance of differential expression analysis of the twelve methods.

Conclusion

- The proposed covariate selection method control FDR well
- The selected model has good performance in identifying DE genes in terms of
 - FDR control
 - Ability to distinguish EE genes and DE genes
- The proposed method is available at github.com/ntyet/csrnaseq
- Contact: Yet Nguyen, ynguyen@odu.edu

Thank you!

References I

- Law, C. W., Chen, Y., Shi, W., and Smyth, G. K. (2014), “Voom: Precision weights unlock linear model analysis tools for RNA-seq read counts,” *Genome Biology*, 15, R29. <https://doi.org/10.1186/gb-2014-15-2-r29>.
- Nguyen, Y., and Nettleton, D. (2024+), “Identifying relevant covariates in RNA-seq analysis by pseudo-variable augmentation,” *Journal of Agricultural, Biological and Environmental Statistics*, accepted.
- Nguyen, Y., Nettleton, D., Liu, H., and Tuggle, C. K. (2015), “Detecting differentially expressed genes with RNA-seq data using backward selection to account for the effects of relevant covariates,” *Journal of Agricultural, Biological, and Environmental Statistics*, 20, 577–597. <https://doi.org/10.1007/s13253-015-0226-1>.
- Wu, Y., Boos, D. D., and Stefanski, L. A. (2007), “Controlling variable selection by the addition of pseudovariables,” *Journal of the American Statistical Association*, Taylor & Francis, 102, 235–243. <https://doi.org/10.1198/016214506000000843>.