

Solving the Empirical Bayes Normal Means Problem with Correlated Noise

Lei Sun¹ and Matthew Stephens^{1, 2}

¹Department of Statistics, University of Chicago

²Department of Human Genetics, University of Chicago

Abstract

The Normal Means problem plays a fundamental role in many areas of modern high-dimensional statistics, both in theory and practice. And the Empirical Bayes (EB) approach to solving this problem has been shown to be highly effective, again both in theory and practice. However, almost all EB treatments of the Normal Means problem assume that the observations are independent. In practice correlations are ubiquitous in real-world applications, and these correlations can grossly distort EB estimates. Here, exploiting theory from Schwartzman (2010), we develop new EB methods for solving the Normal Means problem that take account of *unknown* correlations among observations. We provide practical software implementations of these methods, and illustrate them in the context of large-scale multiple testing problems and False Discovery Rate (FDR) control. In realistic numerical experiments our methods compare favorably with other commonly-used multiple testing methods.

1 Introduction

We consider the Empirical Bayes (EB) approach to the Normal Means problem (Efron and Morris, 1973; Johnstone and Silverman, 2004):

$$X_j \mid \theta_j, s_j \sim N(\theta_j, s_j^2), \quad j = 1, \dots, p. \quad (1)$$

Here $N(\mu, \sigma^2)$ denotes the normal distribution with mean μ and variance σ^2 ; $X := (X_1, \dots, X_p)$ are observations; $s := (s_1, \dots, s_p)$ are standard deviations that are assumed known; and $\theta := (\theta_1, \dots, \theta_p)$ are unknown means to be estimated. The EB approach assumes that θ_j are independent and identically distributed (iid) from some “prior” distribution,

$$\theta_j \stackrel{\text{iid}}{\sim} g(\cdot), \quad j = 1, \dots, p; \quad (2)$$

and performs inference for θ_j in two steps: first obtain an estimate of g , \hat{g} say, and second compute the posterior distributions $p(\theta_j \mid X_j, s_j, \hat{g})$. We refer to the two-step process as “solving the Empirical Bayes Normal Means (EBNM) problem.” The first step, estimating g , is sometimes of direct interest in itself, and is an example of a “deconvolution” problem (e.g. Kiefer and Wolfowitz, 1956; Laird, 1978; Stefanski and Carroll, 1990; Fan, 1991; Cordy and Thomas, 1997; Bovy et al., 2011; Efron, 2016).

First named by Robbins (1956), EB methods have seen extensive theoretical study (e.g. Robbins, 1964; Morris, 1983; Efron, 1996; Jiang and Zhang, 2009; Brown and Greenshtein, 2009; Scott and Berger, 2010; Petrone et al., 2014; Rousseau and Szabo, 2017; Efron, 2018), and are becoming widely used in practice. Indeed, according to Efron and Hastie (2016), “large parallel data sets are a hallmark of twenty-first-century scientific investigation, promoting the popularity of empirical Bayes methods.”

The EB approach provides a particularly attractive solution to the Normal Means problem. For example, the posterior means of θ provide shrinkage point estimates, with all the accompanying risk-reduction benefits (Efron and Morris, 1972; Berger, 1985). And the posterior distributions for θ provide corresponding

“shrinkage” interval estimates, which can have good coverage properties even “post-selection” (Dawid, 1994; Stephens, 2017). Further, by estimating g , EB methods “borrow strength” across observations, and automatically determine an appropriate amount of shrinkage from the data (Johnstone and Silverman, 2004). Because of these benefits, methods for solving the EBNM problem – and related extensions – are increasingly used in data applications (e.g. Clyde and George, 2000; Johnstone and Silverman, 2005b; Brown, 2008; Koenker and Mizera, 2014; Xing and Stephens, 2016; Urbut et al., 2018; Wang and Stephens, 2018; Dey and Stephens, 2018). One application of EBNM methods that we pay particular attention to later is large-scale multiple testing, and estimation/control of the False Discovery Rate (FDR; Benjamini and Hochberg, 1995; Efron, 2010b; Muralidharan, 2010; Stephens, 2017; Gerard and Stephens, 2018).

Almost all existing treatments of the EBNM problem assume that the observations X in (1) are independent given θ, s . However, this assumption can be grossly violated in practice. Non-negligible correlations are common in real world data sets. Further, as we discuss later, EB approaches to the Normal Means problem are particularly vulnerable to being misled by pervasive correlation. Specifically, when the average strength of pairwise correlations among observations is non-negligible, the estimate \hat{g} of g can be very inaccurate, and this adversely affects inference for *all* θ . Ironically then, the attractive “borrowing strength” property of the EB approach becomes, in the presence of pervasive correlation, its Achilles’ heel.

In this paper we introduce methods for solving the EBNM problem *allowing for unknown correlations* among the observations. More precisely, rewriting (1) as

$$\begin{aligned} X_j &= \theta_j + s_j Z_j \\ Z_j &\sim N(0, 1) \end{aligned} \tag{3}$$

we develop methods that allow for unknown correlations among the “noise” $Z := (Z_1, \dots, Z_p)$. Our methods are built on elegant theory from Schwartzman (2010), who shows, in essence, that the limiting empirical distribution, f say, of correlated $N(0, 1)$ random variables can be represented using a basis of the standard Gaussian density and its derivatives of increasing order. We use this result, combined with an assumption that Z are exchangeable, to frame solving this “EBNM with correlated noise” problem as a two-step process: first *jointly estimate f and g* from all observations; and second compute the posterior distribution of θ_j given the estimated \hat{f}, \hat{g} (and X_j, s_j). Although many possible assumptions on g are possible, here we assume g to be a scale mixture of zero-mean Gaussians, following the flexible “adaptive shrinkage” approach in Stephens (2017). We have implemented these methods in an R package, **cashr** (“correlated adaptive shrinkage in R”), available from <https://github.com/LSun/cashr>.

The rest of the paper is organized as follows. In Section 2, we illustrate how correlation can derail existing EBNM methods, and review Schwartzman (2010)’s representation of the empirical distribution of correlated $N(0, 1)$ random variables. In Section 3 we introduce the exchangeable correlated noise (ECN) model, and describe methods to solve the EBNM with correlated noise problem. Section 4 provides numeric examples with realistic simulations and real data illustrations. Section 5 concludes and discusses future research directions.

2 Motivation and Background

2.1 Correlation distorts empirical distribution and misleads EBNM methods

In essence, the reason correlation can cause problems for EBNM methods is that, even with large samples, the empirical distribution of correlated variables can be quite different from their marginal distribution (e.g. Efron, 2007a). To illustrate this, we generated realistic correlated $N(0, 1)$ z -scores using a framework similar to Gerard and Stephens (2017, 2018); Lu (2018). Specifically, we took RNA-seq gene expression data on the 10^4 most highly expressed genes in 119 human liver tissues (The GTEx Consortium, 2015, 2017). In each simulation we randomly drew two groups of five samples (without replacement), and applied a standard RNA-seq analysis pipeline, using the software packages **edgeR** (Robinson et al., 2010), **voom** (Law et al., 2014), and **limma** (Ritchie et al., 2015), to compute, for each gene $j = 1, \dots, 10^4$, an estimate of the \log_2 -fold

difference in mean expression, X_j , and a corresponding p -value, p_j , testing the null hypothesis that the difference in mean is 0. We converted p_j, X_j to a z -score $z_j := -\text{sign}(X_j)\Phi^{-1}(p/2)$, where Φ is the CDF of $N(0, 1)$. We also computed an “effective” standard deviation $s_j := X_j/z_j$ for later use (Figure 2 and Section 4).

In these simulations, because samples are randomly assigned to the two groups, there are no genuine differences in mean expression. Therefore the z -scores should have marginal distribution $N(0, 1)$. And, indeed, empirical checks confirm that the analysis pipeline produces well-calibrated marginally $N(0, 1)$ z -scores (Appendix A). However, the 10^4 z scores in each simulated data set are correlated, due to correlations among genes, and such correlations can distort the empirical distribution so that it is very different from $N(0, 1)$ (Efron, 2007a, 2010a,b). Figure 1 shows four examples, which were chosen to highlight some common patterns. Panels (a-c) all exhibit a feature we call *pseudo-inflation*, where the empirical distribution is *more* dispersed than $N(0, 1)$. Conversely, panel (d) exhibits *pseudo-deflation*, where the empirical distribution is *less* dispersed than $N(0, 1)$. Panel (b) also exhibits skew.

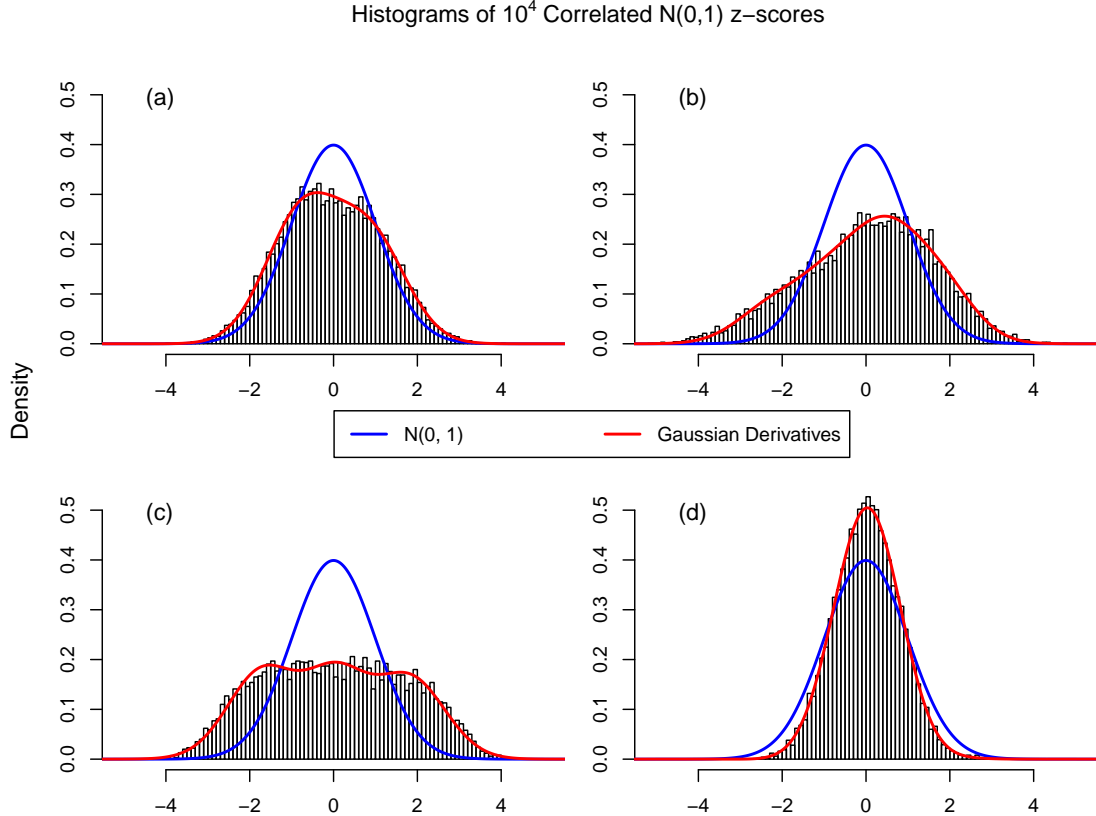


Figure 1: Illustration that the empirical distribution of a large number of correlated and marginally $N(0, 1)$ null z -scores can deviate substantially from $N(0, 1)$. The red lines are fitted densities obtained using our “Exchangeable Correlated Noise” model (Section 3.1) which uses a linear combination of the standard Gaussian density and its standardized derivatives.

Such *correlation-induced distortion* of the empirical distribution, if not appropriately addressed, can have serious consequence for EBNM methods. To illustrate this we applied several EBNM methods to five data sets simulated according to (2)-(3) as follows:

- The $p = 10^4$ normal means θ are iid samples from the mixture $g(\cdot) = 0.6\delta_0(\cdot) + 0.3N(\cdot; 0, 1) +$

$0.1N(\cdot; 0, 3^2)$, where $\delta_0(\cdot)$ denotes a point mass on 0 whose coefficient (0.6) is the null proportion, and $N(\cdot; \mu, \sigma^2)$ denotes the Gaussian density with mean μ and variance σ^2 . The same θ are used in all five data sets.

- In the first four data sets, the noise variables, Z , are the correlated null z -scores from the four panels of Figure 1. In the fifth data set Z are iid $N(0, 1)$ samples.
- The standard deviations s are simulated from the GTEx data as described above, and s are scaled to have $\frac{1}{p} \sum_j s_j^2 = 1$.

We provide the simulated X, s values to four existing EBNM methods – **EbayesThresh** (Johnstone and Silverman, 2004, 2005a), **REBayes** (Koenker and Mizera, 2014; Koenker and Gu, 2017), **ashr** (Stephens, 2017), and **deconvolveR** (Efron, 2016; Narasimhan and Efron, 2016) – that all ignore correlation and assume independence among observations. (For **deconvolveR** we set $s_j \equiv 1$ as its current implementation supports only homoskedastic noise.)

The estimates of g obtained by each method are shown in Figure 2. All methods do reasonably well in the fifth data set where Z are indeed independent (panel (e)). However, in the correlated data sets (panels (a-d)) the methods all misbehave in a similar way: over-estimating the dispersion of g under pseudo-inflation, and under-estimating it under pseudo-deflation. Their estimates of the null proportion are particularly inaccurate. These adverse effects are visible even when the distortion appears not severe (e.g. Figure 1(a)).

As a taster for what is to come, Figure 2 also shows the results for our new method, **cashr**, described later. This new method can account for both pseudo-inflation and pseudo-deflation, and in these examples estimates g consistently well.

