

Assessing signal/noise ratio before and after normalization

Wolfgang Huber

March 28, 2006

Contents

1	Introduction	1
2	Selection of PM and background features	2
3	Normalization	3
4	Alternative normalization methods	3
4.1	Without dropping the worst 5% probes	3
4.2	Background correction by MM	3
5	Assessment	4
5.1	Visually	4
5.2	Quantitatively	5
5.3	Dependency of the signal on GC content before and after normalization	8

1 Introduction

The purpose of this document is to assess the performance of the probe-response normalization by the function `normalizeByReference` in the *tilingArray* package. We use the example data from the David et al. [1] paper, which is provided in the *davidTiling* package.

```
> library("tilingArray")
> library("davidTiling")
> if (!exists("davidTiling")) data("davidTiling")
```

It contains 8 arrays with 6553600 features each. Three of them were hybridized to genomic DNA, which will use as a reference for the normalization, and five to RNA.

```
> dim(exprs(davidTiling))
```

```
[1] 6553600      8
```

```
> sampleNames(davidTiling)
```

```
[1] "09_11_04_S96_genDNA_16hrs_45C_noDMS0.cel"
[2] "041119_S96genDNA_re-hybe.cel"
[3] "041120_S96genDNA_re-hybe.cel"
```

```
[4] "05_04_27_2xpolyA_NAP3.cel"
[5] "05_04_26_2xpolyA_NAP2.cel"
[6] "05_04_20_2xpolyA_NAP_2to1.cel"
[7] "050409_totcDNA_14ug_no52.cel"
[8] "030505_totcDNA_15ug_affy.cel"
```

Some of the computations in this vignette will take a long time and require Gigabytes of RAM. I have sprinkled `cache()` and `if(exists(...))` statements around the most massive computations – this is simply so that during the writing of the vignette I can rerun `Sweave` on it in an incremental fashion without always having to redo all computations from scratch. Eventually they will go away, and for the comprehension of the computations they are irrelevant.

2 Selection of PM and background features

The design of the chip includes about 3 Mio *perfect match probes*, about 3 Mio *mismatch probes*, and a few thousand controls. The `probeAnno` object contains annotations for each probe. It was obtained by aligning the probe sequence to the genomic sequence of *S. cerevisiae* in August 2005. Please see its manual page and the script `makeProbeAnno.R` in the *davidTiling* package for details.

```
> if (!exists("probeAnno")) data("probeAnno")
```

We define two logical vectors: `isPM` is TRUE if a probe has a perfect match anywhere in the genome, `isBG` is TRUE if that match is outside any annotated feature on either strand.

```
> isPM = logical(nrow(exprs(davidTiling)))
> for (j in probeAnno$probeReverse) isPM[as.character(j) != ""] = TRUE
> isBG = (probeAnno$probeReverse$no_feature == "no" & probeAnno$probeDirect$no_feature ==
+       "no")
```

```
> sum(isPM)
```

```
[1] 3049103
```

```
> tab = table(isPM, isBG)
```

	isBG	
isPM	FALSE	TRUE
FALSE	3504497	0
TRUE	2219165	829938

There are 3504497 mismatch probes, whose signal we will not try to normalize, since we only have meaningful reference intensities for the 3049103 perfect match probes. For the background estimation, we will use the 829938 features for which we expect no specific signal. Some of them will indeed have specific signal, due to new transcripts that were not yet annotated in SGD in August 2005; but this is not a concern, since we will only use statistical properties of the signal distribution that are insensitive to a small set of outliers. Alternatively, as we show below, we could also use the signal from the MM probes for background estimation.

