

# APPENDIX: NORMALIZATION OF TWO-CHANNEL MICROARRAYS ACCOUNTING FOR EXPERIMENTAL DESIGN AND INTENSITY-DEPENDENT RELATIONSHIPS

Alan R. Dabney

Department of Statistics

Texas A&M University, College Station, TX 77843

adabney@stat.tamu.edu

John D. Storey

Department of Biostatistics

Department of Genome Sciences

University of Washington, Seattle, WA 98195

jstorey@u.washington.edu

## Translating the eCADS model

While we do not know the  $x_{il}$ , we are able to estimate monotone-increasing functions of them, as follows. Let  $x_{il}^* = d(x_{il}) = (d(x_{il}) + \delta_1(x_{il}) + d(x_{il}) + \delta_2(x_{il}))/2$ . Now define functions  $d^*$ ,  $\delta_k^*$ , and  $a^*$

$$d^*(x_{il}^*) = x_{il}^* = d(x_{il})$$

$$\delta_k^*(x_{il}^*) = \delta_k(x_{il})$$

$$a_j^*(x_{il}^*) = a_j(x_{il}).$$

The new dye functions  $d^*(x_{il}^*) + \delta_k^*(x_{il}^*)$  are also monotone-increasing, and preserving this monotonicity is more important than estimating the true dye functions. These definitions also preserve the constraints  $\sum_{k=1}^2 \delta_k^*(x) = 0$  and  $\sum_{j=1}^n a_j^*(x) = 0$  for any argument  $x$ . We can therefore translate

the model in equation ((2)) from the main document into

$$\begin{aligned} y_{ijkl} &= d^*(x_{il}^*) + \delta_k^*(x_{il}^*) + a_j^*(x_{il}^*) + \epsilon_{ijkl} \\ &= x_{il}^* + \delta_k^*(x_{il}^*) + a_j^*(x_{il}^*) + \epsilon_{ijkl}, \end{aligned} \tag{a}$$

where the second equality follows by the definition of  $d^*$ . We proceed with the two-stage approach: (i) estimate the  $x_{il}^*$ , (ii) fit model (a) using the  $\hat{x}_{il}^*$  as plug-in estimates of the  $x_{il}^*$ .

## Estimating the warped RNA amounts

In general, to estimate the  $x_{il}^*$ , we use fitted values from gene-specific ANOVA or regression models. One simple model choice is based on ANOVA model in equation ((1)) of the main document:

$$y_{ijkl} = \mu_i + d_k + t_{il} + \epsilon_{ijkl}. \tag{b}$$

The spot effects  $a_{ij}$  are not necessary for our purposes here and are excluded from the model. Additional covariates or sources of bias could be added to the gene-specific model as necessary. Recall that  $x_{il}^* = d(x_{il})$  in the eCADS model corresponds to  $\mu_i + t_{il}$  in the ANOVA model. Hence, we estimate the  $x_{il}^*$  with the fitted values  $\hat{\mu}_i + \hat{t}_{il}$  from model (b). The estimated ‘‘warped RNA amounts’’ can thus be thought of as group-specific means, adjusted gene-by-gene for bias.

We note that in our previous work on the CADS method (Dabney & Storey 2007), we defined all intensity-dependent functions to be in terms of *subject-specific* RNA amounts,  $x_{ijl}$ . Since CADS targets technical dye-swap experiments, where comparisons are made between paired arrays, it is natural to think of its arguments as being subject-specific. In the present setting, we could also use subject-specific RNA amounts. This requires some further assumptions. However, importantly, the choice between population- or subject-specific RNA amounts does not affect the targets of interest,  $d(x_{il})$ , and hence does not affect the operating characteristics stated here for eCADS. For details, see Dabney (2006).