2.2 Pseudo-inflation is non-Gaussian

In a series of pioneering papers (Efron, 2004, 2007a,b, 2008, 2010a), Efron studied the impact of correlations among z -scores on EB approaches to multiple testing. To account for the effects of correlation (pseudo-inflation, pseudo-deflation, and skew) on the empirical distribution of null z scores he introduced the concept of an “empirical null.” In his **locfdr** method (Efron, 2005), the empirical null is assumed to be Gaussian $N(\mu_0, \sigma_0^2)$. However, theory suggests that pseudo-inflation is not well modelled by a Gaussian distribution (Schwartzman, 2010, reviewed in Section 2.3), and a closer look at our empirical results here supports this conclusion.

To illustrate, Figure 3 shows more detailed analysis of the empirical distribution of Figure 1(c) z -scores. The central part of this z -score distribution could perhaps be modelled by a Gaussian distribution with inflated variance – for example, it matches more closely to a $N(0, 1.6^2)$ than to $N(0, 1)$. However, in the tails, the empirical distribution has much shorter tails than $N(0, 1.6^2)$. For example, 10^4 iid $N(0, 1.6^2)$ samples would be expected to yield 43 p -values exceeding the Bonferroni threshold of $0.05/10^4$, whereas in fact we observe none here. In short, the effects of pseudo-inflation are primarily in the “shoulders” of the distribution, where $|z|$ -scores are only moderately large, and not in the tails. (Incidentally, this behavior is far more evident in the histogram of z -scores than in the histogram of corresponding p -values, and we find the z -score histogram generally more helpful for diagnosing potential correlation-induced distortion.)

With hindsight this shoulder-but-not-tail inflation pattern should perhaps be expected. If one views the effect of correlation as to reduce the effective sample size, the number of extreme values of a sample with a smaller effective sample size should indeed be smaller. There are also relevant discussions on “asymptotic independence” in the extreme value theory (Sibuya, 1960; De Carvalho and Ramos, 2012). However, this property of pseudo-inflation does suggest that using a Gaussian to describe correlation-induced distortion, as in **locfdr**, is not ideal (more discussion in Section 4).

2.3 Empirical distribution of correlated $N(0, 1)$ random variables

We now summarize an elegant result of Schwartzman (2010), which characterizes the empirical distribution of a large number of correlated $N(0, 1)$ z -scores. This result plays a key role in our work.

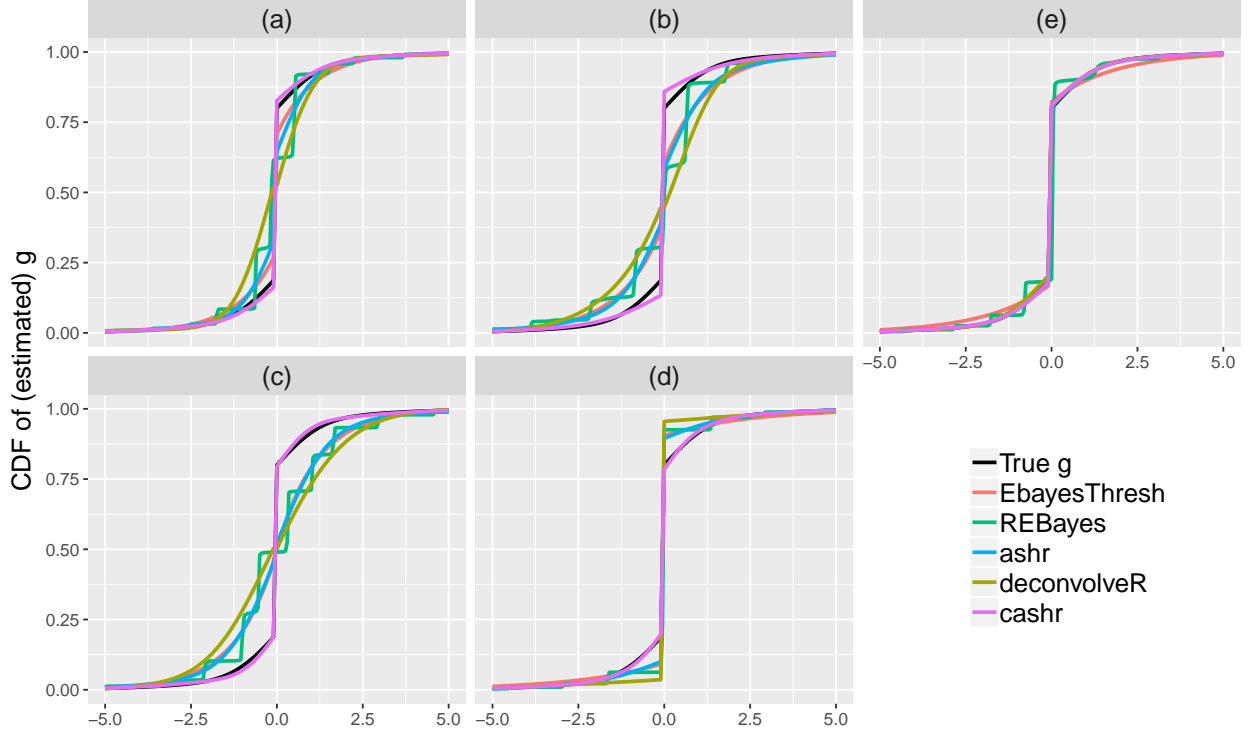


Figure 2: Illustration of how correlation can distort estimates of g obtained by EBNM methods. Each panel compares the true g with the estimated g from several EBNM methods applied to the same simulated dataset (see main text for details). In panels (a-d) Z are the correlated null z -scores from the corresponding panels of Figure 1. In panel (e) Z are iid $N(0, 1)$ samples. Existing EBNM methods (`EbayesThresh`, `REBayes`, `ashr`, `deconvolveR`), which ignore correlation among observations, do reasonably well with iid noise (e). However they do much less well in the correlated cases (a-d): over-estimating the dispersion of g under pseudo-inflation (a-c) and under-estimating it under pseudo-deflation (d). In contrast, our new method `cashr` (Section 3) estimates g consistently well.

On notation: let Φ and φ denote the CDF and PDF of $N(0, 1)$, and $\varphi^{(l)}$ denote the l^{th} derivative of φ . We refer to the collection of functions $\left\{ \frac{1}{\sqrt{l!}} \varphi^{(l)} \right\}_{l=1}^{\infty}$ as the (standardized) Gaussian derivatives. (Here “standardized” means that they are scaled to be orthonormal with respect to the weight function φ .)

Let $Z := \{Z_1, \dots, Z_p\}$ be p identically distributed, but not necessarily independent, $N(0, 1)$ random variables. Let F_p denote their empirical CDF:

$$F_p(\cdot) := \frac{1}{p} \sum_{j=1}^p \mathcal{I}(Z_j \leq \cdot), \quad (4)$$

where the indicator function $\mathcal{I}(Z_j \leq \cdot) := \begin{cases} 1 & Z_j \leq \cdot \\ 0 & Z_j > \cdot \end{cases}$. Since Z are random variables, F_p is a random function on $\mathbb{R} \rightarrow [0, 1]$. Also, because Z are marginally $N(0, 1)$, the expectation of F_p is Φ .

Schwartzman (2010) studies the distribution of F_p , and how its deviation from the expectation Φ depends on the correlations among Z . Specifically, assuming that each pair $\{Z_i, Z_j\}$ is bivariate normal with correlation ρ_{ij} (which is weaker than the common assumption that all $\{Z_1, \dots, Z_p\}$ are joint multivariate

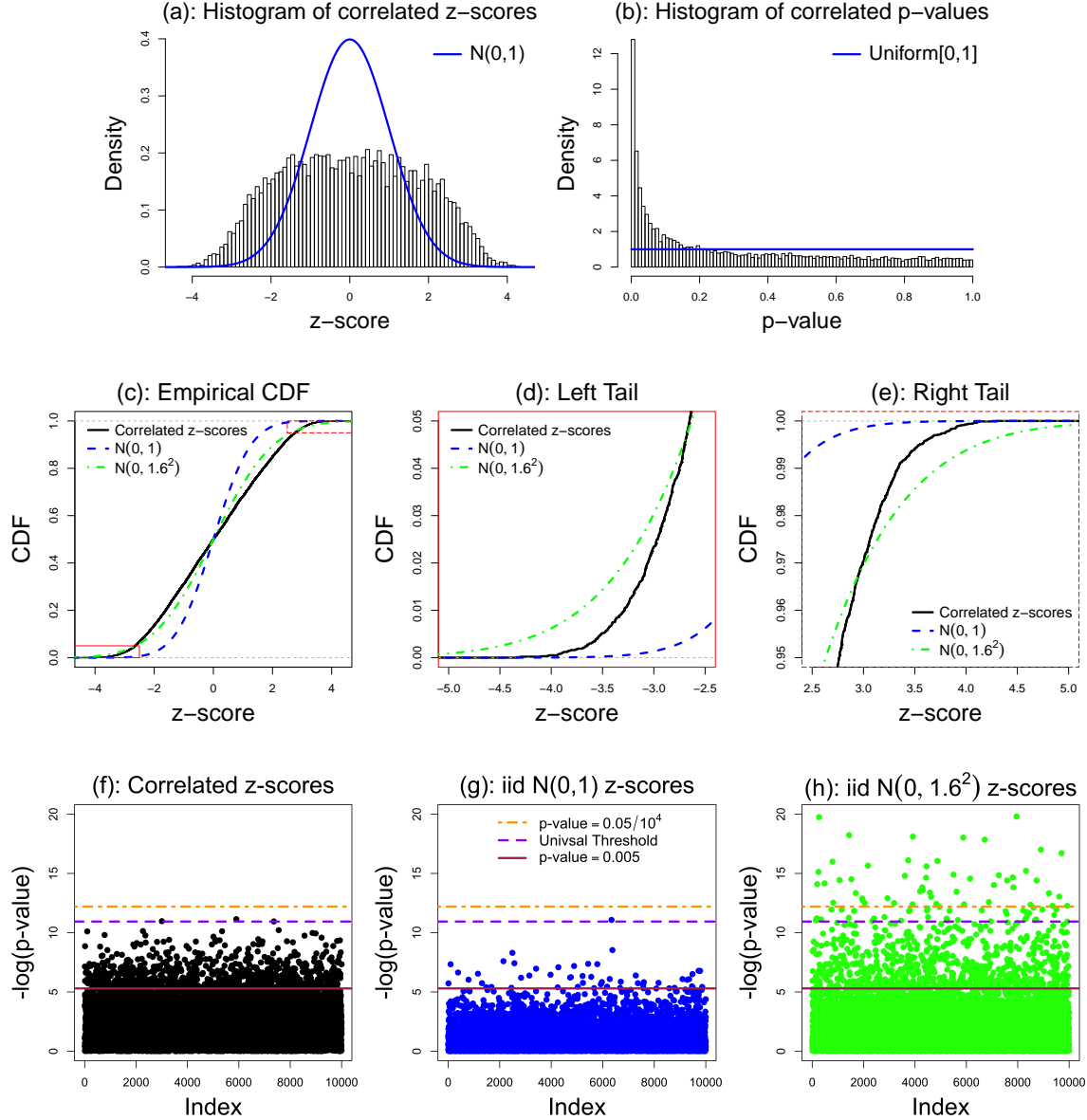


Figure 3: Illustration that the effects of pseudo-inflation are primarily in the “shoulders” of the distribution of null z -scores, and not in the tails. Panels (a-b): Histograms of correlated z -scores (from Figure 1(c)) and their corresponding p -values. Note that the “shoulder-but-not-tail” inflation is evident in the histogram of z -scores (a) but not in the oft-used histogram of p -values (b). Panels (c-e): Comparison of the empirical CDF of correlated z -scores with the CDF of $N(0, 1)$ and $N(0, 1.6^2)$. The z -score distribution is closer to $N(0, 1.6^2)$ in the center, but closer to $N(0, 1)$ in the tails. Panels (f-h): Comparison of correlated p -values with p -values obtained from 10^4 iid $N(0, 1)$ and $N(0, 1.6^2)$ z -scores. The number of correlated p -values ≤ 0.005 is closer to z scores from $N(0, 1.6^2)$, but the number in the extreme tail (e.g. clearing the Bonferroni or universal thresholds) is closer to $N(0, 1)$.

normal), Schwartzman (2010) provides the following representation of F_p when p is large:

$$F_p(\cdot) \approx F(\cdot) := \Phi(\cdot) + \sum_{l=1}^{\infty} W_l \frac{1}{\sqrt{l!}} \varphi^{(l-1)}(\cdot), \quad (5)$$

where W_1, W_2, \dots are uncorrelated random variables with $E[W_l] = 0$, and

$$\text{var}(W_l) = \overline{\rho^l} := \frac{1}{p(p-1)} \sum_{i,j:i \neq j} \rho_{ij}^l. \quad (6)$$

(Although uncorrelated, W_1, W_2, \dots are not independent; they must have higher-order dependence to guarantee that F is non-decreasing.)

Since F is a CDF, its derivative defines a corresponding density:

$$f(\cdot) := F'(\cdot) = \varphi(\cdot) + \sum_{l=1}^{\infty} W_l \frac{1}{\sqrt{l!}} \varphi^{(l)}(\cdot). \quad (7)$$

Intuitively, (7) characterizes how the (limiting) empirical distribution (histogram) of Z is likely to randomly deviate from the expectation φ , using standardized Gaussian derivatives as basis functions.

