

# Shim and Stephens (2015) with HMT - Beta properties

Brendan Law

April 2, 2019

Just some properties of the  $\beta$  estimator under the HMT version of [?]. Refer to the supporting documents (the ones titled ‘EM\_algo\_Shim\_Stephens...’, etc) for some of the workings. Also, a lot of the closed form workings and derivations were done in the supplementary material to [?].

## 1 Properties of the base case - no HMT

### 1.1 $\beta$ Coefficient transform

This section is all about the  $\beta$  coefficient transform from the wavelet space to the data space. This is important, as the  $\beta$  coefficient is estimated in the wavelet space, and then transformed into the data space in order to provide interpretable estimates of effect sizes at different base locations along the genome.

It’s worth noting that there are also some practicalities involved with its estimation:

- Even though quantile regression is performed (to make the model less sensitive to extreme values in the dependent variable), this is not performed for the purposes of estimating effect sizes.
- Otherwise, the data  $\rightarrow$  wavelet  $\rightarrow$  inverse wavelet  $\rightarrow$  data space transform would be complicated by the fact that we would need to go data  $\rightarrow$  wavelet  $\rightarrow$  quantile transform  $\rightarrow$  inverse quantile transform  $\rightarrow$  inverse wavelet  $\rightarrow$  data, which is infeasible.
- Hence, estimation here occurs without a quantile transforms, and the translation of estimates from the wavelet back to the data space is a simple process due to the fact that the wavelet transform is *linear* and *one-to-one*.

### 1.2 Sketching out the wavelet linear transform

To represent the model in the wavelet space, we have that:

$$\mathbf{Y} = \mathbf{M} + \beta g + \mathbf{E}$$

So we have that:

$$\begin{bmatrix} y_{1,1}^1 & \cdots & y_{1,1}^n \\ \vdots & \ddots & \vdots \\ y_{J,L_J}^1 & \cdots & y_{J,L_J}^n \end{bmatrix} = \begin{bmatrix} \mu_{1,1}^1 & \cdots & \mu_{1,1}^n \\ \vdots & \ddots & \vdots \\ \mu_{J,L_J}^1 & \cdots & \mu_{J,L_J}^n \end{bmatrix} + \begin{bmatrix} \beta_{1,1} \\ \vdots \\ \beta_{J,L_J} \end{bmatrix} [g^1 \quad \cdots \quad g^n] + \begin{bmatrix} \epsilon_{1,1}^1 & \cdots & \epsilon_{1,1}^n \\ \vdots & \ddots & \vdots \\ \epsilon_{J,L_J}^1 & \cdots & \epsilon_{J,L_J}^n \end{bmatrix}$$

Where the  $\mathbf{Y}, \mathbf{M}, \mathbf{E}$  matrices are all of dimension  $B \times n$ . This represents the  $B$  base pair locations (which also correspond to the number of wavelet coefficients the transform creates –  $B/2$  at the lowest/finest level, all the way to 1 at the highest/coarsest level), and the  $n$  individuals in the dataset. The  $\beta$  vector is  $B \times 1$ , and  $g$  is  $1 \times n$ , with  $\beta g$  becoming a  $B \times n$  matrix:

$$\begin{bmatrix} \beta_{1,1}g^1 & \cdots & \beta_{1,1}g^n \\ \vdots & \ddots & \vdots \\ \beta_{J,L_J}g^1 & \cdots & \beta_{J,L_J}g^n \end{bmatrix}$$

Now, to represent the linear transform from the data space to the wavelet space, we note that:

$$\begin{aligned} y &= Wd \\ \Rightarrow d &= W^{-1}y \end{aligned}$$

$$\begin{aligned} D &= W^{-1}Y \\ \Rightarrow D &= W^{-1}M + W^{-1}\beta g + W^{-1}E, \end{aligned}$$

where  $W^{-1}\beta := \alpha$

The original data,  $D$ , is a  $B \times n$  matrix as follows:

$$D = \begin{bmatrix} d_1^1 & \dots & d_1^n \\ \vdots & \ddots & \vdots \\ d_B^1 & \dots & d_B^n \end{bmatrix}$$

