CQN (Conditional Quantile Normalization)

Kasper Daniel Hansen

Zhijin Wu

khansen@jhsph.edu

zhijin_wu@brown.edu

Modified: August 8, 2012. Compiled: October 24, 2025

Introduction

This package contains the CQN (conditional quantile normalization) method for normalizing RNA-seq datasets. This method is described in [1].

```
> library(cqn)
```

Data

As an example we use ten samples from Montgomery [2]. The data has been processed as described in [1]. First we have the region by sample count matrix

```
> data(montgomery.subset)
```

```
[1] 23552 10
```

> montgomery.subset[1:4,1:4]

	NA06985	NA06994	NA07037	NA10847
ENSG00000000419	69	54	67	70
ENSG00000000457	53	37	27	41
ENSG00000000460	12	25	33	22
ENSG00000000938	168	270	140	103

> colnames(montgomery.subset)

> library(scales)

> dim(montgomery.subset)

```
[1] "NA06985" "NA06994" "NA07037" "NA10847" "NA11920" "NA11918" [7] "NA11931" "NA12003" "NA12006" "NA12287"
```

Because of (disc) space issues, We have removed all genes that have zero counts in all 10 samples. Next we have the *sizeFactors* which simply tells us how deep each sample was sequenced:

```
> data(sizeFactors.subset)
> sizeFactors.subset[1:4]

NA06985 NA06994 NA07037 NA10847
3107420 2388948 3087234 2852972
```

Finally, we have a matrix containing length and GC-content for each gene.

```
> data(uCovar)
> head(uCovar)
```

	length	gccontent
ENSG00000000419	1207	0.3976802
ENSG00000000457	2861	0.4606781
ENSG00000000460	4912	0.4338355
ENSG00000000938	3524	0.5749149
ENSG00000000971	8214	0.3613343
ENSG0000001036	2590	0.4312741

Note that the row ordering of the count matrix is the same as the row ordering of the matrix containing length and GC-content and that the sizeFactor vector has the same column order as the count matrix. We can formally check this

```
> stopifnot(all(rownames(montgomery.subset) == rownames(uCovar)))
> stopifnot(colnames(montgomery.subset) == names(sizeFactors.subset))
```

Normalization

The methodology is described in [1]. The main workhorse is the function cqn which fits the following model

$$\log_2(\text{RPM}) = s(x) + s(\log_2(\text{length}))$$

where x is some covariate, s are smooth functions (specifically natural cubic splines with 5 knots), and RPM are "reads per millions". It is also possible to just fit a model like

$$\log_2(\mathsf{RPKM}) = s(x)$$

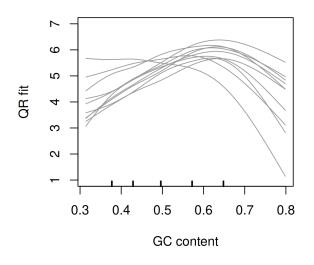
In this model gene length is included as a known offset. This is done by using the cqn (lengthMethod = "fixed"). If this is done, and lengths is equal to 1000, it is equivalent to not using gene length at all.

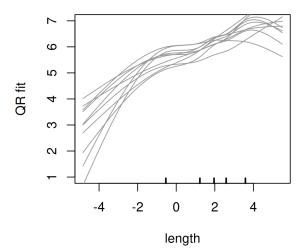
The basic call to cqn is relatively easy, we need the count matrix, a vector of lengths, a vector of GC content and a vector of sizeFactors. Make sure that they all have the same ordering.

This normalized matrix is similar, but not equivalent, to the data examined in [1]. The main differences are (1) in [1] we normalize 60 samples together, not 10 and (2) we have removed all genes with zero counts in all 10 samples.

We can examine plots of systematic effects by using canplot. The n argument refers to the systematic effect, n=1 is always the covariate specified by the x argument above, while n=2 is lengths.

```
> par(mfrow=c(1,2))
> cqnplot(cqn.subset, n = 1, xlab = "GC content", lty = 1, ylim = c(1,7))
> cqnplot(cqn.subset, n = 2, xlab = "length", lty = 1, ylim = c(1,7))
```





The normalized expression values are

```
> RPKM.cqn <- cqn.subset$y + cqn.subset$offset
> RPKM.cqn[1:4,1:4]
```

```
NA06985 NA06994 NA07037 NA10847
ENSG00000000419 5.762645 5.569588 5.548111 5.976500
ENSG00000000457 4.436670 4.110060 3.393547 4.139232
ENSG00000000460 2.602654 3.443828 3.776863 3.067707
ENSG00000000938 5.152698 6.084822 4.698430 4.281723
```

These values are on the log_2 -scale.