The representation (7) is crucial to our work here, and provides some remarkable insights. We highlight particularly the following:

1. The expected deviations of f from φ are determined by the variances of the coefficients W_l , which are determined by the mean and higher moments of the pairwise correlations, $\overline{\rho^l}$. In the special case where Z are independent all these terms are 0, and $f = \varphi$.
2. Following from 1, to create a tangible deviation from φ , $\overline{\rho^l}$ must be non-negligible (for some l). This requires *pervasive, but not necessarily strong*, pairwise correlations. For example, pervasive correlations occur if there is an underlying low-rank structure in the data, where all Z are influenced by a small number of underlying random factors, and so are all correlated with one another. In this case $\overline{\rho^l}$ will be non-negligible, and f may deviate from φ . In contrast, there can exist very strong pairwise correlations with negligible effect on f . For example, suppose p is even, and let Z be in $p/2$ pairs, with each pair having correlation one but different pairs being independent. The histogram of Z will look very much like $N(0, 1)$, because $\overline{\rho^l} = \frac{1}{p-1} \approx 0$ for large p . In other words, not all kinds of correlation, even large ones, distort the empirical distribution of Z .
3. Barring special cases such as $\rho_{ij} = 1$, the moments $\overline{\rho^l}$, and hence the expected magnitude of W_l , will decay quickly with l . Consequently the sum in (7) will typically be dominated by the first few terms, and the shape of the first few basis functions will determine the typical deviation of f from φ . Of the first four basis functions (Figure 4), the 2nd and 4th correspond to pseudo-inflation or pseudo-deflation in the shoulders of φ , depending on the signs of their coefficients, whereas the 1st and 3rd correspond to mean shift and skewness. This explains the empirical observation that correlation-induced pseudo-inflation tends to focus in the shoulders, and not the tails. (Also see Appendix B for the special case $\rho_{ij} = 1$.)

In discussing Efron (2010a), Schwartzman (2010) used this result to argue that “a wide unimodal histogram (of z -scores) may be indication of the presence of true signal, rather than an artifact of correlation.” Specifically, by discarding terms for $l \geq 4$ in (7), he found that the largest central spread (standard deviation) for f in (7) being a unimodal density is approximately 1.3. Along similar lines, we can show (Appendix B) that the maximum standard deviation for f being a Gaussian density is $\sqrt{2} \approx 1.4$. The key point here is that the effects of correlation are different from the effects of true signals, so the two can (often) be separated. Our methods here are designed to do exactly that.

3 Empirical Bayes Normal Means with Correlated Noise

3.1 The Exchangeable Correlated Noise model

As a first step towards allowing for correlated noise in the EBNM problem, we develop methods to fit the representation (7) to correlated null z -scores. We do this by treating Z as conditionally iid samples from f

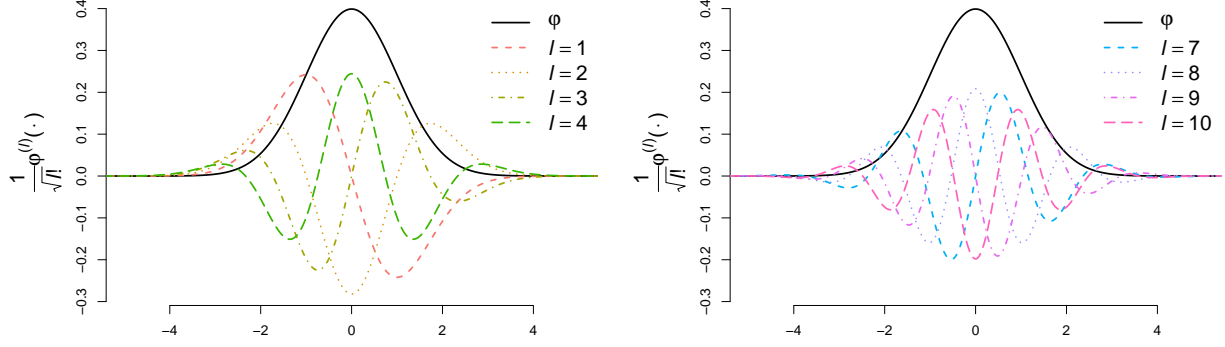


Figure 4: Illustration of the standard Gaussian density, ϕ , and its standardized derivatives. The left panel shows ϕ and its first four standardized derivatives. The 2nd and 4th derivatives correspond to pseudo-inflation or pseudo-deflation in the shoulders; the 1st and 3rd derivatives correspond to mean shift or skewness. The right panel shows ϕ and its 7th-10th derivatives. Even for these higher-order derivatives, tails are short, implying that correlation-induced distortion is unlikely to have long tails.

in (7), parameterized by $\omega := \{\omega_1, \omega_2, \dots\}$ which are realizations of $W := \{W_1, W_2, \dots\}$:

$$Z_j \mid \{W = \omega\} \stackrel{\text{iid}}{\sim} f(\cdot; \omega) := \phi(\cdot) + \sum_{l=1}^{\infty} \omega_l \frac{1}{\sqrt{l!}} \phi^{(l)}(\cdot). \quad (8)$$

It may seem perverse to model correlated random variables as conditionally iid. However, this treatment can be motivated by assuming Z are exchangeable and appealing to de Finetti's representation theorem (De Finetti, 1937), which says that (infinitely) exchangeable random variables can be represented as being conditionally iid from their empirical distribution. We therefore refer to the model (8) as the *exchangeable correlated noise (ECN)* model. We also refer to f as the *correlated noise distribution*.

To fit the ECN model (8) with observed Z , we estimate ω , essentially by maximum likelihood, but with a couple of additional complications that we now describe. First, since f is a density, we must constrain the parameters ω to ensure that $f(\cdot; \omega)$ is non-negative (note that (8) integrates to one for any ω , but is not guaranteed to be non-negative). Ideally f should be non-negative on the whole real line, but this constraint is difficult to work with, so we approximate it using a discrete approximation: we constrain $f(z_i; \omega) \geq 0$ on a fine grid $\{z_1, \dots, z_m\}$ such as $\{-10, -9.999, -9.998, \dots, +9.998, +9.999, +10\}$, in addition to $f(z_j; \omega) \geq 0$ for all j .

Second, to incorporate the prior expectation that the absolute value of ω_l should decay quickly with l (because $\text{var}(W_l) = \overline{\rho^l}$) we introduce a penalty on ω ,

$$h(\omega) := \sum_l \gamma_l |\omega_l|, \quad (9)$$

where we take the penalty parameters γ_l to be

$$\gamma_l = \begin{cases} 0 & l \text{ is odd} \\ \gamma / \rho^{l/2} & l \text{ is even} \end{cases}, \quad (10)$$

where γ represents a common penalty, and ρ some notion of average pairwise correlation. For computational convenience we use only the first $L = 10$ Gaussian derivatives (see Figure 4 for 7th-10th standardized Gaussian derivatives) and set $\omega_l = 0$ for $l > 10$. (Recall that $\text{var}(W_l) = \overline{\rho^l}$, so W_l 's realization ω_l will generally be negligible in practice for $l > 10$.) Of course a full Bayesian treatment would attempt to account for uncertainty in ω ; in ignoring that here we are making the usual EB compromise.

In numerical simulations we experimented with different combinations of $\gamma \in \{1, 5, 10, 50, 100\}$ and $\rho \in \{0.10, 0.25, 0.50, 0.75, 0.90\}$, and found that $\gamma = 10, \rho = 0.5$ performed well in a variety of situations, although results were not very sensitive to the choice of γ and ρ . All results in this paper were obtained with $\gamma = 10, \rho = 0.5$.

In summary, we estimate ω by solving:

$$\begin{aligned} \max_{\omega} \quad & \sum_j \log f(Z_j; \omega) - h(\omega) \\ \text{s.t.} \quad & f(Z_j; \omega) \geq 0, \quad j = 1, \dots, p, \\ & f(\mathbf{z}_i; \omega) \geq 0, \quad i = 1, \dots, m. \end{aligned} \quad (11)$$

This is a convex optimization and can be solved efficiently and stably using an interior point method; we implemented this using the **R** package **Rmosek** to interface to the MOSEK commercial solver (MOSEK ApS, 2018). With $p = 10^4$, the problem is solved on average within 0.50 seconds on a personal computer (Apple iMac, 3.2 GHz, Intel Core i5).

Figure 1 shows the fitted distributions from the ECN model, $\hat{f}(\cdot; \hat{\omega}) := \varphi(\cdot) + \sum_{l=1}^L \hat{\omega}_l \frac{1}{\sqrt{l!}} \varphi^{(l)}(\cdot)$, on the four illustrative sets of correlated null z -scores.

3.2 The EBNM model with correlated noise

To allow for correlated noise in the EBNM problem, we combine the standard EBNM model (2)-(3) with the ECN model (8):

$$X_j = \theta_j + s_j Z_j \quad (12)$$

$$\theta_j \sim g(\cdot) \quad (13)$$

$$Z_j \mid \omega \sim f(\cdot; \omega) = \varphi(\cdot) + \sum_{l=1}^L \omega_l \frac{1}{\sqrt{l!}} \varphi^{(l)}(\cdot). \quad (14)$$

Note that in this model the observations are conditionally independent given f and g .

Following Stephens (2017) we model the prior distribution g by a finite mixture of zero-mean Gaussians:

$$g(\cdot; \pi) = \pi_0 \delta_0(\cdot) + \sum_{k=1}^K \pi_k N(\cdot; 0, \sigma_k^2), \quad (15)$$

where π_0 is the null proportion. Here the mixture proportions $\pi := \{\pi_0, \pi_1, \dots, \pi_K\}$ are non-negative and sum to 1, and are to be estimated, whereas the component standard deviations $\sigma_1 < \sigma_2 < \dots < \sigma_K$ are a fixed pre-specified grid of values. By using a sufficiently wide and dense grid of standard deviations this finite mixture can approximate, to arbitrary accuracy, any scale mixture of zero-mean Gaussians.

The marginal log-likelihood for π, ω , integrating out θ, Z , is given by the following Theorem.

Theorem 1. *Combining (12)-(15), the marginal log-likelihood of π, ω is*

$$L(\pi, \omega) := \log \left(\prod_{j=1}^n p(X_j \mid \pi, \omega) \right) = \sum_{j=1}^n \log \left(\sum_{k=0}^K \pi_k \left(p_{jk0} + \sum_{l=1}^L \omega_l p_{jkl} \right) \right), \quad (16)$$

where

$$p_{jkl} = \frac{s_j^l}{\sqrt{\sigma_k^2 + s_j^2}^{l+1}} \frac{1}{\sqrt{l!}} \varphi^{(l)} \left(\frac{X_j}{\sqrt{\sigma_k^2 + s_j^2}} \right). \quad (17)$$

Proof. See Appendix C. □

3.3 Fitting the model

Following the usual EB approach, we fit the model (12)-(15) in two steps, first estimating g, f by estimating π, ω and then basing inference for θ on the (estimated) posterior distribution $p(\theta|X, s, \hat{\pi}, \hat{\omega})$. Note that under the model (12)-(15) $\theta_1, \dots, \theta_p$ are conditionally independent given f, g, X, s , so this posterior distribution $p(\theta|X, s, \hat{\pi}, \hat{\omega})$ factorizes, and is determined by its marginal distributions $p(\theta_j|X_j, s_j, \hat{\pi}, \hat{\omega})$. The intuition here is that, under the exchangeability assumption, the effects of correlation are captured entirely by the (realized) correlated noise distribution f . Once this distribution is estimated the inferences for each θ_j become independent, just as in the standard EBNM problem.

The usual EBNM approach to estimating π, ω would be to maximize the likelihood $L(\pi, \omega)$. Here we modify this approach using maximum penalized likelihood. Specifically we use the penalty on ω as in (9), and the penalty on π used by Stephens (2017) to encourage conservative (over-)estimation of the null proportion π_0 (to induce conservative estimation of false discovery rates). Thus, we solve

$$\hat{\pi}, \hat{\omega} = \arg \max_{\pi, \omega} \sum_{j=1}^n \log \left(\sum_{k=0}^K \pi_k \left(p_{jk0} + \sum_{l=1}^L \omega_l p_{jkl} \right) \right) + \sum_{k=0}^K \lambda_k \log(\pi_k) - \sum_{l=1}^L \gamma_l |\omega_l| \quad (18)$$

subject to the constraints

$$\sum_{k=0}^K \pi_k = 1 \quad (19)$$

$$\pi_k \geq 0, \quad k = 0, 1, \dots, K \quad (20)$$

$$\varphi(\mathbf{z}_i) + \sum_{l=1}^L \omega_l \frac{1}{\sqrt{l!}} \varphi^{(l)}(\mathbf{z}_i) \geq 0, \quad i = 1, \dots, m. \quad (21)$$

In (21) we used the same device as in (11) to capture non-negativity of f . We set γ_l as in (10), use only the first $L = 10$ Gaussian derivatives, and set $\lambda_0 = 10, \lambda_1 = \dots = \lambda_K = 0$ as in Stephens (2017).

Problem (18) is biconvex. That is, given a feasible $\hat{\pi}$, the optimization over ω is convex; and given a feasible $\hat{\omega}$, the optimization over π is convex. The optimization over π can be solved using the EM algorithm, or more efficiently using convex optimization methods (Koenker and Mizera, 2014; Koenker and Gu, 2017; Kim et al., 2018). To optimize over ω we use the same approach as in solving (11). To solve (18) we simply iterate between these two steps until convergence.