3 Normalization

We select the 5 arrays that had RNA hybridized to them and the 3 arrays that had DNA hybridized to them. We want to normalize the intensity readings from the former by using the latter as a reference.

```
> isRNA = davidTiling$nucleicAcid %in% c("poly(A) RNA", "total RNA")
> isDNA = davidTiling$nucleicAcid %in% "genomic DNA"
> stopifnot(sum(isRNA) == 5, sum(isDNA) == 3)

> pfn = sprintf("assessNorm-normalize%d.pdf", seq(along = which(isRNA)))
> if (!exists("xn2")) xn2 = cache("xn2", normalizeByReference(davidTiling[,
+   isRNA], davidTiling[, isDNA], pm = isPM, background = isBG,
+   plotFileNames = pfn))
```

4 Alternative normalization methods

4.1 Without dropping the worst 5% probes

For comparison, we also compare to the situation in which we do not throw out the weakest features, by setting `cutoffQuantile=0`.

```
> if (!exists("xn1")) xn1 = cache("xn1", normalizeByReference(davidTiling[,
+   isRNA], davidTiling[, isDNA], pm = isPM, background = isBG,
+   cutoffQuantile = 0))
```

4.2 Background correction by MM

Instead of the background correction based on “similar” no-target probes that is done in `normalizeByReference`, we can also consider background correction by MM probes.

There is a slight complication: the set of probes that we are considering as PM in this document is somewhat different from those that went into the design of PM/MM pairs on the array, since different versions of the yeast genome sequence were used. We take the intersection. For this, we determine the indices of the designed PM/MM pairs. The array has 2560 rows and 2560 columns. If we count the rows and columns from 0 to 2559, then the linear indices of the features (e.g. in the expression matrix of the *eSet* `davidTiling`) are given by $r \times 2560 + c$. The PM features lie in rows 1, 3, ..., 2557, their corresponding MM features in rows 2, 4, ..., 2558. We are going to use this information, as well as the probe sequences that are provided in the *Scerevisiaetilingprobe* package.

```
> library("Scerevisiaetilingprobe")
> nc = as.integer(2560)
> stopifnot(nc * nc == nrow(davidTiling))
> ipm = rep(as.integer(seq(1, nc - 3, by = 2)), each = nc) * nc +
+   (1:nc)
> seqPM = Scerevisiaetilingprobe$sequence[ipm]
> seqMM = Scerevisiaetilingprobe$sequence[ipm + nc]
> ispair = (!is.na(seqPM) & !is.na(seqMM) & (seqPM == complementSeq(seqMM,
+   start = 13, stop = 13)))
```

We see that out of the 3274240 probes in odd-numbered rows, listed by `ipm`, most do indeed have a corresponding mismatch probe in the row below.

```
> length(ispair)
```

```
[1] 3274240
```

```
> table(ispair)
```

```
ispair
  FALSE    TRUE
129718 3144522
```

We define vectors `PMind` and `MMind` that contain the indices of PM/MM pairs,

```
> PMind = ipm[ispair]
```

```
> MMind = PMind + nc
```

and a function `normalizeByReferenceWithMM` that is similar to `normalizeByReference` in the *tilingArray* package, except for that for each PM probe it uses the corresponding MM value for the background correction rather than the background estimator that is used in `normalizeByReference`.

```
> normalizeByReferenceWithMM = function(x, reference, pm, mm, cutoffQuantile = 0.05) {
+   refSig = 2^rowMeans(log2(exprs(reference)[pm, ]))
+   xn = (exprs(x)[pm, ] - exprs(x)[mm, ])/refSig
+   vsnres = vsn(xn, lts.quantile = 0.95, subsample = 2e+05,
+     verbose = FALSE)
+   yn = exprs(vsnres)/log(2)
+   throwOut = (refSig < quantile(refSig, probs = cutoffQuantile))
+   yn[throwOut, ] = NA
+   res = matrix(as.numeric(NA), nrow = nrow(x), ncol = ncol(yn))
+   res[pm, ] = yn
+   return(res)
+ }
> xwmm = cache("xwmm", normalizeByReferenceWithMM(davidTiling[,
+   isRNA], davidTiling[, isDNA], PMind, MMind))
```

5 Assessment

5.1 Visually

We would like to visualize the data along genomic coordinates. We select the features that map to the “-” strand of chromosome 9. The integer vectors `sta` and `end` contain the start and end coordinate of their match, `ind` their indices in the array `exprs(davidTiling)`.

```
> sta = probeAnno$"9.-.start"
> end = probeAnno$"9.-.end"
> ind = probeAnno$"9.-.index"
```

