

Removing batch-effects from expression data

Maarten van Iterson

Leiden University Medical Center
Department of Molecular Epidemiology

June 23, 2017

Batch-Effects

In differential expression analyses there are primary variables of interest and often other nuisance factors, technical or biological, that introduce unwanted variation.

Batch-Effects

In differential expression analyses there are primary variables of interest and often other nuisance factors, technical or biological, that introduce unwanted variation.

- biological nuisance factors:
 - gender, age, white blood cell composition, etc.

Batch-Effects

In differential expression analyses there are primary variables of interest and often other nuisance factors, technical or biological, that introduce unwanted variation.

- biological nuisance factors:
 - gender, age, white blood cell composition, etc.
- technical nuisance factors (batch effects):
 - lab, sequence machine, library generation date, operator, etc.

Batch-Effects

In differential expression analyses there are primary variables of interest and often other nuisance factors, technical or biological, that introduce unwanted variation.

- biological nuisance factors:
 - gender, age, white blood cell composition, etc.
- technical nuisance factors (batch effects):
 - lab, sequence machine, library generation date, operator, etc.
- Often not all factors are known!

Batch-Effects

In differential expression analyses there are primary variables of interest and often other nuisance factors, technical or biological, that introduce unwanted variation.

- biological nuisance factors:
 - gender, age, white blood cell composition, etc.
- technical nuisance factors (batch effects):
 - lab, sequence machine, library generation date, operator, etc.
- Often not all factors are known!

Confounding occurs when there is correlation between primary variable of interest and the outcome

GEUVADIS RNAseq data¹

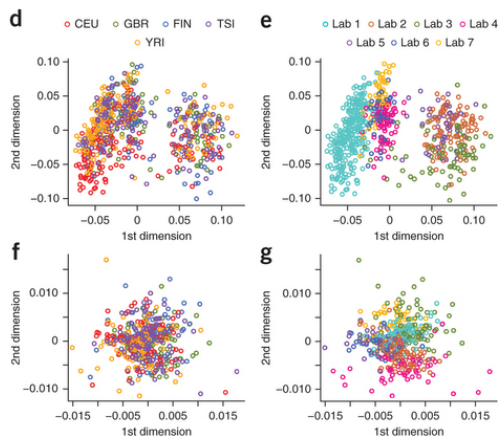


Figure 1 : (d) MDS plot of RNAseq data before batch correction colored by population and (e) colored by laboratory, (f) after batch correction colored by population and (g) colored by laboratory.

¹t Hoen, P. A. et al. (2013). Reproducibility of high-throughput mRNA and small RNA sequencing across laboratories.

Batch correction methods

- normalization methods:
 - quantile normalization, trimmed mean of M-values (TMM)
edgeR

¹Hansen, K. D., Irizarry, R. A., and Wu, Z. (2012). [Removing technical variability in RNA-seq data using conditional quantile normalization.](#)
Biostatistics, 13(2):204–216

²Risso, D., Schwartz, K., Sherlock, G., and Dudoit, S. (2011). [GC-content normalization for RNA-Seq data.](#)
BMC Bioinformatics, 12:480

Batch correction methods

- normalization methods:
 - quantile normalization, trimmed mean of M-values (TMM)
edgeR
- technology specific:
 - within-plate-, print-tip-normalization, etc.
 - GC-bias correction methods *cqn*¹, *EDASeq*²

¹Hansen, K. D., Irizarry, R. A., and Wu, Z. (2012). [Removing technical variability in RNA-seq data using conditional quantile normalization.](#) *Biostatistics*, 13(2):204–216

²Risso, D., Schwartz, K., Sherlock, G., and Dudoit, S. (2011). [GC-content normalization for RNA-Seq data.](#) *BMC Bioinformatics*, 12:480

Batch correction methods

- normalization methods:
 - quantile normalization, trimmed mean of M-values (TMM)
edgeR
- technology specific:
 - within-plate-, print-tip-normalization, etc.
 - GC-bias correction methods *cqn*¹, *EDASeq*²
- Batch correction methods:
 - **Nuisance factors are known:** linear model, ComBat
 - **Nuisance factors are unknown:** estimate batch-effects from the data
 - controls e.g. spike-ins or housekeeping
 - principal components
 - ...