3.4 Posterior Calculations

For each j , the posterior distribution $p(\theta_j | X_j, \hat{\pi}, \hat{\omega})$ is, by Bayes Theorem, given by

$$p(\theta_j | X_j, \hat{\pi}, \hat{\omega}) = \frac{\left[\hat{\pi}_0 \delta_0 + \sum_{k=1}^K \hat{\pi}_k N(\theta_j; 0, \sigma_k^2) \right] \left[\frac{1}{s_j} \varphi \left(\frac{X_j - \theta_j}{s_j} \right) + \sum_{l=1}^L \hat{\omega}_l \frac{1}{s_j} \frac{1}{\sqrt{l!}} \varphi^{(l)} \left(\frac{X_j - \theta_j}{s_j} \right) \right]}{\sum_{k=0}^K \hat{\pi}_k \left(p_{jk0} + \sum_{l=1}^L \hat{\omega}_l p_{jkl} \right)}. \quad (22)$$

Despite the somewhat complex form, some important functionals of this posterior distribution are analytically available.

1. The posterior mean for θ_j

$$E[\theta_j | X_j, \hat{\pi}, \hat{\omega}] = \frac{\sum_{k=0}^K \hat{\pi}_k \left(m_{jk0} + \sum_{l=1}^L \hat{\omega}_l m_{jkl} \right)}{\sum_{k=0}^K \hat{\pi}_k \left(p_{jk0} + \sum_{l=1}^L \hat{\omega}_l p_{jkl} \right)}, \quad (23)$$

where $m_{jkl} = -\frac{s_j^l \sigma_k^2}{\sqrt{\sigma_k^2 + s_j^{2l+2}}} \frac{1}{\sqrt{l!}} \varphi^{(l+1)} \left(\frac{X_j}{\sqrt{\sigma_k^2 + s_j^2}} \right)$.

2. The local FDR (lfdr; Efron, 2008) is

$$\text{lfdr}_j := \Pr(\theta_j = 0 \mid X_j, \hat{\pi}, \hat{\omega}) = \frac{\hat{\pi}_0 \frac{1}{s_j} \varphi \left(\frac{X_j}{s_j} \right) + \sum_{l=1}^L \hat{\omega}_l \frac{1}{s_j} \frac{1}{\sqrt{l!}} \varphi^{(l)} \left(\frac{X_j}{s_j} \right)}{\sum_{k=0}^K \hat{\pi}_k \left(p_{jk0} + \sum_{l=1}^L \hat{\omega}_l p_{jkl} \right)}. \quad (24)$$

From this, the FDR of any discovery set $\Gamma \subseteq \{1, \dots, n\}$ can be estimated as

$$\widehat{\text{FDR}}(\Gamma) = \frac{1}{|\Gamma|} \sum_{j \in \Gamma} \text{lfdr}_j, \quad (25)$$

where $|\Gamma|$ denotes the number of elements in Γ . Storey's q -value (Storey, 2003) for each j is defined as

$$q_j := \widehat{\text{FDR}}(\{k : \text{lfdr}_k \leq \text{lfdr}_j\}). \quad (26)$$

3. Stephens (2017) introduced the term “local false sign rate (lfsr)” to refer to the probability of getting the sign of an effect wrong, as well as the false sign rate (FSR) and the s -value, analogous to the FDR and the q -value, respectively. Making statistical inference about the sign of a parameter, rather than solely focusing on whether the parameter being zero or not, was also discussed in Tukey (1991); Gelman et al. (2012). The value of lfsr_j is defined as

$$\text{lfsr}_j := \min\{\Pr(\theta_j \geq 0 \mid X_j, \hat{\pi}, \hat{\omega}), \Pr(\theta_j \leq 0 \mid X_j, \hat{\pi}, \hat{\omega})\}, \quad (27)$$

which is easily calculated from lfdr_j and

$$\Pr(\theta_j > 0 \mid X_j, \hat{\pi}, \hat{\omega}) = \frac{\sum_{k=1}^K \hat{\pi}_k \left(\hat{\tau}_{jk0} + \sum_{l=1}^L \hat{\omega}_l \tau_{jkl} \right)}{\sum_{k=0}^K \hat{\pi}_k \left(p_{jk0} + \sum_{l=1}^L \hat{\omega}_l p_{jkl} \right)}, \quad (28)$$

where $\tau_{jkl} = \frac{s_j^l}{\sqrt{l!} \sqrt{s_j^2 + \sigma_k^2}^{l+1}} \left(\sum_{m=0}^l \binom{l}{m} \left(\frac{\sigma_k}{s_j} \right)^m \varphi^{(m-1)} \left(\frac{X_j}{\sqrt{s_j^2 + \sigma_k^2}} \frac{\sigma_k}{s_j} \right) \varphi^{(l-m)} \left(\frac{X_j}{\sqrt{s_j^2 + \sigma_k^2}} \right) \right)$. The FSR and s -value are estimated and defined similarly to the FDR and q -value as

$$\widehat{\text{FSR}}(\Gamma) = \frac{1}{|\Gamma|} \sum_{j \in \Gamma} \text{lfsr}_j, \quad s_j := \widehat{\text{FSR}}(\{k : \text{lfsr}_k \leq \text{lfsr}_j\}). \quad (29)$$

3.5 Software

We implemented both the fitting procedure and posterior calculations in an R package `cashr` which is available at <https://github.com/LSun/cashr>. For $p = 10^4$, it takes on average about 6 seconds for model fitting and posterior calculations on a personal computer (Apple iMac, 3.2 GHz, Intel Core i5).

4 Numerical results

We now empirically assess the performance of `cashr` on both simulated and real data. We focus our assessments on the “multiple testing” setting where θ is sparse and the main goal is to identify “significant” non-zero elements θ_j . This problem can be tackled using EB methods (Thomas et al., 1985; Greenland and

Robins, 1991) and here we compare **cashr** with both **locfdr** (Efron, 2005), which attempts to capture effects of correlation through an empirical null strategy discussed in Section 2.2, and **ashr** (Stephens, 2017), which fits the same EBNM model as **cashr** but without allowing for correlation – i.e. **ashr** is equivalent to setting $f = \varphi$ in (14). Multiple testing can also be tackled by attempting to control the FDR in the frequentist sense, and so we also compare with the Benjamini-Hochberg procedure (BH; Benjamini and Hochberg, 1995) and **qvalue** (Storey, 2002, 2003). One advantage of the EBNM approach to multiple testing is that it can provide not only FDR assessments, but also point estimates and interval estimates for the effects θ_j (Stephens, 2017). However, to keep our comparisons simple we focus here only on FDR assessments.

4.1 Realistic simulation with gene expression data

We constructed synthetic data with realistic correlation structure using the simulation framework in Section 2.1. The data are simulated according to the EBNM with correlated noise model (2)-(3) as follows.

- The $p = 10^4$ normal means $\theta_1, \dots, \theta_p$ are iid samples from

$$g(\cdot) = \pi_0 \delta_0(\cdot) + (1 - \pi_0) g_1(\cdot), \quad (30)$$

for six choices of g_1 and three choices of $\pi_0 \in \{0.5, 0.9, 0.99\}$ (Figure 5). The density functions of these six choices of g_1 and other simulation details are in Appendix D.

- To make the correlation structure among noise realistic, in each simulation Z are simulated from real gene expression data as in Section 2.1.
- The standard deviations s are also simulated from real gene expression data using the same pipeline, and are scaled to have $\frac{1}{p} \sum s_j^2 = 1$.
- The observations are constructed as $X_j = \theta_j + s_j Z_j$, $j = 1, \dots, p$.

In each simulated data set, this framework generates p correlated observations X_j of respective normal means θ_j with corresponding standard deviations s_j . The data $\{(X_1, s_1), \dots, (X_p, s_p)\}$ are made available to each method, while the effects θ_j are withheld. The analysis goal is to identify which θ_j are significantly different from 0. We applied each method to formulate a discovery set at nominal FDR = 0.1, and calculated the empirical false discovery proportion (FDP) for each discovery set. We ran 1000 simulations for each g_1 , divided evenly among the three choices of π_0 .

Figure 5 compares the performance of each method in these simulations. Our first result is that, despite the presence of correlation, most of the methods control FDR in the usual frequentist sense under most scenarios: that is, the mean FDP is usually below the nominal level of 0.1. Indeed, BH is notable in never showing a mean FDP exceeding the nominal level, even though, as far as we are aware, no known theory guarantees this under the realistic patterns of correlation used here (Benjamini and Yekutieli (2001) gives relevant theoretical results under more restrictive assumptions on the correlation). The method most prone to lose control is **ashr**, but even its mean FDP is never above 0.2.

However, despite this frequentist control of FDR, for most methods the FDP for individual data sets can often lie far from the nominal level (see also Owen, 2005; Qiu et al., 2005; Blanchard and Roquain, 2009; Friguet et al., 2009, for example). Arguably, then, frequentist control of FDR is insufficient in practice, since we desire – as far as is possible – to make sound statistical inference for each data set. That is, we might consider a method to perform well if its FDP is consistently close to the nominal level, rather than close on average. By this criterion, **cashr** consistently outperforms other methods (Figure 5): it provides uniformly lower root MSE of FDP from the nominal FDR, 0.1, and the whiskers in the boxplots (indicating 5th and 95th percentiles) are narrower. Along with FDP, Figure 5 also shows the empirical true discovery proportion (TDP), defined as the proportion of true discoveries out of the number of all non-zero θ_j , as an indication of statistical power. **cashr** maintains good power in that it produces higher TDP than most methods in most scenarios. In some scenarios, **ashr** sometimes finds more true signals than **cashr**, but at the cost of severely losing control of FDP.

Nominal FDR = 0.1

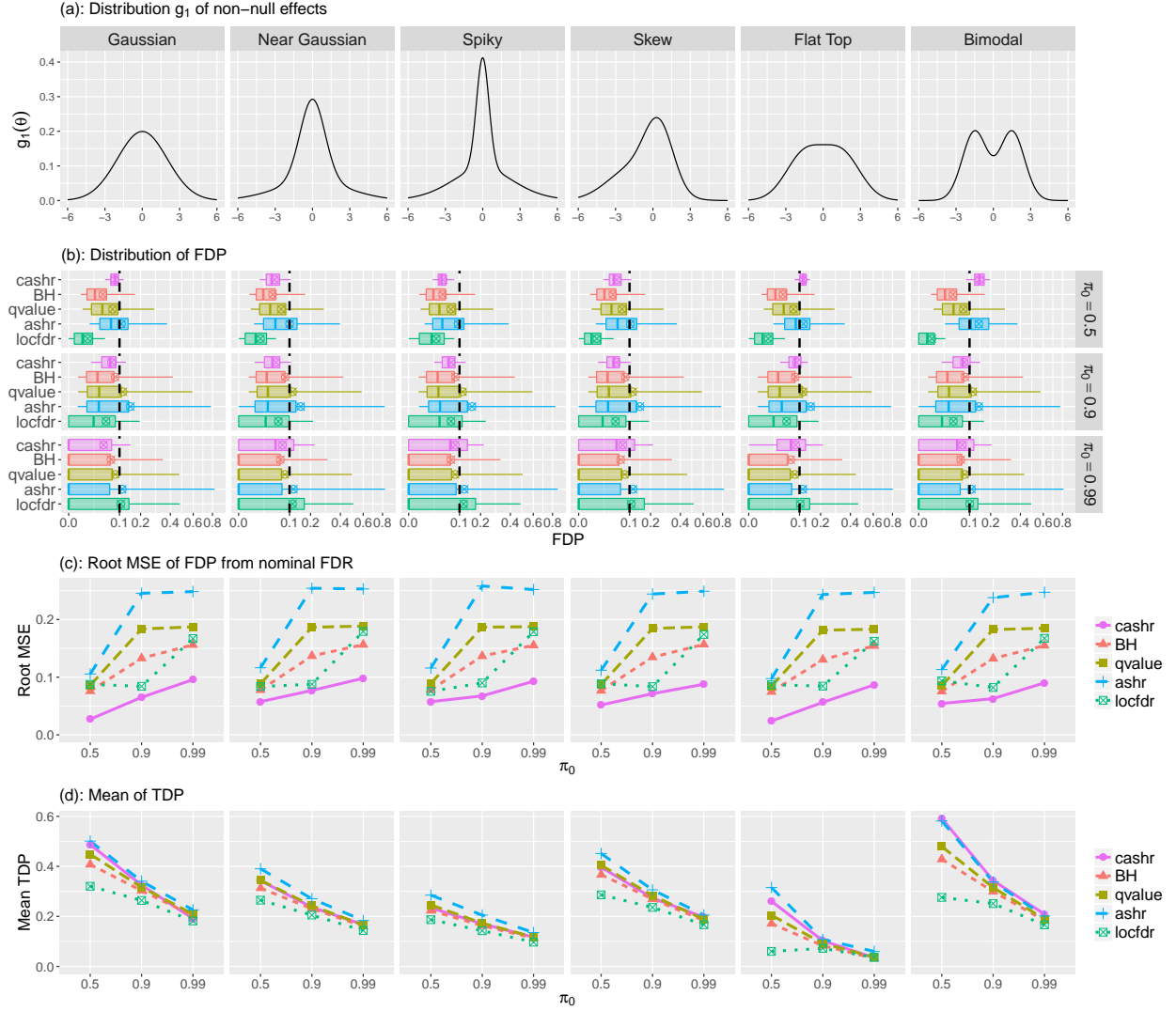


Figure 5: Illustration that **cashr** outperforms other methods in producing discovery sets whose FDP are consistently close to the nominal FDR, while maintaining good statistical power. Simulation results are shown for six different distributions for the non-null effect (g_1 ; panel (a)) and three different values of the null proportion ($\pi_0 \in \{0.5, 0.9, 0.99\}$), stratified by methods. Panel (b): Comparison of the distribution of FDP, summarized as boxplots on square-root scale. The boxplots show the mean (cross), median (line), inter-quartile ranges (box), and 5th and 95th percentiles (whiskers). Panel (c): Comparison of the root MSE of FDP from the nominal FDR of 0.1, defined as $\sqrt{\text{mean}[(\text{FDP} - 0.1)^2]}$. In all scenarios the distribution of FDP for **cashr** is more concentrated near the nominal 0.1 level than other methods. Especially, the root MSE of FDP for **cashr** is uniformly lower than other methods. Panel (d): Comparison of the mean of TDP, as an indication of statistical power. On average, **cashr** maintains good power, only worse than **ashr** in some scenarios, which sometimes finds more true signals at the cost of severely losing control of FDP.