We construct a list of vectors, each containing different versions of the intensity data, in order that corresponds to `sta` and `ind` from above.

```

> dat = vector(mode = "list", length = 5)
> dat[[1]] = log2(exprs(davidTiling)[ind, which(isDNA)[1]])
> dat[[2]] = log2(exprs(davidTiling)[ind, which(isRNA)[1]])
> dat[[3]] = dat[[2]] - dat[[1]]
> dat[[4]] = exprs(xn1)[ind, 1]
> dat[[5]] = exprs(xn2)[ind, 1]
> dat[[6]] = xwmm[ind, 1]
> for (j in 3:length(dat)) dat[[j]] = dat[[j]] - quantile(dat[[j]],
+ 0.05, na.rm = TRUE)
> names(dat) = letters[seq(along = dat)]

```

We select a 10kB region around the highly expressed genes RPN2 and SER33 to fit on a plot, and set the y -axis limits:

```

> sel = (sta >= 216600 & end <= 227000)
> ysc = sapply(dat, function(py) quantile(py, probs = c(0, 1),
+ na.rm = TRUE))
> ysc[, 3:6] = c(-3, 8)

```

Now we are ready to plot:

```

> anno = data.frame(start = c(217860, 221078), end = c(220697,
+ 222487), name = I(c("RPN2", "SER33")))
> ticks = c(217, 223, 224, 225, 226)
> comparisonPlot((sta + end)[sel]/2, lapply(dat, "[" , sel), yscale = ysc,
+ anno = anno, ticks = ticks, cex = 0.2)

```

The result is shown in Figure 1. It shows scatterplots of different types of signal (y -axis) along genomic coordinates (x -axis). Each dot corresponds to a microarray feature. Note how the signal distributions in panels b)-f) vary both with respect to the variability within transcribed and untranscribed segments, and in the amplitude between them. a) signal from one of the DNA hybridizations (logarithmic scale, base 2). The y -coordinate of each dot is also encoded using a pseudo-color scheme. Dark red corresponds to features that have a very weak response, dark blue to those with the strongest response. The same coloring is also used in panels b)-f). b) unnormalized intensities from one of the poly(A) RNA hybridizations (logarithmic scale, base 2). c) Divide RNA-signal by DNA-signal then take logarithm (base 2). d) Background subtraction of the RNA-signal, divide by DNA-signal, then variance stabilizing normalization (vsb, glog base 2). e) In addition to d), drop the 5% weakest features in the DNA hybridization. f) Similar to e), but using the MM probes for background correction instead.

Now here's a bit of a hack: the plot symbols are too big if the plot is produced as above, and somehow the PDF and EPS drivers of R ignore the `cex` parameter. However, one can open the EPS file and replace all occurrences of the string "3.00" by "1.00".

```

> writeLines(sub(" 3.00", " 1.50", readLines("assessNorm-alongChromPlot.eps")),
+ fixed = TRUE), con = ("assessNorm-tmp.eps"))
> system("ps2pdf assessNorm-tmp.eps assessNorm-alongChromPlot.pdf")

```

5.2 Quantitatively

In order to assess to quantitatively assess the results of normalization, we consider a signal/noise ratio. Note that only looking at one or another (signal, or noise) by itself could be misleading.