¹Hansen, K. D., Irizarry, R. A., and Wu, Z. (2012). [Removing technical variability in RNA-seq data using conditional quantile normalization.](#)
Biostatistics, 13(2):204–216

²Risso, D., Schwartz, K., Sherlock, G., and Dudoit, S. (2011). [GC-content normalization for RNA-Seq data.](#)
BMC Bioinformatics, 12:480

Normalization does not remove batch-effects

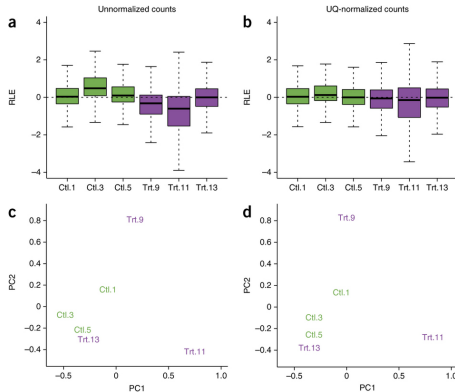


Figure 2 : Raw vs upper-quartile-normalized data ¹

¹Risso, D., Ngai, J., Speed, T. P., and Dudoit, S. (2014). [Normalization of RNA-seq data using factor analysis of control genes or samples.](#)

Removing batch-effects using RUV

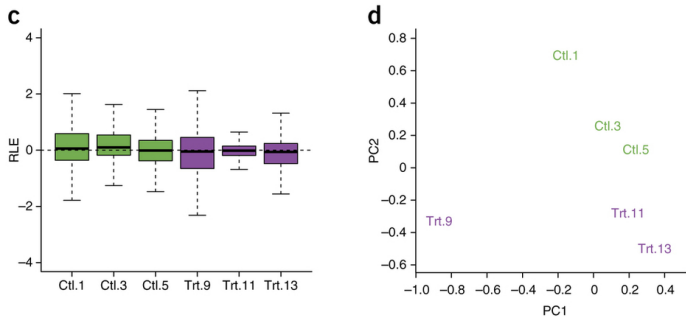


Figure 3 : RUV estimate and corrected

ComBat¹

Usage:

- Input: Known batches
- Output: Batch corrected expression matrix

¹Johnson, W. E., Li, C., and Rabinovic, A. (2007). [Adjusting batch effects in microarray expression data using empirical Bayes methods.](#)
Biostatistics, 8(1):118–127

ComBat¹

Usage:

- Input: Known batches
- Output: Batch corrected expression matrix

Method briefly:

- Mean center and standardize the variance of each batch for each gene independently
- Use an empirical Bayes approach to estimate robust mean and variance

¹Johnson, W. E., Li, C., and Rabinovic, A. (2007). [Adjusting batch effects in microarray expression data using empirical Bayes methods.](#) *Biostatistics*, 8(1):118–127

ComBat¹

Usage:

- Input: Known batches
- Output: Batch corrected expression matrix

Method briefly:

- Mean center and standardize the variance of each batch for each gene independently
- Use an empirical Bayes approach to estimate robust mean and variance

Remarks:

- Specially suited for small sample microarray studies
- Method is based on the same idea's for hypothesis testing as implemented in *limma*

R implementation available within the *sva* package

¹ Johnson, W. E., Li, C., and Rabinovic, A. (2007). [Adjusting batch effects in microarray expression data using empirical Bayes methods.](#) *Biostatistics*, 8(1):118–127

Surrogate variable analysis¹

Usage:

- Input: Does not use known factors but estimates a set of *surrogate variables*
- Optimize: number of surrogate variables
- Output: Estimated surrogate variables
- Testing: Include surrogate variables in a (generalized) linear model

¹Leek, J. T. and Storey, J. D. (2007). [Capturing heterogeneity in gene expression studies by surrogate variable analysis.](#)

Surrogate variable analysis¹

Usage:

- Input: Does not use known factors but estimates a set of *surrogate variables*
- Optimize: number of surrogate variables
- Output: Estimated surrogate variables
- Testing: Include surrogate variables in a (generalized) linear model