We note that **cashr** performs well even in settings that do not fully satisfy its underlying assumptions (e.g. where g_1 is asymmetric or multimodal). Note also that for our choices of g_1 , $\pi_0 = 0.99$ is a highly

sparse setting, as a large portion of the non-zero θ_j are close to zero. For example, when g_1 is Gaussian, only about 3 out of 10^4 $|\theta_j|$ are expected to be larger than $\sqrt{2\log p} \approx 4.3$. Therefore, it is understandable that no methods perform particularly well in this difficult setting. But even for this $\pi_0 = 0.99$ setting, although first impressions from the plot may be that **cashr** and BH perform similarly, closer visual inspection shows **cashr** to be better, in that its median FDP tends to be closer to 0.1.

The reason that **cashr** produces more consistently reliable FDP is that, by design, it adapts itself to the particular correlation-induced distortion present in each data set. As illustrated in Figure 1, correlation can lead to pseudo-inflation in some data sets and pseudo-deflation in others. **cashr** is able to recognize which pattern is present, and correspondingly modify its behavior – becoming more conservative in the former case and less conservative in the latter. This is illustrated in Figure 6, which stratifies the realized data sets according to sample standard deviation of the realized correlated $N(0,1)$ noise Z in each data set (for the setting where g_1 is Gaussian, $\pi_0 = 0.9$). The bottom 1/3 are categorized as pseudo-deflation, top 1/3 pseudo-inflation, and the others “in-between.”

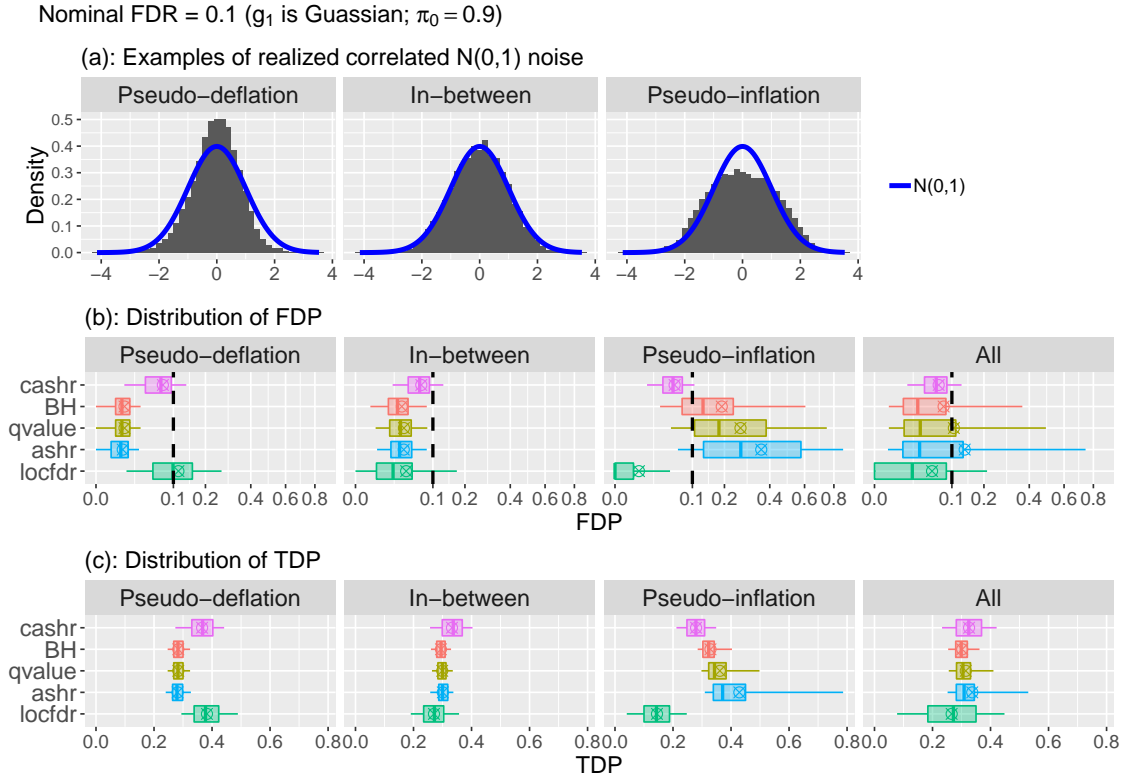


Figure 6: Illustration that **cashr** consistently produces reliable FDP under different types of correlation-induced distortion. Here we take the results from a single simulation scenario (g_1 is Gaussian, $\pi_0 = 0.9$) and stratify them into three groups of equal size according to the sample standard deviations of the realized correlated $N(0,1)$ noise. Methods that ignore correlations among observations (BH, **qvalue**, **ashr**) are generally too conservative under pseudo-deflation and too anti-conservative under pseudo-inflation; **locfdr** tends to be too conservative under pseudo-inflation and consequently lose power; **cashr** maintains good FDR control in all settings. The boxplots show the mean (cross), median (line), inter-quartile ranges (box), and 5th and 95th percentiles (whiskers). FDP are plotted on square-root scale. Other choices of g_1 and π_0 give qualitatively similar results (not shown here).

For data sets where Z show no strong distortion (“in-between”) all methods give similar and reasonable

results, with **cashr** showing only a small improvement. However, when Z are pseudo-inflated, methods ignoring correlation, such as BH, **qvalue**, **ashr**, tend to be anti-conservative; that is, they form discovery sets whose FDP are often much larger than the nominal FDR. In contrast, **cashr** produces conservative FDP near the nominal value; and **locfdr** is too conservative, consequently losing substantial power (discussed further in Section 4.2). Conversely, with pseudo-deflation, methods ignoring correlation are too conservative, producing FDP much smaller than the nominal FDR, losing power compared with **cashr** and **locfdr**.

4.2 Real data illustrations

We now use two real data examples to illustrate some of the features of **cashr** (and other methods) that we observed in simulated data. The first example is a well-studied data set from a leukemia study (Golub et al., 1999), comparing gene expression in 47 acute myeloid leukemia vs 25 acute lymphoblastic leukemia samples, which was discussed extensively in Efron (2010a) as a prime example of how correlation can distort empirical distributions. The second example comes from a study on embryonic mouse hearts (Smemo, 2012), comparing gene expression in 2 left ventricle samples vs 2 right ventricle samples. (The number of samples is small, but each sample is a pool of ventricles from 150 mice – necessary to obtain sufficient tissue for the experiments to work well – and so this experiment involved dissection of 300 mouse hearts.)

For each data set we let θ_j denote the true \log_2 -fold change in gene expression between the two groups for gene j . We use a standard analysis protocol (based on Smyth (2004); see Appendix D for details) to obtain an estimate X_j for θ_j , and a corresponding p -value p_j . As in Section 2.1, we convert the p -value to the corresponding z -score z_j and use this to compute an effective standard deviation s_j .

Figure 7 shows the empirical distribution of the z -scores for each data set, together with the fitted correlated noise distribution from **cashr** and the fitted empirical null from **locfdr**. In both cases the histogram is substantially more dispersed than $N(0,1)$. However the two data sets have otherwise quite different patterns of inflation: the leukemia data show inflation in both the shoulders and tails of the distribution, whereas the mouse data show inflation only in the shoulders. This indicates the presence of some strong signals in the leukemia data, whereas the inflation in the mouse data may be primarily pseudo-inflation caused by correlation. Consistent with this, both **locfdr** and **cashr** identify hundreds of significant signals in the leukemia data (at nominal FDR = 0.1), but no significant signals in the mouse data (Table 1).

Method	Number of discoveries	
	Leukemia data	Mouse data
cashr	385	0
locfdr	282	0
BH	1579	4130
qvalue	1972	6502
ashr	3346	17191

Table 1: Numbers of discoveries from different methods at nominal FDR = 0.1. We analyzed 7128 genes in the leukemia data and 17191 genes in the mouse data. In both data sets, the z -score distributions appear to have correlation-induced inflation, and the numbers of significant discoveries declared by methods accounting for correlation (**cashr** and **locfdr**) are much smaller than those ignoring correlation (BH, **qvalue**, **ashr**). For the leukemia data, **cashr** finds 37% more significant genes than **locfdr**.

Although the conclusions from **cashr** and **locfdr** are, here, qualitatively similar, there are some notable differences in their results. First, in the mouse data, the **cashr** correlated noise distribution gives, visually, a much better fit than the **locfdr** empirical null, particularly in the tails (Figure 7). This is because the **cashr** correlated noise distribution is ideally suited to capture this “shoulder-but-not-tail” inflation pattern that is symptomatic of correlation-induced inflation. The Gaussian empirical null distribution assumed by **locfdr** is simply inconsistent with these data. Indeed, this inconsistency is reflected in the null proportion estimated by **locfdr** (1.08) which exceeds the theoretical upper bound of 1.

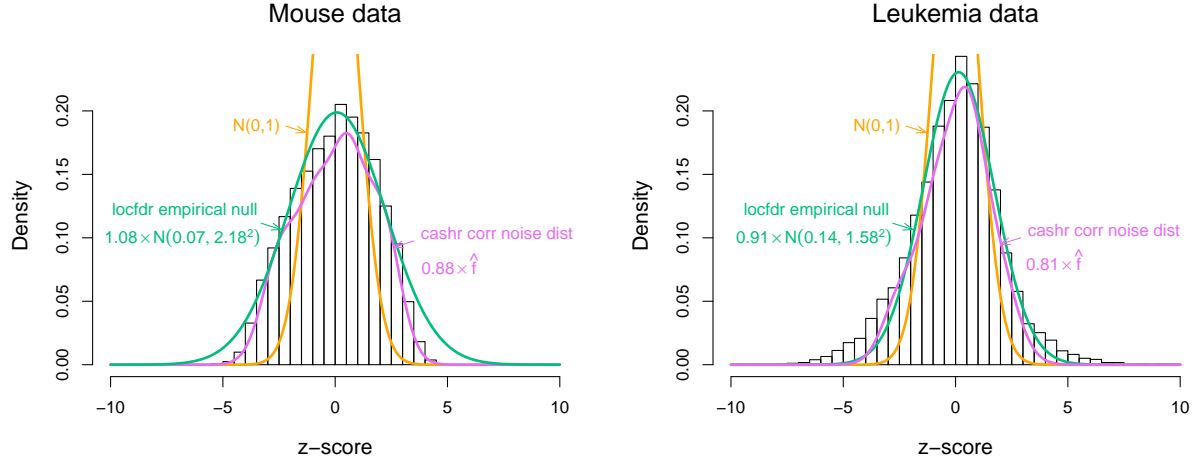


Figure 7: Distribution of z -scores from analyzing gene differential expression in two real data sets. In both data sets, for each gene j , a z -score z_j is computed, and $z_j \sim N(0, 1)$ under the null hypothesis of no differential expression. Then we compare the histogram of z -scores with $N(0, 1)$, the fitted correlated noise distribution from **cashr**, and the fitted empirical null from **locfdr**, scaled by respective estimated null proportions. Both histograms are substantially more dispersed than $N(0, 1)$. The mouse data show inflation primarily in the shoulders of the distribution, and the fitted correlated noise distribution from **cashr** appears to be a much better fit than the fitted empirical null from **locfdr**, particularly in the tails. The leukemia data show inflation in both shoulders and tails of the distribution, indicating the presence of some strong signals. Although otherwise similar, the fitted correlated noise distribution from **cashr** has a noticeably shorter right tail than the fitted empirical null from **locfdr**, improving power.

Second, in the leukemia data, **cashr** identifies 37% more significant results than **locfdr** (385 vs 282). This is consistent with the greater power of **cashr** vs **locfdr** in our simulations. One reason that **locfdr** can lose power is that its Gaussian empirical null distribution tends to overestimate inflation in the tails when it tries to fit inflation in the shoulders. We see this feature in the mouse data, and although less obvious, this appears to also be the case for the leukemia data: the estimated standard deviation of the empirical null is 1.58, which is almost certainly too large: a pseudo-inflated Gaussian correlated noise distribution is unlikely to have standard deviation exceeding 1.4 (Appendix B). In comparison the fitted correlated noise distribution from **cashr** has a noticeably shorter right tail (e.g. $z \in [4, 5]$) which leads it to categorize more z -scores in the right tail as significant (Figure 7). On a side note, **cashr** also experiences the benefits of **ashr** highlighted in Stephens (2017), which can also help increase power. For example, the unimodal assumption on the effects – which allows that some of the z -scores around zero may correspond to true, albeit non-significant, signals – can help improve estimates of π_0 , and hence improve power.

Another feature of **cashr**, which distinguishes it from **locfdr**, is that, by estimating g while accounting for correlation-induced distortion, it can provide an estimate on the effect size distribution, g_1 . For the mouse data, **cashr** estimates $\hat{\pi}_0 = 0.88$, or 12% of genes may be differentially expressed to some extent, although it is not able to pin down any clear example of a differentially expressed gene: no gene has an estimated local FDR less than 0.80. One possible explanation for the lack of significant results in this case is lack of power. However, the estimated g_1 from **cashr** suggests that there may simply not exist any large effects to be discovered: 99% of the probability mass of the estimated g_1 is on effect size ≤ 0.26 , or a mere 1.2-fold change in gene expression. Thus the signals here, if any, are too weak to be discerned from noise and pseudo-inflation.