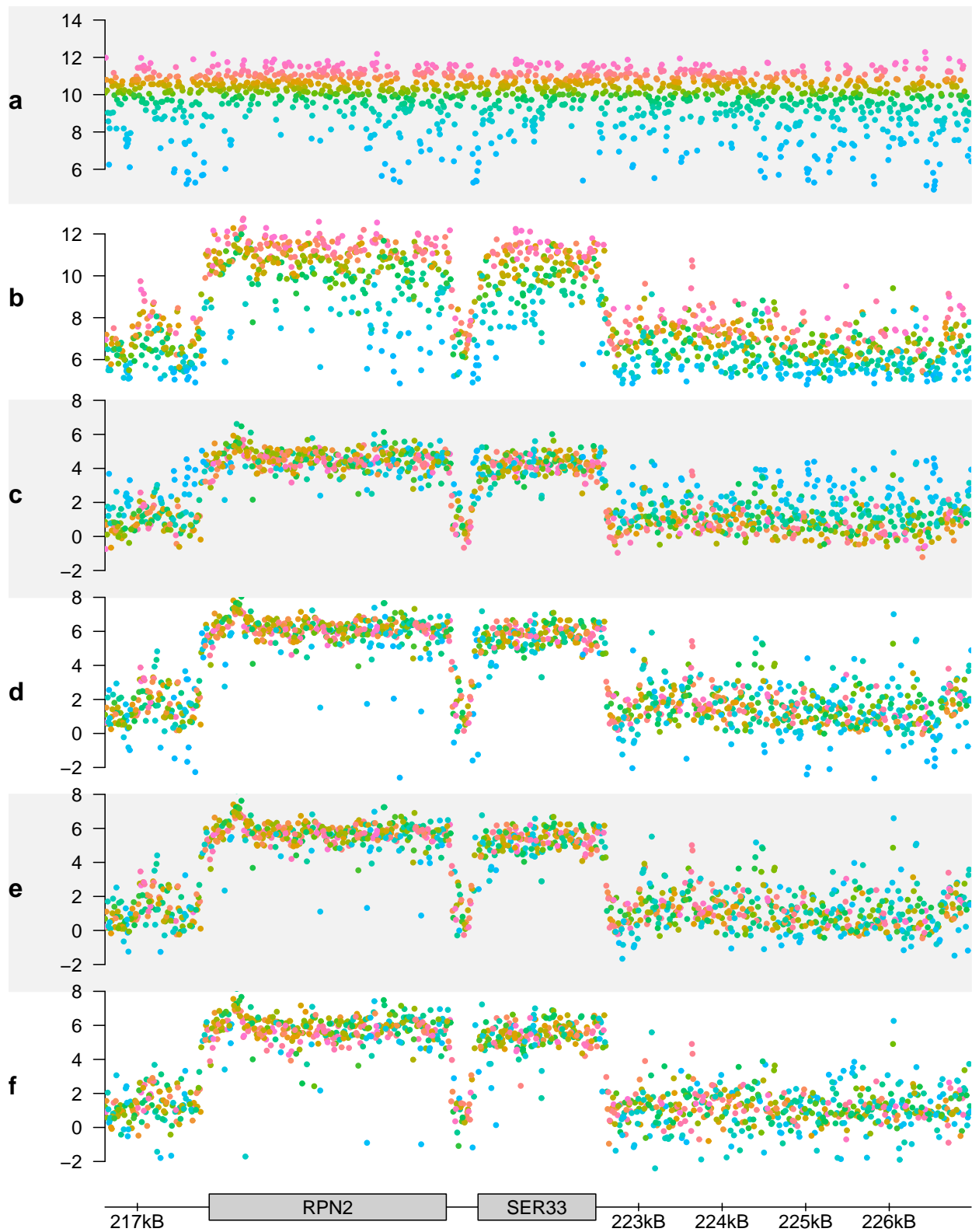


Figure 1: Along-chromosome plots of the data for different normalization methods. Please see text.

We define a set of control regions, which correspond to either highly expressed transcripts (**positiveCtrls**) or untranscribed intergenic regions (**negativeCtrls**). The assumption is that the signal within a region should be constant, and deviations from that are “noise”, while the difference between positive and negative controls should be large, and is counted as “signal”.

```
> positiveCtrls = cbind(c(217860, 220697), c(221078, 222487))
> negativeCtrls = cbind(c(216800, 217700), c(222800, 227000))
```

Noise σ is calculated as the average of the differences between 97.5% and 2.5% quantiles of the data within each of the control regions. The range between the 97.5% and 2.5% quantiles contains 95% of the data, while 5% is outside the range.

$$\sigma = \frac{1}{Q_N^{0.975} - Q_N^{0.025}} \cdot \frac{\sum_{r \in \{\text{pos}, \text{neg}\}} Q_r^{0.975} - Q_r^{0.025}}{|\{\text{pos}, \text{neg}\}|}. \quad (1)$$

Here, the symbol r counts over the different regions. The constant in the denominator is the differences between 95% and 5% quantiles for the standard Normal distribution, hence σ is equal to 1 if the data come from the standard Normal distribution.