## Basis matrix representation of the fANOVA model

Suppose we can write the component functions of (a) in terms of basis matrices:

$$\begin{aligned} \delta_k^*(x_{il}^*) &= \mathbf{B}_{\delta il} \beta_{\delta k} \\ a_j^*(x_{il}^*) &= \mathbf{B}_{a il} \beta_{a j}. \end{aligned}$$

Here, for example,  $\mathbf{B}_{\delta il} = [b_{\delta 1}(x_{il}^*) \ b_{\delta 2}(x_{il}^*) \ \dots \ b_{\delta q_\delta}(x_{il}^*)]$  is the  $q_\delta$ -dimensional basis component (a  $q_\delta$ -dimensional row-vector) for the dye functions at the  $i$ th gene, evaluated at RNA amounts specific to group  $l$ . We can similarly define  $\delta_k^*(\mathbf{x}_l^*)$  to be the vector-valued function with  $i$ th component equal to  $\delta_k^*(x_{il}^*)$ , etc. This allows us to write

$$\begin{aligned}\delta_k^*(\mathbf{x}_l^*) &= \mathbf{B}_{\delta l} \beta_{\delta k} \\ \mathbf{a}_j^*(\mathbf{x}_l^*) &= \mathbf{B}_{a l} \beta_{a_j},\end{aligned}$$

where, for example,  $\mathbf{B}_{\delta l} = [\mathbf{B}_{\delta 1 l}^T \ \mathbf{B}_{\delta 2 l}^T \ \dots \ \mathbf{B}_{\delta m l}^T]^T$  is now a  $m \times q_\delta$  matrix.

Now group the  $y_{ijkl}$  into single-channel chunks,  $\mathbf{y}_h$ ,  $h = 1, 2, \dots, 2n$ . Let  $\mathbf{Z}$  be the model matrix describing the experiment, as defined in the section *Parameterizing the model*. Specifically, let  $\mathbf{Z}$  be the  $2n \times (n + p + 2)$  matrix with  $h$ th row equal to  $[z_{g_1 h} \ \dots \ z_{g_p h} \ z_{d_1 h} \ z_{d_2 h} \ z_{a_1 h} \ \dots \ z_{a_n h}]$ . Here,  $z_{g_l h}$  is a scalar indicator of whether channel  $h$  comes from group  $l$  (1 for yes, 0 for no),  $l = 1, 2, \dots, p$ . Similarly,  $z_{d_1 h}$  and  $z_{d_2 h}$  indicate whether channel  $h$  was labeled with the red or green dyes, respectively, and  $z_{a_j h}$  indicates whether the  $h$ th channel profile comes from array  $j$ ,  $j = 1, 2, \dots, n$ . Define

$$\begin{aligned}\mathbf{G}_{xh} &= \sum_{l=1}^p \mathbf{x}_l^* z_{g_l h} \\ \mathbf{G}_{\delta_k h} &= \left[ \sum_{l=1}^p \mathbf{B}_{\delta l} z_{g_l h} \right] z_{d_k h} \\ \mathbf{G}_{a_j h} &= \left[ \sum_{l=1}^p \mathbf{B}_{a l} z_{g_l h} \right] z_{a_j h},\end{aligned}$$

$k = 1, 2$ ,  $j = 1, 2, \dots, n$ ,  $h = 1, 2, \dots, 2n$ . We can then rewrite (a) as

$$\begin{aligned}\mathbf{y}_h &= \mathbf{G}_{xh} + \sum_{k=1}^2 \mathbf{G}_{\delta_k h} \beta_{\delta_k} + \sum_{j=1}^n \mathbf{G}_{a_j h} \beta_{a_j} + \epsilon_h \\ &= \mathbf{G}_{xh} + \mathbf{H}_h \boldsymbol{\theta} + \epsilon_h,\end{aligned}$$

$h = 1, 2, \dots, 2n$ , where  $\mathbf{H}_h = [\mathbf{G}_{\delta_1 h} \ \mathbf{G}_{\delta_2 h} \ \mathbf{G}_{a_1 h} \ \dots \ \mathbf{G}_{a_n h}]$  and  $\boldsymbol{\theta} = (\beta_{\delta_1}^T, \beta_{\delta_2}^T, \beta_{a_1}^T, \dots, \beta_{a_n}^T)^T$ . Finally, we can stack the  $\mathbf{H}_h$  on top of each other to form  $\mathbf{H} = [\mathbf{G}_{\delta_1} \ \mathbf{G}_{\delta_2} \ \mathbf{G}_{a_1} \ \dots \ \mathbf{G}_{a_n}]$ , where the  $\mathbf{G}$  components are the corresponding stacked versions of the  $\mathbf{G}_h$  defining  $\mathbf{H}_h$ . We similarly

define  $\mathbf{G}_x$  as the stacked version of the  $\mathbf{G}_{xh}$ . This gives

$$\begin{aligned}\mathbf{y} &= \mathbf{G}_x + \sum_{k=1}^2 \mathbf{G}_{\delta_k} \beta_{\delta_k} + \sum_{j=1}^n \mathbf{G}_{a_j} \beta_{a_j} + \boldsymbol{\epsilon} \\ &= \mathbf{G}_x + \mathbf{H} \boldsymbol{\theta} + \boldsymbol{\epsilon}.\end{aligned}\tag{c}$$