We also applied the other methods – BH, **qvalue**, and **ashr** – to both data sets. All three methods

find very large numbers of significant results in both data sets (Table 1). Although we do not know the truth in these real data, there is a serious concern that many of these results could be false positives, since these methods are all prone to erroneously viewing pseudo-inflation as true signal (Figure 6), and Figure 7 suggests that pseudo-inflation may be present in both data sets.

5 Discussion

We have presented a general approach to accounting for correlations among observations in the widely-used Empirical Bayes Normal Means model. Our strategy exploits theoretical results from Schwartzman (2010) to model the impact of correlation on the empirical distribution of correlated $N(0, 1)$ variables, and convex optimization techniques to produce an efficient implementation. We demonstrated through empirical examples that this strategy can both improve estimation of the underlying distribution of true effects (Figure 2) and – in the multiple testing setting – improve estimation of FDR compared with EB methods that ignore correlation (Figures 5, 6). To the best of our knowledge, **cashr** is the first EBNM methodology to deal with correlated noise in this way.

Our empirical results demonstrate some benefits of the EB approach to multiple testing compared with traditional methods. In particular, **cashr** provides, on average, more accurate estimates of the FDP than either BH or **qvalue**. However, although we find these empirical results encouraging, we do not have theoretical guarantees of (frequentist) FDR control. That said, theoretical guarantees of FDR control under arbitrary correlation structure are lacking even for the widely-studied BH method. BH has been shown to control FDR under certain correlation structures (e.g. “positive regression dependence on subsets”; Benjamini and Yekutieli, 2001). The Benjamini-Yekutieli procedure (Benjamini and Yekutieli, 2001) is proved to control FDR under arbitrary dependence, but at the cost of being excessively conservative, and is consequently rarely used in practice.

A key feature of **cashr** is that it requires no information about the actual correlations among observations. This has the important advantage that it can be applied wherever EBNM methods that ignore correlation can be applied. On the other hand, when additional information on correlations is available it clearly may be helpful to incorporate it into analyses. Within our approach such information could be used to estimate the moments of the pairwise correlations, and thus inform estimates of ω in the correlated noise distribution $f(\cdot; \omega)$. Alternatively, one could take a more ambitious approach: explicitly model the whole $p \times p$ correlation matrix, and use this to help inform inference (e.g. Benjamini and Heller, 2007; Wu, 2008; Sun and Cai, 2009; Friguet et al., 2009; Fan et al., 2012). Modeling correlation is likely to provide more efficient inferences when it can be accurately achieved (Hall and Jin, 2010). However, in many situations – particularly involving small sample sizes – reliably modeling correlation may be impossible. Under what circumstances this more ambitious approach produces better inferences could be one area for future investigation.

The main assumptions underlying **cashr** are that the correlated noise is marginally $N(0, 1)$, and that the standard deviations are reliably computed. In the multiple testing setting this corresponds to assuming that the test statistics are (marginally) well calibrated. If these conditions do not hold – for example, due to failure of asymptotic theory underlying test statistic computations, or due to confounding factors (such as batch effects in gene expression studies), then **cashr** could give unreliable results. Of course **cashr** is not unique in this regard – methods like BH and **qvalue** similarly assume that test statistics are well calibrated. Dealing with confounders in gene expression studies is an active area of research, and several approaches exist, many of them based on factor analysis (e.g. Leek and Storey, 2007; Sun et al., 2012; Gagnon-Bartsch and Speed, 2012; Wang et al., 2017; Gerard and Stephens, 2017, 2018). Again, the possibility of combining these ideas with our methods could be a future research direction.

Appendix

A The marginal distribution of the simulated null z -scores

Figures A.1 and A.2 offer support for the claim that the z -scores simulated in Section 2.1 are marginally $N(0, 1)$ -distributed.

Figure A.1 compares z -scores simulated as in Section 2.1 with z -scores simulated under a modified framework that removes gene-gene correlations, and with iid $N(0, 1)$ samples. The modified framework uses exactly the same simulation and analysis pipeline as the original framework of Section 2.1, with one important difference: in each simulation, *for each gene independently* we randomly selected two groups of five samples without replacement, hence removing gene-gene correlations.

The empirical CDF of 10^4 data sets simulated as in Section 2.1 show a huge amount of variability (panel (a)), presumably due to correlations among genes. In the modified framework, correlation-induced distortion disappears: the empirical CDF of all 10^4 data sets are almost exactly the same as $N(0, 1)$ (panel (b)), just as with the iid $N(0, 1)$ samples (panel (c)). This demonstrates that without gene-gene correlations, the analysis pipeline used here produces uncorrelated $N(0, 1)$ z -scores.

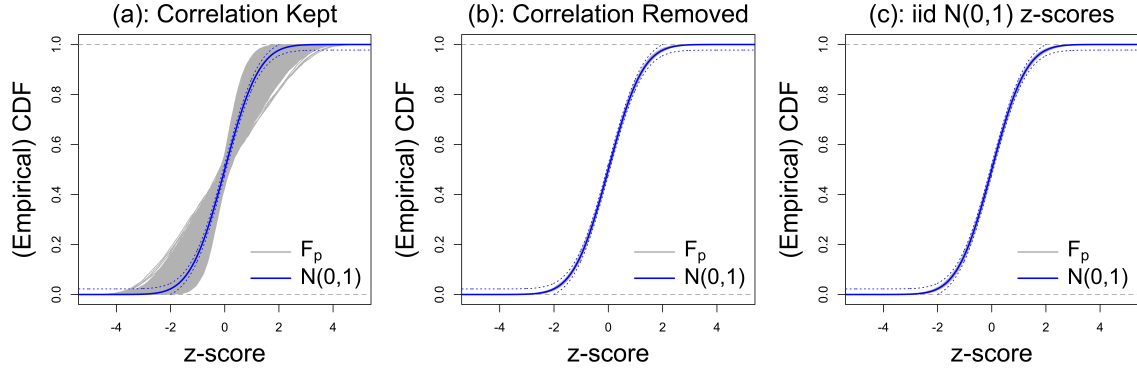


Figure A.1: Comparison of 10^4 empirical CDF of z -scores (F_p) obtained by applying the same analysis pipeline to data simulated by two different frameworks: the original framework in Section 2.1 which keeps gene-gene correlations (panel (a)); and the modified framework to remove gene-gene correlations by randomizing samples for each gene (panel (b)). We also plot 10^4 empirical CDF of iid $N(0, 1)$ samples for comparison (panel (c)). The z -scores obtained under the original framework show clear correlation-induced distortion – the variability of empirical CDF is huge. In contrast, when gene-gene correlations are removed under the modified framework, distortion disappears: empirical CDF are almost exactly the same as $N(0, 1)$ and the variability is essentially invisible; indeed, they are indistinguishable from 10^4 empirical CDF of iid $N(0, 1)$ z -scores. It shows clear evidence that the analysis pipeline can produce well-calibrated null z -scores if no gene-gene correlations. Dotted lines are Dvoretzky-Kiefer-Wolfowitz bounds with $\alpha = 1/10^4$.

In addition, Figure A.2 shows that the mean empirical CDF of the 10^4 data sets simulated from the original framework – the average of empirical CDF of Figure A.1(a) – is very close to $N(0, 1)$. Possible deviation happens only in the far tails ($|\cdot| \in \{5, 6\}$). Compared with $N(0, 1.05^2)$ and $N(0, 1.1^2)$, the deviation is very small even on the logarithmic scale (panels (b-c)), probably caused by numerical constraints as one or two z -scores in this area in a few data sets can make a visible difference.

B Decompose Gaussian by standardized Gaussian derivatives

Proposition 1. *The PDF of $N(\mu, \sigma^2)$ can be decomposed by standard Gaussian and its derivatives in the form of (7) if and only if $\sigma^2 \leq 2$.*

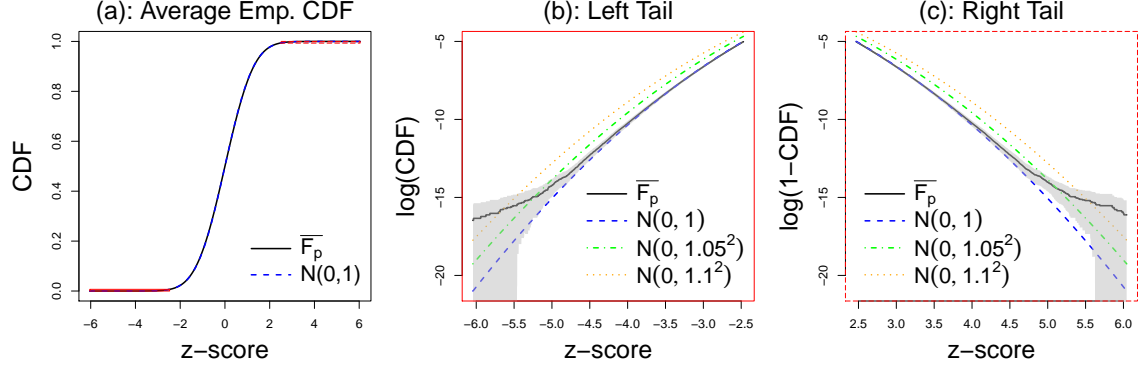


Figure A.2: Illustration that the average empirical CDF of z -scores (\bar{F}_p) simulated as in Section 2.1 closely matches $N(0,1)$, aggregated over 10^4 data sets. Left: the average of all empirical CDF in Figure A.1(a). The average empirical CDF is extremely close to $N(0,1)$. Center and Right: the left and right tails of the average empirical CDF on logarithmic scale. Shaded areas are 99.9% confidence bands. Compared with $N(0, 1.05^2)$ and $N(0, 1.1^2)$, possible deviation from $N(0,1)$ is light even in the far tails.

Proof. Let $h_l(\cdot)$ denote the l^{th} probabilists' Hermite polynomial. The orthogonality and completeness of Hermite polynomials in $L^2(\mathbb{R}, d\Phi)$ (e.g. Szegő, 1975) leads to the following fact

$$\int_{\mathbb{R}} \frac{1}{\sqrt{m!}} h_m(x) \frac{1}{\sqrt{n!}} \varphi^{(n)}(x) dx = (-1)^n \delta_{mn}, \quad \forall m, n = 0, 1, 2, \dots, \quad (31)$$

where $\delta_{mn} = \begin{cases} 1 & m = n \\ 0 & \text{otherwise} \end{cases}$. Therefore, if any PDF f can be decomposed in the form of (7), the coefficient of the l^{th} -order standardized Gaussian derivative has to be

$$w_l = (-1)^l \int_{\mathbb{R}} \frac{1}{\sqrt{l!}} h_l(x) f(x) dx = \frac{(-1)^l}{\sqrt{l!}} E_f[h_l], \quad (32)$$

where $E_f[h_l]$, sometimes called ‘‘Hermite moment,’’ is the expected value of $h_l(\cdot)$ when the PDF of the random variable is f . If f is $N(\mu, \sigma^2)$, we can obtain analytic expressions of these Hermite moments

$$E_f[h_l] = \mu^l + \sum_{k=1}^{\lfloor l/2 \rfloor} \binom{l}{2k} \mu^{l-2k} (\sigma^2 - 1)^k (2k-1)!! := M_l(\mu, \sigma^2 - 1), \quad (33)$$

where $n!!$ denotes the double factorial of n , and $M_l(x, y)$ denotes the function of l^{th} -order moment of a Gaussian with mean x and variance y . Putting (32)-(33) together, the coefficients in (7) become

$$w_l = \frac{(-1)^l}{\sqrt{l!}} M_l(\mu, \sigma^2 - 1). \quad (34)$$

Note that w_l is not exploding if and only if $|\sigma^2 - 1| \leq 1$ or equivalently, $\sigma^2 \leq 2$. \square

This result suggests that a pseudo-inflated Gaussian correlated noise distribution is not likely to have standard deviation greater than $\sqrt{2} \approx 1.4$.