Method briefly:

- Constructs *surrogate variables* from a set of genes that are not associated with the biological factor of interest but are affected by unknown batches: principal component analysis on the residuals
- R implementation available [sva](#) package

¹Leek, J. T. and Storey, J. D. (2007). [Capturing heterogeneity in gene expression studies by surrogate variable analysis.](#)

Removing unwanted variation (RUV)¹

Usage:

- Input: Does not use known factors but estimates a set of factors describing the *unwanted variation*
- Optimize: Number of unknown factors
- Output: Estimated batch-effects
- Testing: Include estimated batch-effects in a (generalized) linear model

Method briefly:

¹Risso, D., Ngai, J., Speed, T. P., and Dudoit, S. (2014). [Normalization of RNA-seq data using factor analysis of control genes or samples.](#)
Nat. Biotechnol., 32(9):896–902

Removing unwanted variation (RUV)¹

Usage:

- Input: Does not use known factors but estimates a set of factors describing the *unwanted variation*
- Optimize: Number of unknown factors
- Output: Estimated batch-effects
- Testing: Include estimated batch-effects in a (generalized) linear model

Method briefly:

- Factor analysis (PC) on the residuals of the control genes
- R implementation available *RUVseq*

¹Risso, D., Ngai, J., Speed, T. P., and Dudoit, S. (2014). [Normalization of RNA-seq data using factor analysis of control genes or samples.](#) *Nat. Biotechnol.*, 32(9):896–902

Usage:

- Input: Does not use known factors but estimates a set of *latent factors* describing the *unobserved confounding factors*
- Optimize: Number of latent factor
- Output: Estimated latent factors
- Testing: hypotheses testing included (robust regression)

¹Wang, J., Zhao, Q., Hastie, T., and Owen, A. B. (2015). [Confounder Adjustment in Multiple Hypothesis Testing](#).
ArXiv e-prints

CATE¹

Usage:

- Input: Does not use known factors but estimates a set of *latent factors* describing the *unobserved confounding factors*
- Optimize: Number of latent factor
- Output: Estimated latent factors
- Testing: hypotheses testing included (robust regression)

Method briefly:

- Factor analysis on *residuals*
- R implementation available [cate](#)

¹Wang, J., Zhao, Q., Hastie, T., and Owen, A. B. (2015). [Confounder Adjustment in Multiple Hypothesis Testing](#).
[ArXiv e-prints](#)

Comparison from Leek¹

Simulated data with one group (Case/Control) and one batch

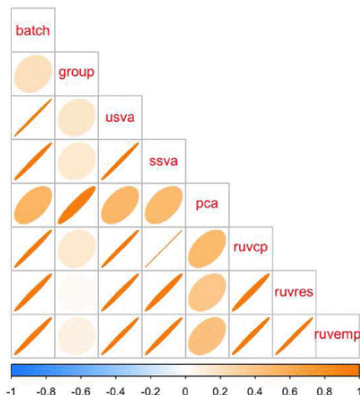


Figure 4 : Correlation between simulated batch and group variables and various batch estimates.

¹Leek, J. T. (2014). [svaseq: removing batch effects and other unwanted noise from sequencing data](#). *Nucleic Acids Res.*, 42(21)

Comparison from Leek

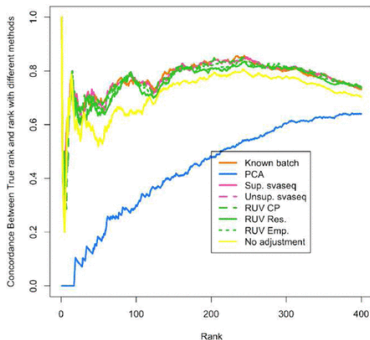


Figure 5 : Differential expression results for simulated data. A concordance at the top plot (CAT plot) shows the fraction of DE results that are concordant between the analysis with the true batch and the analyses using different batch estimates.