The wavelet transform represents a linear transform from this  $B \times n$  matrix,  $D$  into a  $B \times n$  matrix of coefficients,  $Y$ , and hence can be represented as a  $B \times B$  matrix:

$$W = \begin{bmatrix} w_{1,(1,1)} & \dots & w_{B,(1,1)} \\ \vdots & \ddots & \vdots \\ w_{1,(J,L_J)} & \dots & w_{B,(J,L_J)} \end{bmatrix}$$

where  $w_{ij}$  represents the 'amount' by which the data point at base  $i$  contributes to the wavelet coefficient indexed by  $j$ , where  $j$  represents a scale-location pair,  $(s, l)$ . This is for any given individual. Furthermore, we have that  $W$  is orthogonal (from suppl mtl) hence  $W^T = W^{-1}$

Therefore for the matrix  $\alpha$ :

$$\begin{aligned} \alpha &= W^{-1}\beta \\ &= W^T\beta \\ \Rightarrow \alpha &= \begin{bmatrix} w_{1,(1,1)} & \dots & w_{1,(J,L_J)} \\ \vdots & \ddots & \vdots \\ w_{B,(1,1)} & \dots & w_{B,(J,L_J)} \end{bmatrix} \begin{bmatrix} \beta_{1,1} \\ \vdots \\ \beta_{J,L_J} \end{bmatrix} \\ &= \begin{bmatrix} \sum_{s,l} w_{1,(s,l)}\beta_{s,l} \\ \vdots \\ \sum_{s,l} w_{B,(s,l)}\beta_{s,l} \end{bmatrix} = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_B \end{bmatrix} \end{aligned}$$

which yields us the  $\alpha_i$ 's – the estimated data space effect at each base location  $i$ .

### 1.3 Estimating statistical properties of $\alpha$

We're also interested in some statistical properties of  $\alpha$ , such as its mean and variance. This requires us to perform some operations on the posterior distribution of  $\beta$ , (i.e.  $P(\beta \mid \text{Data, hyper-params})$ , with all other parameters integrated out). In the no-HMT case, we have from the supplementary material, a closed form of the posterior of  $\beta$ , and thus we can derive its mean and variance properties, and use inverse wavelet transforms to derive means and variances of  $\alpha$ .

#### 1.3.1 Posterior distribution

From the no-HMT supplementary material, we have a posterior of  $\beta_{s,l} = P(\beta_{s,l} \mid y_{s,l}, g)$  as follows:

$$\begin{aligned} P(\beta_{s,l} \mid y_{s,l}, g, \pi) &= P(\beta_{s,l} \mid y_{s,l}, g, \pi, \gamma_{s,l} = 1)P(\gamma_{s,l} = 1 \mid y_{s,l}, g, \pi) + P(\beta_{s,l} \mid y_{s,l}, g, \pi, \gamma_{s,l} = 0)P(\gamma_{s,l} = 0 \mid y_{s,l}, g, \pi) \\ &= P(\beta_{s,l} \mid y_{s,l}, g, \pi, \gamma_{s,l} = 1)\phi_{s,l} + P(\beta_{s,l} \mid y_{s,l}, g, \pi, \gamma_{s,l} = 0)(1 - \phi_{s,l}) \end{aligned}$$

where:

- $\boldsymbol{\pi}$  is a vector of hyperparameters (in this case, the proportion of non-zero wavelet coefficients at a given scale)
- $P(\beta_{s,l} \mid y_{s,l}, g, \boldsymbol{\pi}, \gamma_{s,l} = 0)$  is a point mass at 0 (wp 1)
- 