In the special case when $\rho_{ij} = 1$, f becomes δ_z , a point mass on $Z \equiv z$, with z randomly sampled from $N(0,1)$. It is interesting to note that δ_z can be decomposed in the form of (7) as

$$\delta_z(\cdot) = \varphi(\cdot) + \sum_{l=1}^{\infty} \left[\frac{(-1)^l}{\sqrt{l!}} h_l(z) \right] \left[\frac{1}{\sqrt{l!}} \varphi^{(l)}(\cdot) \right], \quad \forall z \in \mathbb{R}.$$

C Proof of Theorem 1

Proof. The marginal distribution of X_j , denoted as $p(X_j)$, is obtained by integrating out θ_j

$$\begin{aligned}
p(X_j) &= \int_{\mathbb{R}} g(\theta_j) p(X_j | \theta_j, s_j) d\theta_j = \int_{\mathbb{R}} g(\theta_j) \frac{1}{s_j} f\left(\frac{X_j - \theta_j}{s_j}\right) d\theta_j \\
&= \int_{\mathbb{R}} \left[\pi_0 \delta_0(\theta_j) + \sum_{k=1}^K \pi_k \frac{1}{\sigma_k} \varphi\left(\frac{\theta_j}{\sigma_k}\right) \right] \left[\frac{1}{s_j} \varphi\left(\frac{X_j - \theta_j}{s_j}\right) + \frac{1}{s_j} \sum_{l=1}^L \omega_l \frac{1}{\sqrt{l!}} \varphi^{(l)}\left(\frac{X_j - \theta_j}{s_j}\right) \right] d\theta_j \\
&= \sum_{k=0}^K \pi_k \left(p_{jk0} + \sum_{l=1}^L \omega_l p_{jkl} \right), \tag{35}
\end{aligned}$$

where $p_{jkl} = \int_{\mathbb{R}} \frac{1}{\sigma_k} \varphi\left(\frac{\theta_j}{\sigma_k}\right) \frac{1}{s_j} \frac{1}{\sqrt{l!}} \varphi^{(l)}\left(\frac{X_j - \theta_j}{s_j}\right) d\theta_j$ is essentially a convolution of φ and $\varphi^{(l)}$ and has an analytic form

$$p_{jkl} = \frac{s_j^l}{\sqrt{\sigma_k^2 + s_j^2}^{l+1}} \frac{1}{\sqrt{l!}} \varphi^{(l)}\left(\frac{X_j}{\sqrt{\sigma_k^2 + s_j^2}}\right).$$

This form is also valid for $k = 0, l = 0$. Following (35), the marginal log-likelihood of π, ω is given by

$$\log \left(\prod_{j=1}^n p(X_j) \right) = \sum_{j=1}^n \log \left(\sum_{k=0}^K \pi_k \left(p_{jk0} + \sum_{l=1}^L \omega_l p_{jkl} \right) \right).$$

□

D Simulation details

D.1 Six choices of the non-null effect distribution

g_1	PDF	$E[\cdot ^2]$	$\Pr(\cdot \geq \sqrt{2 \log p})$
Gaussian	$N(0, 2^2)$	4	0.032
Near Gaussian	$0.6N(0, 1) + 0.4N(0, 3^2)$	4.2	0.061
Spiky	$0.4N(0, 0.5^2) + 0.2N(0, 2^2) + 0.4N(0, 3^2)$	4.5	0.067
Skew	$0.25N(-2, 2^2) + 0.25N(-1, 2^2) + 0.25N(0, 1) + 0.25N(1, 1)$	4	0.045
Flat Top	$0.5N(-1.5, 1.5^2) + 0.5N(1.5, 1.5^2)$	4.5	0.031
Bimodal	$0.5N(-1.5, 1) + 0.5N(1.5, 1)$	3.25	0.0026

Table 2: Details of six choices of g_1 , the distribution of non-null effects in Section 4.1. The table also shows the average signal strength, $E[|\cdot|^2]$, and the probability of large signal, $\Pr(|\cdot| \geq \sqrt{2 \log p})$, conditioned on g_1 .

D.2 Implementation of methods

The existing methods we use for comparison in this paper mostly use the default settings in their respective R packages. That include **REBayes**, **deconvolveR** for deconvolution (Section 2.1), and **qvalue**, **locfdr** for multiple testing (Section 4). For **EbayesThresh**, we set **a=NA** to allow the scale parameter of the Laplace distribution to be estimated from the data. For **ashr**, we set **mixcompdist="normal"** to use scale mixture of zero-mean Gaussians to approximate g .

D.3 Pipeline for analyzing gene expression data

Let θ_j denote the true \log_2 -fold change in gene expression for each gene j . The analysis pipeline is used to provide, for each θ_j , an estimate X_j with a standard error s_j , such that X_j can be assumed to be $N(\theta_j, s_j^2)$.

For RNA-seq data such as the mouse data, the analysis pipeline is described in Section 2.1.

For microarray data such as the leukemia data, we use a widely-used analysis protocol implemented in the `limma` software (Ritchie et al., 2015). This yield an estimate X_j for θ_j , and a corresponding p -value p_j from a moderated t -statistic (Smyth, 2004). Then as in Section 2.1, we convert the p -value to the corresponding z -score z_j and use it to compute the effective standard deviation s_j .

D.4 Reproducibility

All the code generating the results and plots in this paper are available at https://github.com/LSun/cashr_paper.

The RNA-seq gene expression data from human liver tissues we used in this paper were generated by the GTEx Project, which was supported by the Common Fund of the Office of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. The data used for the analyses described in this paper were obtained from the GTEx Portal at <https://www.gtexportal.org>. In particular, the human liver RNA-seq data for realistic simulation are also available at https://github.com/LSun/cashr_paper.

The leukemia microarray data are available at <http://statweb.stanford.edu/~ckirby/brad/LSI/datasets-and-programs/datasets.html>.

The mouse heart RNA-seq data are available at https://github.com/LSun/cashr_paper.

E The representation of correlated noise distribution

If Z are independent and p is large then F_p will be close to its mean, Φ . This is guaranteed by well established results like the Glivenko-Cantelli theorem and the Dvoretzky-Kiefer-Wolfowitz inequality (e.g. Wasserman, 2006). However, when Z are correlated F_p can be grossly different from Φ , as we have seen in Section 2.1. The covariance of F_p indicates how far it tends to stray from its mean, Φ , and therefore captures the extent of correlation-induced distortion. Schwartzman (2010) provides the following elegant characterization of the covariance of F_p . For completeness we also put it here.

Proposition 2. *(The mean, variance, and covariance functions of F_p ; Schwartzman, 2010)*

Assume $\forall i \neq j$, $\begin{bmatrix} Z_i \\ Z_j \end{bmatrix} \sim N\left(0, \begin{bmatrix} 1 & \rho_{ij} \\ \rho_{ij} & 1 \end{bmatrix}\right)$. Let $\bar{\rho}^l := \frac{1}{p(p-1)} \sum_{i,j:i \neq j} \rho_{ij}^l$. Then $\forall x, y \in \mathbb{R}$,

$$E(F_p(x)) = \Phi(x) \tag{36}$$

$$\text{var}(F_p(x)) = \left(1 - \frac{1}{p}\right) \sum_{l=1}^{\infty} \bar{\rho}^l \left[\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(x) \right]^2 + \frac{1}{p} \Phi(x)(1 - \Phi(x)) \tag{37}$$

$$\text{cov}(F_p(x), F_p(y)) = \left(1 - \frac{1}{p}\right) \sum_{l=1}^{\infty} \bar{\rho}^l \left[\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(x) \right] \left[\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(y) \right] + \frac{1}{p} [\Phi(\min(x, y)) - \Phi(x)\Phi(y)] \tag{38}$$

Proof. The mean function is straightforward. The covariance function

$$\begin{aligned}
\text{cov}(F_p(x), F_p(y)) &= \text{cov} \left(\frac{1}{p} \sum_{i=1}^p \mathcal{I}(Z_i \leq x), \frac{1}{p} \sum_{j=1}^p \mathcal{I}(Z_j \leq y) \right) \\
&= E \left[\left(\frac{1}{p} \sum_{i=1}^p \mathcal{I}(Z_i \leq x) \right) \left(\frac{1}{p} \sum_{j=1}^p \mathcal{I}(Z_j \leq y) \right) \right] - E \left[\frac{1}{p} \sum_{i=1}^p \mathcal{I}(Z_i \leq x) \right] E \left[\frac{1}{p} \sum_{j=1}^p \mathcal{I}(Z_j \leq y) \right] \\
&= \frac{1}{p^2} \sum_{i=1}^p \sum_{j=1}^p E[\mathcal{I}(Z_i \leq x) \mathcal{I}(Z_j \leq y)] - \Phi(x)\Phi(y) \\
&= \frac{1}{p^2} \sum_{i=1}^p \sum_{j=1}^p P(Z_i \leq x, Z_j \leq y) - \Phi(x)\Phi(y) \\
&= \frac{1}{p^2} \sum_{i \neq j} P(Z_i \leq x, Z_j \leq y) + \frac{1}{p} \Phi(\min(x, y)) - \Phi(x)\Phi(y) .
\end{aligned} \tag{39}$$

According to Mehler's identity (Kibble, 1945), under the assumption of $\{Z_i, Z_j\}$ being bivariate normal, the joint PDF can be written as

$$p(x, y) = \varphi(x)\varphi(y) + \sum_{l=1}^{\infty} \rho_{ij}^l \left[\frac{1}{\sqrt{l!}} \varphi^{(l)}(x) \right] \left[\frac{1}{\sqrt{l!}} \varphi^{(l)}(y) \right] , \tag{40}$$

so the joint CDF is

$$P(Z_i \leq x, Z_j \leq y) = \Phi(x)\Phi(y) + \sum_{l=1}^{\infty} \rho_{ij}^l \left[\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(x) \right] \left[\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(y) \right] . \tag{41}$$

(39) and (41) lead to the covariance function (38). Setting $x = y$ gives the variance function (37). \square

Note that $\text{var}(F_p)$ has two parts. The second part $\frac{1}{p}\Phi(z)(1-\Phi(z))$ is the familiar variance function when Z are independent, and it quickly vanishes as p increases. This is why F_p of iid $N(0, 1)$ sample will not deviate much from Φ when p is large. In contrast, the first part

$$\left(1 - \frac{1}{p}\right) \sum_{l=1}^{\infty} \bar{\rho}^l \left(\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(x) \right)^2 \tag{42}$$

demonstrates the effect of correlation. If $\bar{\rho}^l$ is non-negligible for large p , $\text{var}(F_p)$ will be non-vanishing, and so F_p and the histogram of Z are more likely to deviate substantially from $N(0, 1)$.

When p is large,

$$\text{cov}(F_p(x), F_p(y)) \approx \sum_{l=1}^{\infty} \bar{\rho}^l \left[\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(x) \right] \left[\frac{1}{\sqrt{l!}} \varphi^{(l-1)}(y) \right] . \tag{43}$$

(36) and (43) suggest we can characterize F_p as (5) (Schwartzman, 2010), assuming $\bar{\rho}^l \geq 0$ for all $l \in \mathbb{N}$. This assumption should not be too demanding for large p . For example, when $l = 1$,

$$\bar{\rho} = \frac{1}{p(p-1)} \sum_{i \neq j} \rho_{ij} = \frac{1}{p(p-1)} (\mathbf{1}^T \Sigma_Z \mathbf{1} - p) \geq \frac{1}{p(p-1)} (-p) = -\frac{1}{p-1} , \tag{44}$$

following the fact that Σ_Z , the correlation matrix of Z , is positive semi-definite.

References

- Benjamini, Y. and Heller, R. (2007). “False Discovery Rates for Spatial Signals.” *Journal of the American Statistical Association*, 102(480): 1272–1281.
URL <https://doi.org/10.1198/016214507000000941>
- Benjamini, Y. and Hochberg, Y. (1995). “Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing.” *Journal of the Royal Statistical Society. Series B (Methodological)*, 57(1): 289–300.
- Benjamini, Y. and Yekutieli, D. (2001). “The control of the false discovery rate in multiple testing under dependency.” *Ann. Statist.*, 29(4): 1165–1188.
URL <https://doi.org/10.1214/aos/1013699998>
- Berger, J. O. (1985). *Statistical decision theory and Bayesian analysis*. Springer.
- Blanchard, G. and Roquain, E. (2009). “Adaptive False Discovery Rate Control Under Independence and Dependence.” *J. Mach. Learn. Res.*, 10: 2837–2871.
URL <http://dl.acm.org/citation.cfm?id=1577069.1755880>
- Bovy, J., Hogg, D. W., and Roweis, S. T. (2011). “Extreme deconvolution: Inferring complete distribution functions from noisy, heterogeneous and incomplete observations.” *Ann. Appl. Stat.*, 5(2B): 1657–1677.
URL <https://doi.org/10.1214/10-A0AS439>
- Brown, L. D. (2008). “In-season prediction of batting averages: A field test of empirical Bayes and Bayes methodologies.” *Ann. Appl. Stat.*, 2(1): 113–152.
URL <https://doi.org/10.1214/07-A0AS138>
- Brown, L. D. and Greenshtein, E. (2009). “Nonparametric empirical Bayes and compound decision approaches to estimation of a high-dimensional vector of normal means.” *Ann. Statist.*, 37(4): 1685–1704.
URL <https://doi.org/10.1214/08-A0S630>
- Clyde, M. and George, E. I. (2000). “Flexible empirical Bayes estimation for wavelets.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 62(4): 681–698.
URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/1467-9868.00257>
- Cordy, C. B. and Thomas, D. R. (1997). “Deconvolution of a Distribution Function.” *Journal of the American Statistical Association*, 92(440): 1459–1465.
URL <https://doi.org/10.1080/01621459.1997.10473667>
- Dawid, A. P. (1994). *Selection paradoxes of Bayesian inference*, volume Volume 24 of *Lecture Notes–Monograph Series*, 211–220. Hayward, CA: Institute of Mathematical Statistics.
URL <https://doi.org/10.1214/lnms/1215463797>
- De Carvalho, M. and Ramos, A. (2012). “Bivariate extreme statistics, II.” *REVSTAT-Statistical Journal*, 10(EPFL-ARTICLE-180505): 83–107.
- De Finetti, B. (1937). “La prévision: ses lois logiques, ses sources subjectives [Foresight: its logical laws, its subjective sources].” *Annales de l’institut Henri Poincaré*, 7(1): 1–68.
- Dey, K. K. and Stephens, M. (2018). “CorShrink : Empirical Bayes shrinkage estimation of correlations, with applications.” *bioRxiv*.
URL <https://www.biorxiv.org/content/early/2018/07/24/368316>
- Efron, B. (1996). “Empirical Bayes Methods for Combining Likelihoods.” *Journal of the American Statistical Association*, 91(434): 538–550.
URL <https://amstat.tandfonline.com/doi/abs/10.1080/01621459.1996.10476919>

