

Detecting Differentially Expressed Genes with RNA-seq Data Using Backward Selection to Account for the Effects of Relevant Covariates



Yet Nguyen and Dan Nettleton, Department of Statistics, Iowa State University
Haibo Liu and Chris Tuggle, Department of Animal Science, Iowa State University

Introduction

RNA-seq datasets often contain several covariates in addition to the factor of primary scientific interest. Either ignoring relevant covariates or accounting for the effects of irrelevant covariates can result in low power for identifying DE genes. To address the challenge of identifying DE genes with RNA-seq datasets that include covariates, we propose a backward selection algorithm for selecting a subset of covariates whose effects are estimated and adjusted for when testing for differential expression.

Generalized Linear Model for RNA-seq Read Count Data

- Let y_{gi} be the read count for gene g from unit i ($g = 1, \dots, m; i = 1, \dots, n$). $y_{gi} \sim \text{NB}(\mu_{gi}, \omega_g)$, i.e., $\text{Var}(y_{gi}) = \mu_{gi} + \omega_g \mu_{gi}^2$.
- Let $\mathbf{x}_i = (\mathbf{x}'_{i1}, \dots, \mathbf{x}'_{ik})'$ denote a vector of known explanatory variable values for the i th unit.
- Letting \mathcal{S} represent a subset of $\{1, \dots, k\}$ that contains 1, we consider log-linear models of the form
$$\log(\mu_{gi}) = o_i + \beta_{g0|\mathcal{S}} + \sum_{j \in \mathcal{S}} \mathbf{x}'_{ij} \beta_{gj|\mathcal{S}} \quad (1)$$
- Using **QuasiSeq** R package (Lund et al (2012)), we test $H_{0g1}^{\mathcal{S}^*} : \beta_{g1|\mathcal{S}^*} = \mathbf{0}$ to declare DE genes.

The Proposed Backward Selection Algorithm

Let $\mathcal{S} \subset \{1, \dots, k\}$. For any $j \in \mathcal{S}$, let $\mathbf{p}_{j|\mathcal{S}}$ denote the vector of m p -values obtained by testing $H_{0gj}^{\mathcal{S}} : \beta_{gj|\mathcal{S}} = \mathbf{0}$ for $g = 1, \dots, m$. Let $r(\mathbf{p}_{j|\mathcal{S}})$ be a measure of the relevance of \mathbf{x}_j in model (1). Let $\mathcal{S}_1 = \{1, \dots, k\}$:

- 1 Compute $\mathbf{p}_{j|\mathcal{S}_1}$ for all $j \in \mathcal{S}_1$.
- 2 Let \mathbf{q}_ℓ be the vector of q -values obtained from $\mathbf{p}_{1|\mathcal{S}_\ell}$.
- 3 Let $R_\ell(\alpha)$ be the number of q -values in \mathbf{q}_ℓ less than or equal to a user-defined FDR threshold α .
- 4 If $\mathcal{S}_\ell \neq \{1\}$, find j^* so that $r(\mathbf{p}_{j^*|\mathcal{S}}) < r(\mathbf{p}_{j|\mathcal{S}})$ for all $j \in \mathcal{S}_\ell \setminus \{j^*\}$.
- 5 If $j^* = 1$, stop iterating. Otherwise, carry out the $\ell + 1$ st iteration with $\mathcal{S}_{\ell+1} = \mathcal{S}_\ell \setminus \{j^*\}$.

Suppose the iterative procedure concludes after L iterations, and let ℓ^* be the smallest element of $\{1, \dots, L\}$ such that $R_{\ell^*}(\alpha) \geq R_\ell(\alpha)$ for all $\ell \in \{1, \dots, L\}$. We set $\hat{\mathcal{S}} = \mathcal{S}_{\ell^*}$ and base our inference about differential expression on the fit of model (1).