## Least-squares solutions

Our estimation goal is to minimize the least-squares criterion

$$(\mathbf{y} - \mathbf{G}_x - \mathbf{H} \boldsymbol{\theta})^T (\mathbf{y} - \mathbf{G}_x - \mathbf{H} \boldsymbol{\theta})$$

with respect to  $\boldsymbol{\theta}$  subject to the constraints  $\sum_{k=1}^2 \beta_{\delta_k} = 0$  and  $\sum_{j=1}^n \beta_{a_j} = 0$ . This can be done by replacing  $\beta_{\delta_2}$  with  $-\beta_{\delta_1}$  and  $\beta_{a_n}$  with  $-\sum_{j=1}^{n-1} \beta_{a_j}$  in equation (c), leading to

$$\mathbf{y} = \mathbf{G}_x + \mathbf{H}_0 \boldsymbol{\theta}_0 + \boldsymbol{\epsilon}.$$

Here,  $\mathbf{H}_0 = [\mathbf{G}_{\delta_1} - \mathbf{G}_{\delta_2} \quad \mathbf{G}_{a_1} - \mathbf{G}_{a_n} \quad \dots \quad \mathbf{G}_{a_{n-1}} - \mathbf{G}_{a_n}]$  and  $\boldsymbol{\theta}_0 = (\beta_{\delta_1}^T, \beta_{a_1}^T, \dots, \beta_{a_{n-1}}^T)^T$ . Standard least-squares theory leads to the estimate

$$\hat{\boldsymbol{\theta}}_0 = (\mathbf{H}_0^T \mathbf{H}_0)^{-1} \mathbf{H}_0^T (\mathbf{y} - \mathbf{G}_x).$$

## eCADS preserves differential expression relationships

eCADS can be applied to any valid experimental design. For optimal results, however, the design should be balanced with respect to comparison group. When this balance holds, the gene-specific sample averages of the eCADS-normalized data equal the quantities of interest ( $d(x_{il})$ ) in expectation. Thus, since  $d(x_{il})$  is a monotone function of the true RNA amounts, the expected value of any test statistic based on sample averages that compares the expression level of a gene in different groups has the same sign as the true difference in RNA amount. This means that, in expectation, null genes will be called null, overexpressed genes called overexpressed, and underexpressed genes called underexpressed. Based on simulation work, minor imbalances do not affect this result. A detailed proof follows.

Let  $\mathbf{Z}$  be the model matrix describing the experiment. Denote the  $h$ th row of  $\mathbf{Z}$  by  $\mathbf{z}_h = [z_{gh} \ z_{dh} \ z_{ah} \ z_{ch} \ z_{bh}]$ ,  $h = 1, 2, \dots, 2n$ . Here,  $\mathbf{z}_{gh} = (z_{g_1h} \ z_{g_2h} \ \dots \ z_{g_ph})^T$ ,  $\mathbf{z}_{dh} = (z_{d_1h} \ z_{d_2h})^T$ , and  $\mathbf{z}_{ah} = (z_{a_1h} \ z_{a_2h} \ \dots \ z_{a_nh})^T$ . The component  $\mathbf{z}_{ch}$  represents factors of interest, such as confounders, while the component  $\mathbf{z}_{bh}$  represents biases additional to the dyes and arrays used. To be general,

suppose there are  $C$  additional factors of interest, where the  $r$ th factor of interest has  $p_{Cr}$  levels,  $r = 1, 2, \dots, C$ . We associate with the  $u$ th level of the  $r$ th factor of interest the function  $c_{ru}$  and assume that  $\sum_{u=1}^{p_{Cr}} c_{ru} = 0$ ,  $r = 1, 2, \dots, C$ . Similarly, suppose there are  $B$  additional bias terms, where the  $s$ th bias factor has  $p_{Bs}$  levels,  $s = 1, 2, \dots, B$ . We associate with the  $v$ th level of the  $s$ th additional bias factor the function  $b_{sv}$  and assume that  $\sum_{v=1}^{p_{Bs}} b_{sv} = 0$ ,  $s = 1, 2, \dots, B$ .

