

Estimating batch effect in Microarray data with Principal Variance Component Analysis(PVCA) method

Pierre R. Bushel, Jianying Li

April 15, 2025

Contents

1	Introduction	1
2	Installation	2
3	An example run on published Golub data	2
4	Session	3

1 Introduction

Often times "batch effects" are present in microarray data due to any number of factors, including e.g. a poor experimental design or when the gene expression data is combined from different studies with limited standardization. To estimate the variability of experimental effects including batch, a novel hybrid approach known as principal variance component analysis (PVCA) has been developed. The approach leverages the strengths of two very popular data analysis methods: first, principal component analysis (PCA) is used to efficiently reduce data dimension with maintaining the majority of the variability in the data, and variance components analysis (VCA) fits a mixed linear model using factors of interest as random effects to estimate and partition the total variability. The PVCA approach can be used as a screening tool to determine which sources of variability

(biological, technical or other) are most prominent in a given microarray data set. Using the eigenvalues associated with their corresponding eigenvectors as weights, associated variations of all factors are standardized and the magnitude of each source of variability (including each batch effect) is presented as a proportion of total variance. Although PVCA is a generic approach for quantifying the corresponding proportion of variation of each effect, it can be a handy assessment for estimating batch effect before and after batch normalization.

The `pvca` package implements the method described in the book *Batch Effects and Noise in Microarray Experiment*, chapter 12 "Principal Variance Components Analysis: Estimating Batch Effects in Microarray Gene Expression Data": *Jianying Li, Pierre R Bushel, Tzu-Ming Chu, and Russell D Wolfinger 2010*

The `pvca` method was applied in the paper: *Michael J Boedigheimer, Russell D Wolfinger, Michael B Bass, Pierre R Bushel, Jeff W Chou, Matthew Cooper, J Christopher Corton, Jennifer Foster, Susan Hester, Janice S Lee, Fenglong Liu, Jie Liu, Hui-Rong Qian, John Quackenbush, Syril Pettit and Karol L Thompson* 2008 "Sources of variation in baseline gene expression levels from toxicogenomics study control animals across multiple laboratories" *BMC Genomics* 2008, June 12, 9:285

2 Installation

Simply skip this section if one has been familiar with the usual Bioconductor installation process. Assume that a recent version of R has been correctly installed.

Install the packages from the Bioconductor repository, using the `biocLite` utility. Within R console, type

```
> if (!requireNamespace("BiocManager", quietly=TRUE))
+   install.packages("BiocManager")
> BiocManager::install("pvca")
```

Installation using the `biocLite` utility automatically handles the package dependencies. The `pvca` package depends on the packages like `lme4` etc., which can be installed when `pvca` package is stalled.

3 An example run on published Golub data

We use `Golub` dataset in the package `golubEsets` as an example to illustrate the PVCA batch effect estimation procedure. This dataset contains 7129 genes from Microarray

data on 72 samples from a leukemia study. It is a merged dataset and we are testing variability from each factor and their two-ways interactions. The package performs PVCA on this merged data and produces an estimation of the proportion (as possible batch effect) each factor and interaction contribute. The figure shows the proportion in a bar chart. .

```
> library(golubEsets)
> library(pvca)
> data(Golub_Merge)
> pct_threshold <- 0.6
> batch.factors <- c("ALL.AML", "BM.PB", "Source")
> pvcaObj <- pvcaBatchAssess (Golub_Merge, batch.factors, pct_threshold)
>
```

We can plot the source of potential batch effect in proportioning as shown in Figure 1.

```
> bp <- barplot(pvcaObj$dat, xlab = "Effects",
+             ylab = "Weighted average proportion variance",
+             ylim= c(0,1.1), col = c("blue"), las=2,
+             main="PVCA estimation bar chart")
> axis(1, at = bp, labels = pvcaObj$label, xlab = "Effects", cex.axis = 0.5, las=2)
> values = pvcaObj$dat
> new_values = round(values , 3)
> text(bp,pvcaObj$dat,labels = new_values, pos=3, cex = 0.8)
>
```

4 Session

```
> print(sessionInfo())
```

```
R version 4.5.0 RC (2025-04-03 r88103 ucrt)
Platform: x86_64-w64-mingw32/x64
Running under: Windows Server 2022 x64 (build 20348)
```

```
Matrix products: default
LAPACK version 3.12.1
```

```
locale:
[1] LC_COLLATE=C
```