Analysis of RFI RNA-seq Dataset

The dataset contains RNA-seq count from 31 animals of 2 genetic lines which has 12280 genes with average count greater than 8. We test those genes for differential expression between two lines (\mathbf{x}_1) in the presence of the thirteen categorical and continuous covariates: RFI (\mathbf{x}_2 , continuous), Diet (\mathbf{x}_3 , categorical, 2 levels), Baso (\mathbf{x}_4 , continuous), Eosi (\mathbf{x}_5 , continuous), Lymph (\mathbf{x}_6 , continuous), Mono (\mathbf{x}_7 , continuous), Neut (\mathbf{x}_8 , continuous), Concb (\mathbf{x}_9 , continuous), Conca (\mathbf{x}_{10} , continuous), RINb (\mathbf{x}_{11} , continuous), RINa (\mathbf{x}_{12} , continuous), Block (\mathbf{x}_{13} , categorical, 4 levels), and Order (\mathbf{x}_{14} , categorical, 4 levels). The best subset of relevant covariates is obtained at iteration $\ell = 7$, and number of DE genes is 448 when FDR is controlled at 5%.

RFI RNA-seq Data Analysis - pvalue05



RFI RNA-seq Data Analysis - ks



Simulation Study

1. Simulation 1: Using the set of covariates at iteration $\ell = 7$. $\mathcal{S}_\ell = \{1, 4, 5, 7, 8, 9, 12, 13\}$, the count data y_{gi} of gene g ($g = 1, \dots, 12280$) and pig i ($i = 1, \dots, 31$) has $\text{NB}(\mu_{gi}, \omega_g)$, where $\text{Var}(y_{gi}) = \mu_{gi} + \omega_g \mu_{gi}^2$, and

$$\log(\mu_{gi}) = o_i + \beta_{g0|\mathcal{S}_\ell} + \sum_{j \in \mathcal{S}_\ell} \mathbf{x}'_{ij} \beta_{gj|\mathcal{S}_\ell}.$$

Suppose that

$$\hat{\omega}_{g|\mathcal{S}_\ell}, \hat{\beta}_{g0|\mathcal{S}_\ell}, \hat{\beta}_{gj|\mathcal{S}_\ell} \ (j \in \mathcal{S}_\ell) \text{ are the estimates of } \omega_g, \beta_{g0|\mathcal{S}_\ell}, \beta_{gj|\mathcal{S}_\ell},$$

respectively, for $g = 1, \dots, 12280$. Let \hat{m}_0 be the estimate of m_0 , the number of true null hypotheses among all $m = 12280$ hypotheses tested (for line effect, i.e., \mathbf{x}_1). Let $\mathbf{q}_\ell = (q_{1|\mathcal{S}_\ell}, \dots, q_{m|\mathcal{S}_\ell})$ be the vector of q -values obtained from testing $H_{0g1}^{\mathcal{S}_\ell} : \beta_{g1|\mathcal{S}_\ell} = \mathbf{0}$, $g = 1, \dots, m$. Suppose that the order statistic of components of \mathbf{q}_ℓ is $q_{(1)|\mathcal{S}_\ell} \leq \dots \leq q_{(m)|\mathcal{S}_\ell}$. Set

$$\tilde{\beta}_{g1|\mathcal{S}_\ell} = \begin{cases} \mathbf{0} & \text{if } g \in \mathcal{G}_{0|\mathcal{S}_\ell} := \{g : q_g \geq q_{(m-\hat{m}_0+1)|\mathcal{S}_\ell}\} \\ \hat{\beta}_{g1|\mathcal{S}_\ell} & \text{if } g \in \mathcal{G}_{1|\mathcal{S}_\ell} := \{g : q_g < q_{(m-\hat{m}_0+1)|\mathcal{S}_\ell}\}. \end{cases}$$

$$\log(\tilde{\mu}_{gi|\mathcal{S}_\ell}) = o_i + \hat{\beta}_{g0|\mathcal{S}_\ell} + \mathbf{x}'_{i1} \tilde{\beta}_{g1|\mathcal{S}_\ell} + \sum_{j \in \mathcal{S}_\ell \setminus \{1\}} \mathbf{x}'_{ij} \hat{\beta}_{gj|\mathcal{S}_\ell}.$$

We simulate 100 datasets, each dataset consists of $m^* = 5000$ genes with different π , the proportion of EE genes, $\pi \in \{0.6, 0.7, 0.8, 0.9\}$. First, we randomly sample $\mathcal{G}_0^* = \{g_1, \dots, g_{m^*\pi}\} \subset \mathcal{G}_{0|\mathcal{S}_\ell}$, $\mathcal{G}_1^* = \{g_{m^*\pi+1}, \dots, g_{m^*}\} \subset \mathcal{G}_{1|\mathcal{S}_\ell}$. Finally, for each $h \in \{1, \dots, m^*\}$, $i \in \{1, \dots, 31\}$, we simulate \tilde{y}_{hi} from $\text{NB}(\tilde{\mu}_{gh|\mathcal{S}_\ell}, \hat{\omega}_{gh|\mathcal{S}_\ell})$. The collection of \tilde{y}_{hi} is the simulated count data we need.