Signal $\Delta\mu$ is calculated as the difference between the averages of the means of positive and negative control regions.

$$\Delta\mu = \frac{\sum_{r \in \{\text{pos}\}} \mu_r}{|\{\text{pos}\}|} - \frac{\sum_{r \in \{\text{neg}\}} \mu_r}{|\{\text{neg}\}|}. \quad (2)$$

I have explored many variations of this calculation, using different definitions of σ , $\Delta\mu$, and of the control regions. The ranking (relative order) of the methods was always the same as shown in the following.

```
> fac = 2 * qnorm(0.975)
> withinAndBetween = function(x, ...) {
+   meanAndSd = function(region, dohist = FALSE) {
+     d = x[(sta >= region[1]) & (end <= region[2]) & !is.na(x)]
+     res = c(mean(d), diff(quantile(d, c(0.025, 0.975)))/fac,
+       length(d))
+     if (dohist)
+       hist(d, 20, main = paste(signif(res, 3)), col = "orange")
+     res
+   }
+   p = apply(positiveCtrls, 2, meanAndSd, ...)
+   n = apply(negativeCtrls, 2, meanAndSd, ...)
+   dmu = sum(p[1, ] * p[3, ])/sum(p[3, ]) - sum(n[1, ] * n[3,
+     ])/sum(n[3, ])
+   sig = (sum(p[2, ] * p[3, ]) + sum(n[2, ] * n[3, ]))/(sum(p[3,
+     ] + sum(n[3, ]))
+   return(c("S/N" = dmu/sig, S = dmu, N = sig))
+ }
```

```
> sn = sapply(dat[2:6], withinAndBetween, dohist = TRUE)
```

	b	c	d	e	f
S/N	3.22	3.473	4.04	4.583	4.36
S	3.67	3.109	4.36	4.333	4.43
N	1.14	0.895	1.08	0.945	1.02

5.3 Dependency of the signal on GC content before and after normalization

We calculate the GC content of each probe,

```
> bc = basecontent(seqPM)
> nrGC = bc[, "C"] + bc[, "G"]
> rggc = c(4, 18)
> sel = (nrGC >= rggc[1] & nrGC <= rggc[2])
> nrGC = nrGC[sel]
> myPlot = function(x, tit) {
+   maxval = quantile(x, probs = 0.9999, na.rm = TRUE)
+   minval = maxval - 10.5
+   mycol = colorRampPalette(brewer.pal(9, "GnBu"))(diff(rggc) +
+     1)
+   boxplot(x[sel] ~ nrGC, ylim = c(minval, maxval), main = "",
+     col = mycol, outline = FALSE)
+   text(rggc[1] - 5, maxval + 1, tit, pos = 2, font = 2, xpd = NA)
+ }
```

Remember that we defined `seqPM=Scerevisiaetilingprobe$sequence[ipm]` in Section 4.2.

```
> par(mfrow = c(1, 2))
> dat1 = log2(exprs(davidTiling)[ipm, which(isRNA)[1]])
> dat2 = exprs(xn2)[ipm, 1] - quantile(exprs(xn2)[ipm, 1], 0.05,
+   na.rm = TRUE)
> myPlot(dat1, "a")
> myPlot(dat2, "b")
```

The result is shown in Figure 2.