A few other methods

1. PEER¹ cran R package *peer*
2. isva² cran R package *isva*
3. RUV-4, RUV-inv, and RUV-rinv³ cran R package *ruv*

These methods can also be applied to other omics-data e.g. 450k DNA methylation data

¹Stegle, O., Parts, L., Piipari, M., Winn, J., and Durbin, R. (2012). [Using probabilistic estimation of expression residuals \(PEER\) to obtain increased power and interpretability of gene expression analyses.](#) *Nat Protoc*, 7(3):500–507

²Teschendorff, A. E., Zhuang, J., and Widschwendter, M. (2011). [Independent surrogate variable analysis to deconvolve confounding factors in large-scale microarray profiling studies.](#) *Bioinformatics*, 27(11):1496–1505

³Gagnon-Bartsch, J., Jacob, L., and Speed, T. (2013). [Removing unwanted variation from high dimensional data with negative controls.](#) *Tech Report*.

Background: General formulation of the problem

$$Y_{n \times p} = Z_{n \times q} \alpha_{q \times p} + W_{n \times k} \beta_{k \times p} + \epsilon_{n \times p} \quad (1)$$

- $Y_{n \times p}$: observed expr. data on n samples and p genes
- $Z_{n \times q}$: q known cov. including phenotype of interest
- $W_{n \times k}$: k unobserved cov.
- α, β : represent the unknown effects of the obs. and unobs. cov. on gene expr.

Background: Two extremes

Consider only known covariates or no covariates at all

1. Correct for known technical or biological covariates using a linear model:

$$Y = Z\alpha + \epsilon \quad (2)$$

(e.g. all β 's are zero)

2. Estimate batch effects using SVD/PC on $Y = U\Sigma V^T$

$$Y = W\beta + \epsilon, \quad (3)$$

$W_{n \times k} = V_{n \times k}$: are the first k principal components, the optimal k needs to be determined (e.g. now all α 's are zero)

Background: SVA

Step 1: fit covariates of interest

$$Y_{n \times p} = Z_{n \times 1} \alpha_{q \times p} + \epsilon_{n \times p} \quad (4)$$

$$R_{n \times p} = Y_{n \times p} - Z_{n \times 1} \hat{\alpha}_{q \times p} \quad (5)$$

Estimate unobserved batch effects using SVD/PC on
 $R = U \Sigma V^T$

Step 2: fit covariates of interest plus a few estimated unobserved batch effects

$$Y_{n \times p} = Z_{n \times q} \alpha_{q \times p} + W_{n \times r} \beta_{r \times p} + \epsilon_{n \times p} \quad (6)$$

$W_{n \times k} = V_{n \times k}$: are the first k principal components, the optimal k is determined using permutation test on the eigenvalues

Background: RUV-2

Step 1: Consider general model for p' control genes

$$Y_{n \times p'} = Z_{n \times 1} \alpha_{q \times p'} + W_{n \times k} \beta_{k \times p'} + \epsilon_{n \times p'} \quad (7)$$

the 'control gene assumption' is that $\alpha = 0$

$$Y_{n \times p'} = W_{n \times k} \beta_{k \times p'} + \epsilon_{n \times p'} \quad (8)$$

$W_{n \times k} = V_{n \times k}$ SVD/PC of $Y_{n \times p'} = U \Sigma V^T$

Step 2: fit covariates of interest plus a few estimated unobserved batch effects

$$Y_{n \times p} = Z_{n \times q} \alpha_{q \times p} + W_{n \times k} \beta_{k \times p} + \epsilon_{n \times p} \quad (9)$$

the optimal k needs to be determined

if $\alpha \neq 0$ for some covariates use the SVD of $Y - YZ(Z^T Z)^{-1}Z^T$
again only on the control genes

CATE

Background: RUV-4, RUV-inv, and RUV-rinv

RUV-4: a hybrid between RUV-2 and SVA

RUV-inv: includes selection of the number of unobserved batches with a novel method for estimation of the variance . . .

RUV-rinv: ridge regression . . .

Background: isva

i.s.o. using linear uncorrelated eigenvectors in the SVA method
estimate statistically independent variables using independent
component analysis

Background: PEER

similar to SVA estimates unobserved covariates but uses empirical Bayesian statistics

Relation to other methods: LMM