We can do a MA plot of these fold changes, and compare it to fold changes based on standard RPKM. First we compute the standard RPKM (on a log_2 scale):

```
> RPM <- sweep(log2(montgomery.subset + 1), 2, log2(sizeFactors.subset/10^6 > RPKM.std <- sweep(RPM, 1, log2(uCovar$length / 10^3))
```

We now look at differential expression between two groups of samples. We use the same grouping as in [1], namely

```
> grp1 <- c("NA06985", "NA06994", "NA07037", "NA10847", "NA11920")
> grp2 <- c("NA11918", "NA11931", "NA12003", "NA12006", "NA12287")
```

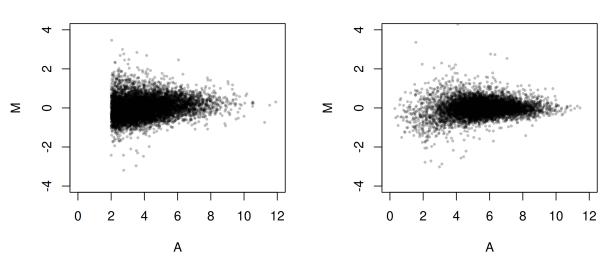
We now do an MA-plot, but we only choose to plot genes with average standard \log_2 -RPKM of $\log_2(5)$ or greater, and we also form the M and A values:

```
> whGenes <- which(rowMeans(RPKM.std) >= 2 & uCovar$length >= 100)
> M.std <- rowMeans(RPKM.std[whGenes, grp1]) - rowMeans(RPKM.std[whGenes, grp1]) - rowMeans(RPKM.std[whGenes, grp1])
> M.cqn <- rowMeans(RPKM.cqn[whGenes, grp1]) - rowMeans(RPKM.cqn[whGenes, grp1]) - rowMeans(RPKM.cqn[whGenes, grp1])</pre>
```

Now we do the MA plots, with alpha-blending

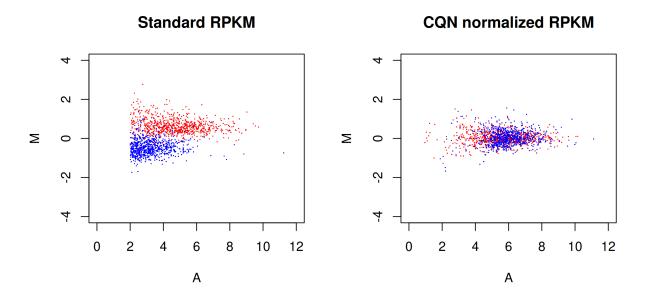
Standard RPKM

CQN normalized RPKM



We can also color the genes according to whether they have high/low GC-content. Here one needs to be careful, because of overplotting. One solution is to leave out all genes with intermediate GC content. We define high/low GC content as the 10% most extreme genes:

```
+ ylim = c(-4,4), xlim = c(0,12), col = "red")
> points(A.std[whLow], M.std[whLow], cex = 0.2, pch = 16, col = "blue")
> plot(A.cqn[whHigh], M.cqn[whHigh], cex = 0.2, pch = 16, xlab = "A",
+ ylab = "M", main = "CQN normalized RPKM",
+ ylim = c(-4,4), xlim = c(0,12), col = "red")
> points(A.cqn[whLow], M.cqn[whLow], cex = 0.2, pch = 16, col = "blue")
```



Note that genes/regions with very small counts should not be relied upon, even if the CQN normalized fold change are big. They should be filtered out using some kind of statistical test, good packages for this are *DESeq*[3] and *edgeR*[4, 5].

Import into edgeR

First we construct a DGEList. In the groups argument we use that the first 5 samples (columns) in montgomery. subset is what we earlier called grp1 and the last 5 samples (columns) are grp2.

```
> library(edgeR)
> d.mont <- DGEList(counts = montgomery.subset, lib.size = sizeFactors.subset
+ group = rep(c("grp1", "grp2"), each = 5), genes = uCoval</pre>
```

In this object we cannot (unfortunately, yet) also store the computed offsets. Since we will use the offsets computed by cqn, there is no need to normalize using the normalization tools from edgeR, such as calcNormFactors. Also, as is clearly described in the edgeR user's guide, the lib.size is unnecessary, since we plan to use the offsets computed from cqn.

However, we need to use the component glm.offset which is on the natural logarithmic scale and also includes correcting for sizeFactors. It is possible to include the offset directly into the DGEList, by post-processing the output like

```
> ## Not run
> d.mont$offset <- cqn.subset$qlm.offset</pre>
```

Using *edgeR* is well described in the user's guide, and we refer to that document for further information. The analysis presented below should be thought of as an example, and not necessarily the best analysis of this data.

The first step is estimating the dispersion parameter(s). Several methods exists, such as estimateGLMCommonD or estimateTagwiseDisp. We also need to setup a design matrix, which is particular simple for this two group comparison. Further information about constructing design matrices may be found in both the *edgeR* user's guide and the *limma* user's guide.

```
> design <- model.matrix(~ d.mont$sample$group)
> d.mont$offset <- cqn.subset$glm.offset
> d.mont.cqn <- estimateGLMCommonDisp(d.mont, design = design)</pre>
```