References

- [1] Lior David, Wolfgang Huber, Marina Granovskaia, Joern Toedling, Curtis J. Palm, Lee Bofkin, Ted Jones, Ronald W. Davis, and Lars M. Steinmetz A high-resolution map of transcription in the yeast genome. *PNAS*, 2006. 1

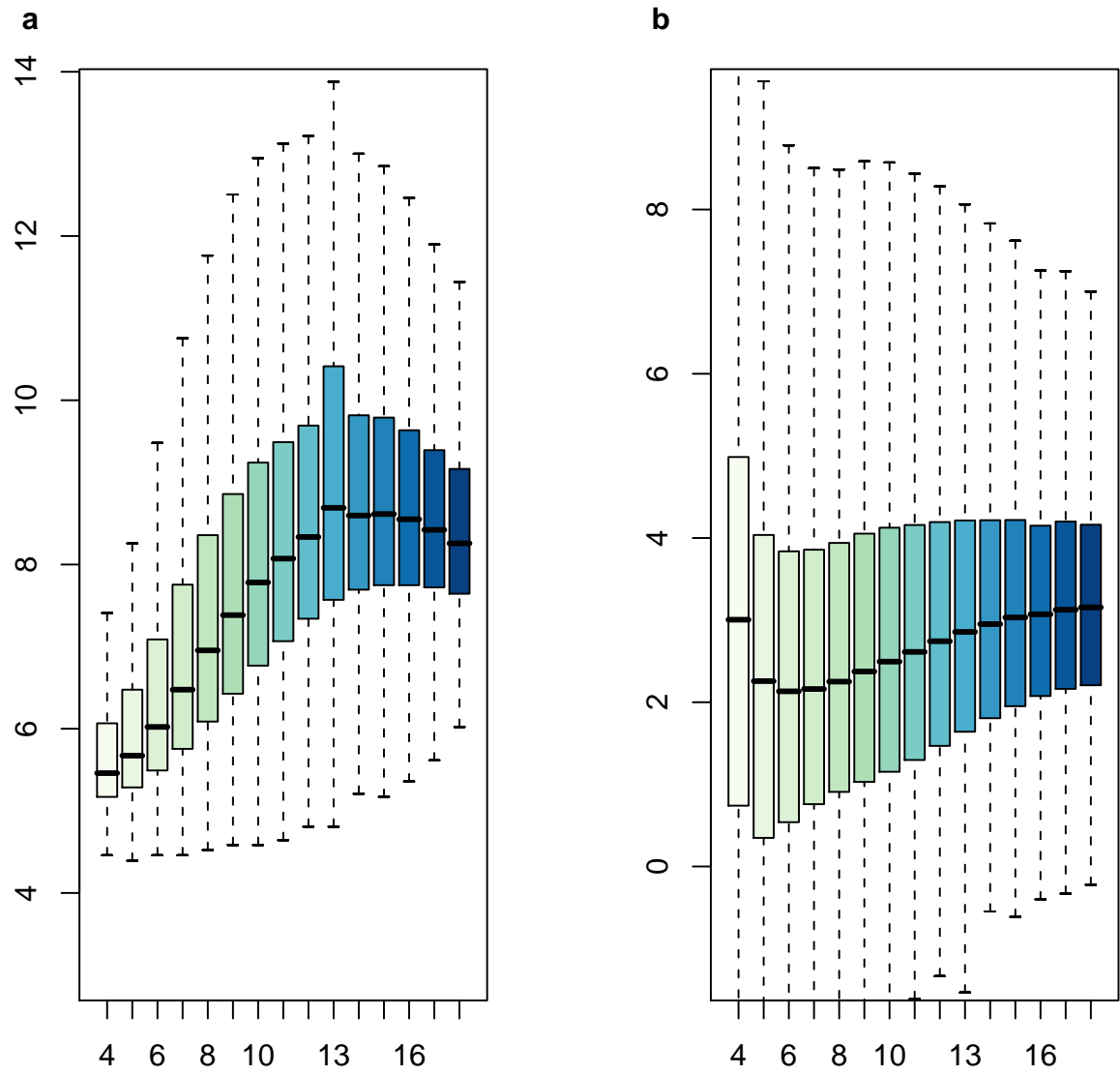


Figure 2: ECDF of \log_2 intensities (a) and normalized intensities (b) stratified by GC content.