Let  $\tilde{y}_{il}$  be the sample average of the eCADS-normalized data for gene  $i$  in comparison group  $l$ . Since we are assuming that test statistics will be based on sample averages, and since the factors of interest are balanced with respect to comparison group, it is sufficient to consider the test statistic  $t_i = \tilde{y}_{il} - \tilde{y}_{il'}$  for comparing the expression level of gene  $i$  in groups  $l$  and  $l'$ . Let  $n_{\delta_k l}$  be the number of times group  $l$  is labeled with dye  $k$ , and let  $n_{a_j l}$  be an indicator of whether group  $l$  is on array  $j$ . Similarly, let  $n_{c_{ru} l}$  be the number of times group  $l$  takes level  $u$  of the  $r$ th additional factor of interest. Finally, let  $n_{b_{sv} l}$  be the number of times group  $l$  takes level  $v$  of the  $s$ th additional bias factor. We can write

$$\begin{aligned} \tilde{y}_{il} = & x_{il}^* + \frac{1}{n} \sum_{k=1}^2 n_{\delta_k l} (\delta_k^*(x_{il}^*) - \hat{\delta}_k^*(\hat{x}_{il}^*)) + \frac{1}{n} \sum_{j=1}^n n_{a_j l} (a_j^*(x_{il}^*) - \hat{a}_j^*(\hat{x}_{il}^*)) \\ & + \frac{1}{n} \sum_{r=1}^C \sum_{u=1}^{p_{Cr}} n_{c_{ru} l} (c_{ru}^*(x_{il}^*) - \hat{c}_{ru}^*(\hat{x}_{il}^*)) + \frac{1}{n} \sum_{s=1}^B \sum_{v=1}^{p_{Bs}} n_{b_{sv} l} (b_{sv}^*(x_{il}^*) - \hat{b}_{sv}^*(\hat{x}_{il}^*)) + \bar{\epsilon}_{il}. \end{aligned}$$

If there is balance with respect to comparison group, then  $n_{\delta_1 l} = n_{\delta_2 l}$ ,  $n_{c_{11} l} = \dots = n_{c_{1p_{C1}} l}$ ,  $\dots$ ,  $n_{c_{C1} l} = \dots = n_{c_{Cp_{CC}} l}$ ,  $n_{b_{11} l} = \dots = n_{b_{1p_{B1}} l}$ ,  $\dots$ , and  $n_{b_{B1} l} = \dots = n_{b_{Bp_{BB}} l}$ . If direct comparisons between two groups were made, then each group appears on every array, and  $n_{a_1 l} = \dots = n_{a_n l}$ . For indirect comparisons or experiments with more than two groups, this will not be the case. However, in these cases, we can further assume that array functions sum to zero within each relevant subset of arrays. In practice, this assumption is not necessary, since array effects are expected to be minor in comparison to dye effects; any residual array effect will not change the monotonicity of  $d$ . Based on the sum-to-zero constraints on the model terms, we therefore have  $\tilde{y}_{il} = x_{il}^* + \bar{\epsilon}_{il} = d(x_{il}) + \bar{\epsilon}_{il}$ . A similar result holds for group  $l'$ , so that  $E(T_i) = d(x_{il}) - d(x_{il'})$ . The average dye function  $d$  is monotone increasing, and so

$$E(T_i) \begin{cases} = 0 & \text{if } x_{il} = x_{il'}, \\ > 0 & \text{if } x_{il} > x_{il'}, \\ < 0 & \text{if } x_{il} < x_{il'}. \end{cases}$$

## References

- Dabney, A. R. (2006). *The Normalization of Two-Channel Microarrays*, Ph.D. dissertation, University of Washington.
- Dabney, A. R. & Storey, J. D. (2007). A new approach to intensity-dependent normalization of two-channel microarrays, *Biostatistics* **8**: 128–139.