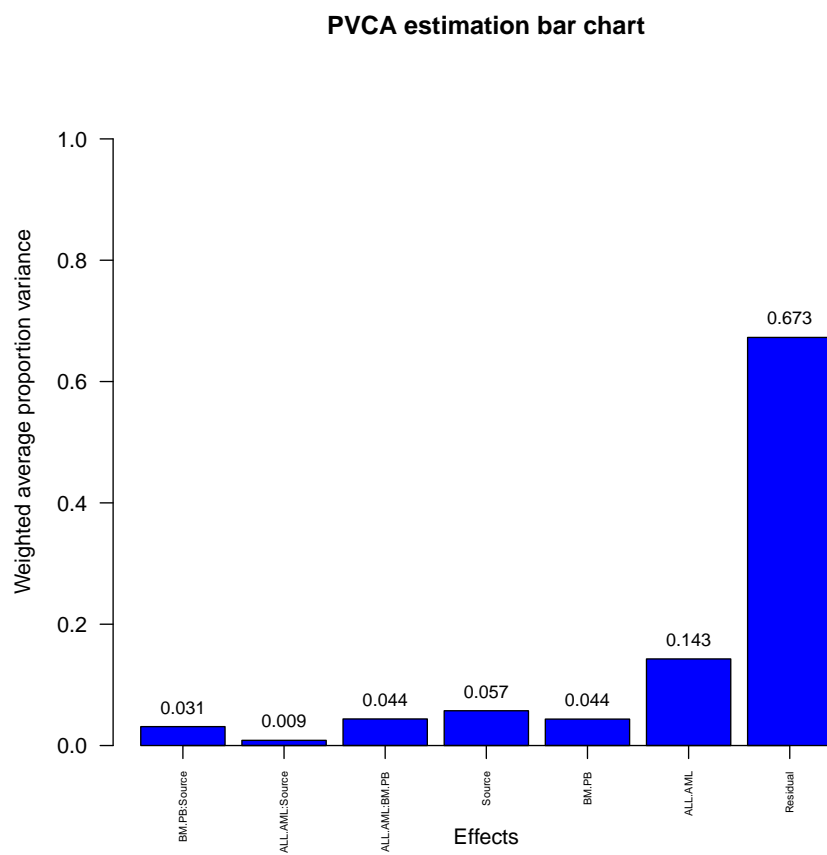


Figure 1: The bar chart shows the proportion of batch effect from possible source.

```
[2] LC_CTYPE=English_United States.utf8
[3] LC_MONETARY=English_United States.utf8
[4] LC_NUMERIC=C
[5] LC_TIME=English_United States.utf8
```

```
time zone: America/New_York
tzcode source: internal
```

```
attached base packages:
```

```
[1] stats      graphics  grDevices  utils      datasets  methods    base
```

```
other attached packages:
```

```
[1] pvca_1.49.0      golubEsets_1.49.0  Biobase_2.69.0
[4] BiocGenerics_0.55.0 generics_0.1.3
```

```
loaded via a namespace (and not attached):
```

```
[1] Matrix_1.7-3      gtable_0.3.6      limma_3.65.0
[4] preprocessCore_1.71.0 dplyr_1.1.4      compiler_4.5.0
[7] BiocManager_1.30.25 tidyselect_1.2.1  Rcpp_1.0.14
[10] splines_4.5.0      scales_1.3.0      boot_1.3-31
[13] statmod_1.5.0      lattice_0.22-7    ggplot2_3.5.2
[16] R6_2.6.1           rbibutils_2.3     MASS_7.3-65
[19] tibble_3.2.1       nloptr_2.2.1      vsn_3.77.0
[22] munsell_0.5.1      affy_1.87.0       minqa_1.2.8
[25] pillar_1.10.2      affyio_1.79.0     rlang_1.1.6
[28] cli_3.6.4          magrittr_2.0.3    Rdpack_2.6.4
[31] grid_4.5.0         lme4_1.1-37       lifecycle_1.0.4
[34] nlme_3.1-168       reformulas_0.4.0  vctrs_0.6.5
[37] glue_1.8.0         colorspace_2.1-1  tools_4.5.0
[40] pkgconfig_2.0.3
```