- (2004). “Large-Scale Simultaneous Hypothesis Testing.” *Journal of the American Statistical Association*, 99(465): 96–104.
URL <https://doi.org/10.1198/016214504000000089>
- (2005). “Local false discovery rates.” Technical report, Department of Statistics, Stanford University.
- (2007a). “Correlation and Large-Scale Simultaneous Significance Testing.” *Journal of the American Statistical Association*, 102(477): 93–103.
URL <https://doi.org/10.1198/016214506000001211>
- (2007b). “Size, power and false discovery rates.” *Ann. Statist.*, 35(4): 1351–1377.
URL <https://doi.org/10.1214/009053606000001460>
- (2008). “Microarrays, Empirical Bayes and the Two-Groups Model.” *Statist. Sci.*, 23(1): 1–22.
URL <https://doi.org/10.1214/07-STS236>
- (2010a). “Correlated z -Values and the Accuracy of Large-Scale Statistical Estimates.” *Journal of the American Statistical Association*, 105(491): 1042–1055.
- (2010b). *Large-Scale Inference: Empirical Bayes Methods for Estimation, Testing, and Prediction*. Institute of Mathematical Statistics Monographs. Cambridge University Press.
- (2016). “Empirical Bayes deconvolution estimates.” *Biometrika*, 103(1): 1–20.
URL <http://dx.doi.org/10.1093/biomet/asv068>
- (2018). “Curvature and inference for maximum likelihood estimates.” *Ann. Statist.*, 46(4): 1664–1692.
URL <https://doi.org/10.1214/17-AOS1598>
- Efron, B. and Hastie, T. (2016). *Computer Age Statistical Inference: Algorithms, Evidence, and Data Science*. Institute of Mathematical Statistics Monographs. Cambridge University Press.
- Efron, B. and Morris, C. (1972). “Limiting the Risk of Bayes and Empirical Bayes Estimators Part II: The Empirical Bayes Case.” *Journal of the American Statistical Association*, 67(337): 130–139.
URL <https://doi.org/10.1080/01621459.1972.10481215>
- (1973). “Stein’s Estimation Rule and its Competitors An Empirical Bayes Approach.” *Journal of the American Statistical Association*, 68(341): 117–130.
URL <https://doi.org/10.1080/01621459.1973.10481350>
- Fan, J. (1991). “On the Optimal Rates of Convergence for Nonparametric Deconvolution Problems.” *Ann. Statist.*, 19(3): 1257–1272.
URL <https://doi.org/10.1214/aos/1176348248>
- Fan, J., Han, X., and Gu, W. (2012). “Estimating False Discovery Proportion Under Arbitrary Covariance Dependence.” *Journal of the American Statistical Association*, 107(499): 1019–1035. PMID: 24729644.
URL <https://doi.org/10.1080/01621459.2012.720478>
- Friguet, C., Kloareg, M., and Causeur, D. (2009). “A Factor Model Approach to Multiple Testing Under Dependence.” *Journal of the American Statistical Association*, 104(488): 1406–1415.
URL <https://doi.org/10.1198/jasa.2009.tm08332>
- Gagnon-Bartsch, J. A. and Speed, T. P. (2012). “Using control genes to correct for unwanted variation in microarray data.” *Biostatistics*, 13(3): 539–552.
URL <http://dx.doi.org/10.1093/biostatistics/kxr034>
- Gelman, A., Hill, J., and Yajima, M. (2012). “Why We (Usually) Don’t Have to Worry About Multiple Comparisons.” *Journal of Research on Educational Effectiveness*, 5(2): 189–211.
URL <https://doi.org/10.1080/19345747.2011.618213>

- Gerard, D. and Stephens, M. (2017). “Unifying and Generalizing Methods for Removing Unwanted Variation Based on Negative Controls.” *ArXiv e-prints*.
- (2018). “Empirical Bayes shrinkage and false discovery rate estimation, allowing for unwanted variation.” *Biostatistics*, kxy029.
URL <http://dx.doi.org/10.1093/biostatistics/kxy029>
- Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J. P., Coller, H., Loh, M. L., Downing, J. R., Caligiuri, M. A., Bloomfield, C. D., and Lander, E. S. (1999). “Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring.” *Science*, 286(5439): 531–537.
URL <http://science.sciencemag.org/content/286/5439/531>
- Greenland, S. and Robins, J. M. (1991). “Empirical-Bayes adjustments for multiple comparisons are sometimes useful.” *Epidemiology*, 244–251.
- Hall, P. and Jin, J. (2010). “Innovated higher criticism for detecting sparse signals in correlated noise.” *Ann. Statist.*, 38(3): 1686–1732.
URL <https://doi.org/10.1214/09-A0S764>
- Jiang, W. and Zhang, C.-H. (2009). “General maximum likelihood empirical Bayes estimation of normal means.” *Ann. Statist.*, 37(4): 1647–1684.
URL <https://doi.org/10.1214/08-A0S638>
- Johnstone, I. and Silverman, B. (2004). “Needles and straw in haystacks: Empirical Bayes estimates of possibly sparse sequences.” *Ann. Statist.*, 32(4): 1594–1649.
URL <https://doi.org/10.1214/009053604000000030>
- (2005a). “EbayesThresh: R Programs for Empirical Bayes Thresholding.” *Journal of Statistical Software, Articles*, 12(8): 1–38.
URL <https://www.jstatsoft.org/v012/i08>
- (2005b). “Empirical Bayes selection of wavelet thresholds.” *Ann. Statist.*, 33(4): 1700–1752.
URL <https://doi.org/10.1214/0090536050000000345>
- Kibble, W. F. (1945). “An extension of a theorem of Mehler’s on Hermite polynomials.” *Mathematical Proceedings of the Cambridge Philosophical Society*, 41(1): 1215.
- Kiefer, J. and Wolfowitz, J. (1956). “Consistency of the Maximum Likelihood Estimator in the Presence of Infinitely Many Incidental Parameters.” *Ann. Math. Statist.*, 27(4): 887–906.
URL <https://doi.org/10.1214/aoms/1177728066>
- Kim, Y., Carbonetto, P., Stephens, M., and Anitescu, M. (2018). “A Fast Algorithm for Maximum Likelihood Estimation of Mixture Proportions Using Sequential Quadratic Programming.” *ArXiv e-prints*.
- Koenker, R. and Gu, J. (2017). “REBayes: An R Package for Empirical Bayes Mixture Methods.” *Journal of Statistical Software, Articles*, 82(8): 1–26.
URL <https://www.jstatsoft.org/v082/i08>
- Koenker, R. and Mizera, I. (2014). “Convex Optimization, Shape Constraints, Compound Decisions, and Empirical Bayes Rules.” *Journal of the American Statistical Association*, 109(506): 674–685.
URL <https://doi.org/10.1080/01621459.2013.869224>
- Laird, N. (1978). “Nonparametric Maximum Likelihood Estimation of a Mixing Distribution.” *Journal of the American Statistical Association*, 73(364): 805–811.
URL <https://amstat.tandfonline.com/doi/abs/10.1080/01621459.1978.10480103>

- Law, C. W., Chen, Y., Shi, W., and Smyth, G. K. (2014). “voom: precision weights unlock linear model analysis tools for RNA-seq read counts.” *Genome Biology*, 15(2): R29.
URL <https://doi.org/10.1186/gb-2014-15-2-r29>
- Leek, J. T. and Storey, J. D. (2007). “Capturing Heterogeneity in Gene Expression Studies by Surrogate Variable Analysis.” *PLOS Genetics*, 3(9): 1–12.
URL <https://doi.org/10.1371/journal.pgen.0030161>
- Lu, M. (2018). “Generalized adaptive shrinkage methods and applications in genomic studies.” Ph.D. thesis, University of Chicago.
- Morris, C. N. (1983). “Parametric Empirical Bayes Inference: Theory and Applications.” *Journal of the American Statistical Association*, 78(381): 47–55.
URL <https://amstat.tandfonline.com/doi/abs/10.1080/01621459.1983.10477920>
- MOSEK ApS (2018). *MOSEK Rmosek Package. Release 8.1.0.51*.
URL <https://docs.mosek.com/8.1/rmosek/index.html>
- Muralidharan, O. (2010). “An empirical Bayes mixture method for effect size and false discovery rate estimation.” *Ann. Appl. Stat.*, 4(1): 422–438.
URL <https://doi.org/10.1214/09-A0AS276>
- Narasimhan, B. and Efron, B. (2016). “A G-modeling Program for Deconvolution and Empirical Bayes Estimation.” Technical report, Department of Statistics, Stanford University.
- Owen, A. B. (2005). “Variance of the number of false discoveries.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(3): 411–426.
- Petrone, S., Rousseau, J., and Scricciolo, C. (2014). “Bayes and empirical Bayes: do they merge?” *Biometrika*, 101(2): 285–302.
URL <http://dx.doi.org/10.1093/biomet/ast067>
- Qiu, X., Klebanov, L., and Yakovlev, A. (2005). “Correlation between gene expression levels and limitations of the empirical Bayes methodology for finding differentially expressed genes.” *Statistical applications in genetics and molecular biology*, 4(1).
- Ritchie, M. E., Phipson, B., Wu, D., Hu, Y., Law, C. W., Shi, W., and Smyth, G. K. (2015). “limma powers differential expression analyses for RNA-sequencing and microarray studies.” *Nucleic Acids Research*, 43(7): e47.
URL <http://dx.doi.org/10.1093/nar/gkv007>
- Robbins, H. (1956). “An Empirical Bayes Approach to Statistics.” In *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, 157–163. Berkeley, Calif.: University of California Press.
URL <https://projecteuclid.org/euclid.bsm/1200501653>
- (1964). “The Empirical Bayes Approach to Statistical Decision Problems.” *Ann. Math. Statist.*, 35(1): 1–20.
URL <https://doi.org/10.1214/aoms/1177703729>
- Robinson, M. D., McCarthy, D. J., and Smyth, G. K. (2010). “edgeR: a Bioconductor package for differential expression analysis of digital gene expression data.” *Bioinformatics*, 26(1): 139–140.
URL <http://dx.doi.org/10.1093/bioinformatics/btp616>
- Rousseau, J. and Szabo, B. (2017). “Asymptotic behaviour of the empirical Bayes posteriors associated to maximum marginal likelihood estimator.” *Ann. Statist.*, 45(2): 833–865.
URL <https://doi.org/10.1214/16-AOS1469>

- Schwartzman, A. (2010). “Comment.” *Journal of the American Statistical Association*, 105(491): 1059–1063.
- Scott, J. G. and Berger, J. O. (2010). “Bayes and empirical-Bayes multiplicity adjustment in the variable-selection problem.” *Ann. Statist.*, 38(5): 2587–2619.
URL <https://doi.org/10.1214/10-AOS792>
- Sibuya, M. (1960). “Bivariate extreme statistics, I.” *Annals of the Institute of Statistical Mathematics*, 11(3): 195–210.
URL <https://doi.org/10.1007/BF01682329>
- Smemo, S. A. (2012). “Regulation of heart development via transcriptional enhancers and epigenetic modifications.” Ph.D. thesis, University of Chicago.
- Smyth, G. K. (2004). “Linear models and empirical Bayes methods for assessing differential expression in microarray experiments.” *Statistical applications in genetics and molecular biology*, 3(1): 1–25.
- Stefanski, L. A. and Carroll, R. J. (1990). “Deconvolving kernel density estimators.” *Statistics*, 21(2): 169–184.
URL <https://doi.org/10.1080/02331889008802238>
- Stephens, M. (2017). “False discovery rates: a new deal.” *Biostatistics*, 18(2): 275–294.
- Storey, J. (2003). “The positive false discovery rate: a Bayesian interpretation and the q -value.” *Ann. Statist.*, 31(6): 2013–2035.
- Storey, J. D. (2002). “A direct approach to false discovery rates.” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 64(3): 479–498.
URL <https://rss.onlinelibrary.wiley.com/doi/abs/10.1111/1467-9868.00346>
- Sun, W. and Cai, T. T. (2009). “Large-Scale Multiple Testing under Dependence.” *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 71(2): 393–424.
URL <http://www.jstor.org/stable/40247580>
- Sun, Y., Zhang, N. R., and Owen, A. B. (2012). “Multiple hypothesis testing adjusted for latent variables, with an application to the AGEMAP gene expression data.” *Ann. Appl. Stat.*, 6(4): 1664–1688.
URL <https://doi.org/10.1214/12-A0AS561>
- Szegő, G. (1975). *Orthogonal Polynomials*. Providence, Rhode Island: American Mathematical Society, 4th edition.
- The GTEx Consortium (2015). “The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans.” *Science*, 348(6235): 648–660.
URL <http://science.sciencemag.org/content/348/6235/648>
- (2017). “Genetic effects on gene expression across human tissues.” *Nature*, 550: 204–213.
URL <https://www.nature.com/articles/nature24277>
- Thomas, D., Siemiatycki, J., Dewar, R., Robins, J., Goldberg, M., and Armstrong, B. (1985). “The problem of multiple inference in studies designed to generate hypotheses.” *American Journal of Epidemiology*, 122(6): 1080–1095.
- Tukey, J. W. (1991). “The Philosophy of Multiple Comparisons.” *Statist. Sci.*, 6(1): 100–116.
URL <https://doi.org/10.1214/ss/1177011945>
- Urbut, S. M., Wang, G., and Stephens, M. (2018). “Flexible statistical methods for estimating and testing effects in genomic studies with multiple conditions.” *Nature Genetics*, in press.
URL <https://www.nature.com/articles/s41588-018-0268-8>

- Wang, J., Zhao, Q., Hastie, T., and Owen, A. B. (2017). “Confounder adjustment in multiple hypothesis testing.” *Ann. Statist.*, 45(5): 1863–1894.
URL <https://doi.org/10.1214/16-AOS1511>
- Wang, W. and Stephens, M. (2018). “Empirical Bayes Matrix Factorization.” *ArXiv e-prints*.
- Wasserman, L. (2006). *All of Nonparametric Statistics (Springer Texts in Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc.
- Wu, W. B. (2008). “On false discovery control under dependence.” *Ann. Statist.*, 36(1): 364–380.
URL <https://doi.org/10.1214/009053607000000730>
- Xing, Z. and Stephens, M. (2016). “Smoothing via Adaptive Shrinkage (smash): denoising Poisson and heteroskedastic Gaussian signals.” *ArXiv e-prints*.