

Hanes Hall, B44  
Chapel Hill, NC 27510  
(919) 265-9114

February 4, 2016

Johns Hopkins University  
Bloomberg School of Public Health  
Office E3624  
615 North Wolfe Street  
Baltimore, MD 21205

**Dear Prof Leek,**

I am Meilei Jiang, a third year Phd student in statistics at UNC Chapel Hill, currently working with Steve Marron. Recently I am doing research on adjusting for batch effects in high dimensional data. I have read your papers and I quite like your approach, Surrogate Variable Analysis(SVA), as a method for batch correction and for addressing the data heterogeneity in the context of multiple testing dependence. However, I am quite puzzled by one aspect of SVA and wonder whether you can explain the reason behind your choice. My modification seems simpler and better than the published version, so I wonder if I am missing something.

In particular, I think we can improve the performance of SVA by modifying one step: weighting the residual matrix  $R$  by the probability  $\Pr(\gamma_i \neq \vec{0} | X, S, \hat{G})$  instead of weighting the data matrix  $X$  by the probability  $\Pr(b_i = \vec{0} \ \& \ \gamma_i \neq \vec{0} | X, S, \hat{G})$ .

Based on a simple simulation study, this new approach to SVA improves the performance of IRW-SVA in the case that there exist no genes (variables) which are strongly associated with unmeasured confounders but are not associated with measured factors. The two approaches have similar performance when such a subset exists.

The simulation study is conducted using your package ‘sva’ in the R software. You can find the simulation study results in the attachment.

Thank you for your time and consideration.

I look forward to your reply.

Sincerely,

Meilei Jiang