

Statistical framework from Flutre, Wen, Pritchard and Stephens (PLoS Genetics, 2013): Halfway to understanding the Bayes Factors

Sarah Urbut

July 30, 2014

Contents

1 Likelihood of the whole data set	1
2 Focus on a single gene-SNP pair	1
3 Answers!	4

This document describes the statistical framework with more details and sometimes a slightly different notation, notably inspired by Wen & Stephens (Annals of Applied Statistics, 2014) and Wen (Biometrics, 2014).

1 Likelihood of the whole data set

2 Focus on a single gene-SNP pair

The likelihood for gene g and SNP p is:

$$Y_g|X_p, B_{gp}, X_c, B_{gc}, \Sigma_{gp} \sim \mathcal{N}_{N \times R}(X_p B_{gp} + X_c B_{gc}, I_N, \Sigma_{gp}) \quad (1)$$

where:

- Y_g is the $N \times R$ matrix of expression levels;
- X_p is the $N \times 1$ matrix of genotypes (assuming the same individuals in all tissues);
- B_{gp} is the unknown $1 \times R$ matrix of genotype effect sizes;
- X_c is the $N \times (1 + Q)$ matrix of known covariates (including a column of 1's for the intercepts);
- B_{gc} is the unknown $(1 + Q) \times R$ matrix of covariate effect sizes (including the μ_s);
- $\mathcal{N}_{N \times R}$ is the matrix Normal distribution;

- Σ_{gp} is the unknown $R \times R$ covariance matrix of the errors.

For mathematical convenience (especially in the case of multiple SNPs), we vectorize the rows of B_{gp} into β_{gp} . Here, as we focus on one SNP at a time, we directly have $\beta_{gp} = B_{gp}^T$.

The conditional posterior of B: $P(B|Y, X, \tau) = \frac{P(B, Y|X, \tau)}{P(Y|X, \tau)}$

Let's neglect the normalization constant for now. We note that if we expand this expression in full, we get the following:

$$\Pr(\mathbf{b}_s | \mathbf{Y}_s, \mathbf{X}_s, \Phi_s^{-1}) \propto \exp((\mathbf{b}_s - \bar{\mathbf{b}})^t \Phi_s^{-1} (\mathbf{b}_s - \bar{\mathbf{b}})) \exp((\mathbf{Y}_s - \mathbf{X}_s \mathbf{b}_s)^t (\mathbf{Y}_s - \mathbf{X}_s \mathbf{b}_s)) \quad (2)$$

Taking terms out of the exponent and distributing terms, we arrive at:

$$\Pr(\mathbf{b}_s | \mathbf{Y}_s, \mathbf{X}_s, \Phi_s^{-1}) \propto \mathbf{b}_s^t \Phi_s^{-1} \mathbf{b}_s - \mathbf{b}_s^t \Phi_s^{-1} \bar{\mathbf{b}} - \bar{\mathbf{b}}^t \Phi_s^{-1} \mathbf{b}_s + \bar{\mathbf{b}}^t \Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{Y}_s^t \mathbf{Y}_s - (\mathbf{X}_s \mathbf{b}_s)^t \mathbf{Y}_s - \mathbf{Y}_s^t \mathbf{X}_s \mathbf{b}_s + (\mathbf{X}_s \mathbf{b}_s)^t (\mathbf{X}_s \mathbf{b}_s) \quad (3)$$

We will first leave out the terms that don't include \mathbf{b}_s , i.e., $\mathbf{Y}_s^t \mathbf{Y}_s$ and $\bar{\mathbf{b}}^t \Phi_s^{-1} \bar{\mathbf{b}}$, and group some terms:

$$\Pr(\mathbf{b}_s | \mathbf{Y}_s, \mathbf{X}_s, \Phi_s^{-1}) \propto \mathbf{b}_s^t (\Phi_s^{-1} + \mathbf{X}_s^t \mathbf{X}_s) \mathbf{b}_s - \mathbf{b}_s^t (\Phi_s^{-1} \bar{\mathbf{b}} - \mathbf{X}_s^t \mathbf{Y}_s)^t - (\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s)^t \mathbf{b}_s \quad (4)$$

In order to aid our completion of the square, let's add a term that doesn't contain \mathbf{b}_s which is a legitimate

$$(\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s)^t (\Phi_s^{-1} - \mathbf{X}_s^t \mathbf{X}_s) (\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s) \quad (5)$$