After fitting the dispersion parameter(s), we need to fit the model, and do a test for significance of the parameter of interest. With this design matrix, there are two coefficients. The first coefficient is just an intercept (overall level of expression for the gene) and it is (usually) not meaningful to test for this effect. Instead, the interesting coefficient is the second one that encodes a group difference.

```
> efit.cqn <- glmFit(d.mont.cqn, design = design)</pre>
> elrt.cqn <- glmLRT(efit.cqn, coef = 2)</pre>
> topTags(elrt.cqn, n = 2)
              d.mont$sample$groupgrp2
Coefficient:
                length gccontent
                                      logFC
                                               logCPM
                    365 0.5835616 -10.29118 6.362866 126.2816
ENSG00000211642
                    411 0.5888078 -10.11052 5.999868 120.5614
ENSG00000211660
                       PValue
                                       FDR
ENSG00000211642 2.668046e-29 6.283782e-25
ENSG00000211660 4.766936e-28 5.613543e-24
```

topTags shows (per default) the "top 10" genes. In this case, since we have biological replicates and just a random group structure, we would expect no differentially expression genes. Instead we get

```
> summary(decideTests(elrt.cqn))
```

```
d.mont$sample$groupgrp2
Down 147
NotSig 22968
Up 437
```

significantly differentially expressed at an FDR (false discovery rate) of 5%. We may contrast this with the result of not using *cqn*:

In this evaluation, it is not clear that using CQN is better.

What is arguably as important is that we achieve a much better estimation of the fold change using *cqn*.

Question and Answers

Can I run cqn() on only 1 sample?

CQN is meant to normalize several samples together. It is not clear that it makes sense at all to use this normalization technique on a single sample. But it is possible.

Can I use this for small RNA-seq (microRNAs)?

We do not have personal experience with using CQN to normalize small RNA sequencing data. However, we believe it might be beneficial. As always, it is *highly* recommended to evaluate whether it is necessary and beneficial.

One special aspect of small RNAs is that they all have very similar length. Fitting a model with a smooth effect of gene length might very well lead to mathematical instability (you get an error). This can be avoided by using the argument lengthMethod = "fixed" which just divides the gene counts by the gene length instead of using a smooth function. Additionally, it may be coupled with setting lengths = 1 which completely removes gene length from the model.

Could it be true that genes with higher GC content are higher expressed?

It has been suggested that genes that are either extremely high or extremely low expressed are under some form of selection leading to "extreme" GC content. What CQN does, is making the effect of GC content comparable across samples, and we show in [1] that this leads to improved inference. It also flattens the effect of GC content on gene expression, but we believe this is better than having the effect of GC content depend on the sample.

Does cqn remove batch effects?

No, unless a batch effect only (or mainly) affects your measurements through GC content. We believe that the sample-specific effect of GC content on gene expression is a kind of batch effect, but is unlikely to be the only one. CQN does normalize your RNA-seq data in the same way that say quantile normalization normalizes microarray data, but such normalization does not remove batch effects.

I don't understand the difference between offset and glm.offset?

This comes from a historical error. In our paper, we use the quantity

```
> cqn$y + cqn$offset
```

as the CQN-corrected estimated expression measures. However, the offset quantity is on the wrong scale for inclusion into a GLM-type model (like edgeR or DEseq2). For this purpose, use glm.offset. We have kept the original naming in order to achieve backwards compatibility.

SessionInfo

- R Under development (unstable) (2025-10-20 r88955), x86_64-pc-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_GB, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Time zone: America/New_York
- TZcode source: system (glibc)
- Running under: Ubuntu 24.04.3 LTS
- Matrix products: default
- BLAS: /home/biocbuild/bbs-3.23-bioc/R/lib/libRblas.so

- LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.12.0
- Base packages: base, datasets, grDevices, graphics, methods, stats, utils
- Other packages: cqn 1.55.0, edgeR 4.7.6, limma 3.65.7, mclust 6.1.1, scales 1.4.0
- Loaded via a namespace (and not attached): MASS 7.3-65, Matrix 1.7-4, MatrixModels 0.5-4, R6 2.6.1, RColorBrewer 1.1-3, SparseM 1.84-2, cli 3.6.5, compiler 4.6.0, dichromat 2.0-0.1, farver 2.1.2, glue 1.8.0, grid 4.6.0, lattice 0.22-7, lifecycle 1.0.4, locfit 1.5-9.12, nor1mix 1.3-3, quantreg 6.1, rlang 1.1.6, splines 4.6.0, statmod 1.5.1, survival 3.8-3, tools 4.6.0

References

- [1] KD Hansen, RA Irizarry, and Z Wu. Removing technical variability in RNA-seq data using conditional quantile normalization. *Biostatistics* 2012, **13**(2), 204–216. DOI: 10.1093/biostatistics/kxr054.
- [2] SB Montgomery *et al.* Transcriptome genetics using second generation sequencing in a Caucasian population. *Nature* 2010, **464**, 773–777. DOI: .10.1038/nature08903
- [3] S Anders and W Huber. Differential expression analysis for sequence count data. *Genome Biology* 2010, **11**(10), R106. DOI: 10.1186/gb-2010-11-10-r106.
- [4] MD Robinson, DJ McCarthy, GK Smyth. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 2010, **26**(1), 139–140. DOI: 10.1093/bioinformatics/btp616.
- [5] DJ McCarthy, Y Chen, GK Smyth. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* 2012, **40**, 4288-4297. DOI: 10.1093/nar/gks042.