- Simulation 2: Using the Set of Covariates from Iteration $\ell \in \{6, 7, 8, 9\}$.
- Simulation 3: Using the set of Covariates from Iteration $\ell \in \{6, 7, 8, 9\}$ with modified RFI values (\mathbf{x}_2).

Simulation Results

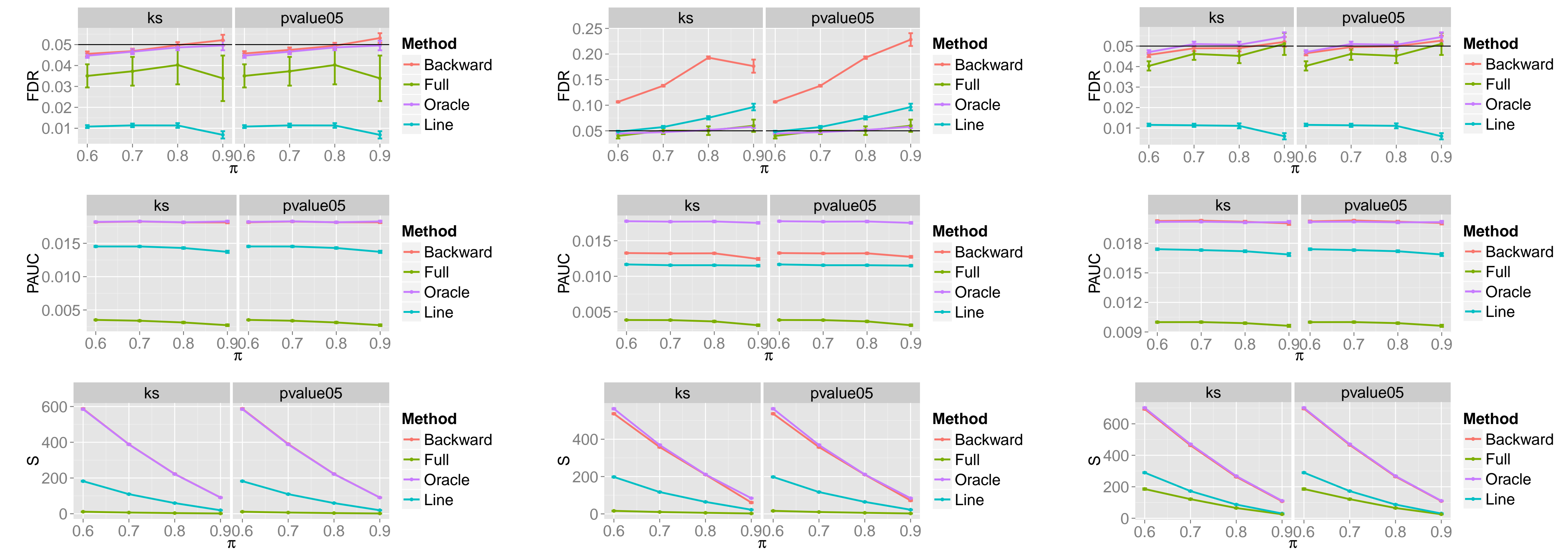


Figure 1: Simulation 1.

Figure 2: Simulation 2.

Figure 3: Simulation 3.

Acknowledgements

This material is based upon work supported by Agriculture and Food Research Initiative Competitive Grant No. 2011-68004-30336 from the USDA National Institute of Food and Agriculture, and the National Science Foundation (NSF) under Grant No. 0922746. This research also was supported by National Institute of General Medical Sciences (NIGMS) of the National Institutes of Health and the joint National Science Foundation / NIGMS Mathematical Biology Program under award number R01GM109458. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health or the National Science Foundation.