$$\begin{aligned}
\phi_{s,l} &= P(\gamma_{s,l} = 1 \mid y_{s,l}, g, \boldsymbol{\pi}) \\
&= \frac{P(y_{s,l} \mid g, \gamma_{s,l} = 1)P(\gamma_{s,l} = 1 \mid \boldsymbol{\pi})}{P(y_{s,l} \mid g, \gamma_{s,l} = 0)P(\gamma_{s,l} = 0 \mid \boldsymbol{\pi}) + P(y_{s,l} \mid g, \gamma_{s,l} = 1)P(\gamma_{s,l} = 1 \mid \boldsymbol{\pi})} \\
&= \frac{BF_{s,l}(y, g)\hat{\pi}_s}{BF_{s,l}(y, g)\hat{\pi}_s + (1 - \hat{\pi}_s)}, \text{ obtaining the BFs by dividing through by the } \gamma_{s,l} = 0 \text{ case}
\end{aligned}$$

- Posterior of  $\beta_{s,l}$  is a mixture (with mixing probabilities  $\phi_{s,l}$ ) of a three-parameter t-distribution, and a point mass at 0.

### 1.3.2 Mean

The desired quantity is  $\mathbb{E}[\alpha_b]$ . Based on the ‘Ash’ material and ‘Wavelets - Supplementary Material’. We start by finding  $\mathbb{E}[\beta_{s,l}]$  and then transforming it back into the data space.

$$\begin{aligned}
\mathbb{E}[\beta_{s,l}] &= \mathbb{E}[\mathbb{E}[\beta_{s,l} \mid \gamma_{s,l}]] \\
&= \mathbb{E}[P(\gamma_{s,l} = 1)\mathbb{E}[\beta_{s,l} \mid \gamma_{s,l} = 1]] + \mathbb{E}[P(\gamma_{s,l} = 0)\mathbb{E}[\beta_{s,l} \mid \gamma_{s,l} = 0]] \\
&= \mathbb{E}[\phi_{s,l}a_{s,l} + (1 - \phi_{s,l})0] \\
&= \phi_{s,l}a_{s,l}
\end{aligned}$$

where  $a_{s,l}$  represents the mean of a three-parameter t-distribution, and in this case, corresponds to  $\mathbf{B}_2^*$ , the 2nd element of  $\mathbf{B} = (\mathbf{D}^{-1} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ , evaluated in the limit (of some of the hyper-parameters,  $\kappa^a, \kappa^b$  and  $\sigma_\mu^2$ ).

$$\begin{aligned}
\therefore \mathbb{E}[\alpha_b] &= \mathbb{E}\left[\sum_{s,l} w_{b,(s,l)} \beta_{s,l}\right] \\
&= \sum_{s,l} w_{b,(s,l)} \mathbb{E}[\beta_{s,l}] && \text{(linearity of expectations)} \\
&= \sum_{s,l} \phi_{s,l} a_{s,l} w_{b,(s,l)}
\end{aligned}$$

### 1.3.3 Variance

$$\begin{aligned}
Var(\beta_{s,l}) &= \mathbb{E}[\beta_{s,l}^2] - \mathbb{E}[\beta_{s,l}]^2 \\
\mathbb{E}[\beta_{s,l}^2] &= \mathbb{E}[\mathbb{E}[\beta_{s,l}^2 \mid \gamma_{s,l}]] \\
\mathbb{E}[\beta_{s,l}^2 \mid \gamma_{s,l}] &= Var[\beta_{s,l} \mid \gamma_{s,l}] + \mathbb{E}[\beta_{s,l} \mid \gamma_{s,l}]^2 \\
\therefore \mathbb{E}[\mathbb{E}[\beta_{s,l}^2 \mid \gamma_{s,l}]] &= \mathbb{E}[\phi_{s,l}(Var[\beta_{s,l} \mid \gamma_{s,l} = 1] + \mathbb{E}[\beta_{s,l} \mid \gamma_{s,l} = 1]^2)] \\
&= \mathbb{E}[\phi_{s,l}(\frac{b_{s,l}v_{s,l}}{v_{s,l} - 2} + a_{s,l}^2)] \\
&= \phi_{s,l}(\frac{b_{s,l}v_{s,l}}{v_{s,l} - 2} + a_{s,l}^2) \\
\therefore Var(\beta_{s,l}) &= \phi_{s,l}(\frac{b_{s,l}v_{s,l}}{v_{s,l} - 2} + a_{s,l}^2) - \phi_{s,l}^2 a_{s,l}^2 = \phi_{s,l} \left( \frac{b_{s,l}v_{s,l}}{v_{s,l} - 2} + a_{s,l}^2 - \phi_{s,l} a_{s,l}^2 \right) \\
\therefore Var(\alpha_b) &= Var(\sum_{s,l} w_{b,(s,l)} \beta_{s,l}) \\
&= \sum_{s,l} w_{b,(s,l)}^2 Var(\beta_{s,l}) \quad (\text{linearity of variance of independent terms}) \\
&= \sum_{s,l} w_{b,(s,l)}^2 \phi_{s,l} \left( \frac{b_{s,l}v_{s,l}}{v_{s,l} - 2} + a_{s,l}^2 - \phi_{s,l} a_{s,l}^2 \right)
\end{aligned}$$

We note here that  $\beta_{s,l}$ 's are (conditionally) independent, given the data and hyperparameters (notation above has been very sloppy – all properties of expectations and variances are of the posterior of the beta, conditional on data and hyperparameters. It has been omitted from the above notation due to brevity. Hence  $Var(\beta_{s,l}) = Var(\beta_{s,l} \mid \text{Data, hyperparams})$ , which are conditionally independent of each other). An explanation of why the  $\beta$  posteriors are conditionally independent can be found in ‘wavelets\_and\_models.pdf’.

## 2 Properties of the base case - with HMT

### 2.1 Estimating statistical properties of $\alpha$

All of the above still applies, except that we now calculate  $\phi_{s,l}$  differently, as  $\phi_{s,l}$  depends on the probabilities of the other  $\gamma_{s,l}$ 's as seen through the upward-downward algorithm required to solve this case.

#### 2.1.1 Posterior distribution

$$\begin{aligned}
\phi_{s,l} &= P(\gamma_{s,l} = 1 \mid \mathbf{Y}, g, \boldsymbol{\pi}) \\
&= \frac{P(\gamma_{s,l} = 1, \mathbf{Y} \mid \boldsymbol{\pi})}{P(\mathbf{Y} \mid \boldsymbol{\pi})} \\
&= \frac{\beta_{s,l}(1)\alpha_{s,l}(1)}{\sum_{m=0}^1 \beta_{s,l}(m)\alpha_{s,l}(m)}
\end{aligned}$$

where  $\beta_{s,l}(m) := P(\mathcal{T}_{s,l} \mid \gamma_{s,l} = m, \boldsymbol{\pi})$  and  $\alpha_{s,l}(m) := P(\gamma_{s,l} = m, \mathcal{T}_{1 \setminus s,l} \mid \boldsymbol{\pi})$ , as found in the other documentation (‘EM\_algo.HMT.pdf’). These are computed as part of the EM algorithm, and will drop out of it nicely (all ‘scaled’ to utilise the closed form of the Bayes Factor, but still with the desired result).

#### 2.1.2 Means and variances

Can't remember how the means work. Don't believe the posterior still has the same 3-parameter t-distribution as the  $\beta$  priors are no longer independent and factorisable (HMT prior imposes dependencies between them). Although we do suspect that  $P(\beta_{s,l} \mid \mathbf{Y}, g, \gamma_{s,l} = m) = P(\beta_{s,l} \mid y_{s,l}, g, \gamma_{s,l} = m)$ . We'll have to calculate this one later. How did we get to this conclusion, and why is it significant? Does it mean that we postulate our joint posterior factorises into the

product of individual 3-parameter t-distributions?

A note about the variances:

- Dependence (even conditionally) between  $\beta_{s,l}$ 's, meaning that the linearity and additivity of variances doesn't apply as there will be covariance cross-terms to be calculated.
- Hence we'll just simulate these properties in this case.
- It's also likely we'll end up simulating the means too.

### 2.1.3 Simulation methodology