Now, let's define Ω_s^{-1} as: $(\Phi_s^{-1} - \mathbf{X}_s^t \mathbf{X}_s)$ and μ_s as $(\Phi_s^{-1} - \mathbf{X}_s^t \mathbf{X}_s) (\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s)$

Then it becomes apparent that we can rewrite $\Pr(\mathbf{b}_s | \mathbf{Y}_s, \mathbf{X}_s, \Phi_s^{-1})$ as:

$$\Pr(\mathbf{b}_s | \mathbf{Y}_s, \mathbf{X}_s, \Phi_s^{-1}) \propto (\mathbf{b}_s - \mu_s)^t \Omega_s^{-1} (\mathbf{b}_s - \mu_s) \quad (6)$$

So that when we compute the marginal probability of Y:

$$\Pr(\mathbf{Y}_s | \mathbf{X}_s) = \frac{\Pr(\mathbf{b}_s | \mathbf{Y}_s, \mathbf{X}_s, \tau) \Pr(\mathbf{b}_s | \tau, \bar{\mathbf{b}})}{\Pr(\mathbf{b}_s | \mathbf{Y}_s, \mathbf{X}_s, \tau)} \quad (7)$$

It is obvious that the numerator will contain the two terms we neglected (because they didn't contain \mathbf{b}_s) in (3) and the denominator will contain the term we added in (4), $(\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s)^t (\Phi_s^{-1} - \mathbf{X}_s^t \mathbf{X}_s) (\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s)$.

Thus the marginal likelihood is:

$$\Pr(\mathbf{Y}_s | \mathbf{X}_s) = \frac{2\pi^{-n_s/2}}{\tau_s} \left| \Phi_s^{-1} \right|^{\frac{-1}{2}} \left| (\Phi_s^{-1} + \mathbf{X}_s^t \mathbf{X}_s) \right|^{\frac{-1}{2}} * \exp\left(\frac{1}{2} (\mathbf{Y}_s^t \mathbf{Y}_s - (\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s)^t (\Phi_s^{-1} - \mathbf{X}_s^t \mathbf{X}_s)^{-1} (\Phi_s^{-1} \bar{\mathbf{b}} + \mathbf{X}_s^t \mathbf{Y}_s) - \bar{\mathbf{b}}^t \Phi_s^{-1} \bar{\mathbf{b}})\right) \quad (8)$$

When we integrate over $\bar{\mathbf{b}}$, we could have just done the integral over \mathbf{b}_s rather than $\mathbf{b}_s \bar{\mathbf{b}}$ which would be $\mathbf{b}_s \sim (\mathbf{0}, \mathbf{W})$, where \mathbf{W} is the prior matrix of covariance effects, with $\phi + \omega$ on the diagonal and ω on the off-diagonal. This is akin to what is done in the biometric paper, where:

$$\begin{aligned}
p(\mathbf{b}_s | \mathbf{X}_s 1) &\propto \exp\left(\frac{-1}{2}[(\mathbf{b}_s - \hat{\beta}_s)^t \mathbf{V}^{-1}(\mathbf{b}_s - \hat{\beta}_s) + (\mathbf{b}_s^t \mathbf{W}^{-1} \mathbf{b}_s)]\right) \\
&\propto \exp\left(\frac{-1}{2}[(\mathbf{b}_s^t (\mathbf{V}^{-1} + \mathbf{W}^{-1}) \mathbf{b}_s)] - (\mathbf{b}_s^t (\mathbf{V}^{-1} \hat{\beta}_s) + (\hat{\beta}_s^t \mathbf{V}^{-1} \mathbf{b}_s))\right) (9)
\end{aligned}$$

And under the Null Hypothesis,

$$\begin{aligned}
p(\mathbf{b}_s | \mathbf{X}_s 0) &\propto \exp\left(\frac{-1}{2}[(\mathbf{b}_s - \hat{\beta}_s)^t \mathbf{V}^{-1}(\mathbf{b}_s - \hat{\beta}_s) + (\mathbf{b}_s^t \mathbf{W}^{-1} \mathbf{b}_s)]\right) \\
p(\mathbf{b}_s | \mathbf{X}_s 0) &\propto \exp\left(\frac{-1}{2}[(\mathbf{b}_s - \hat{\beta}_s)^t \mathbf{V}^{-1}(\mathbf{b}_s - \hat{\beta}_s) + (\mathbf{b}_s^t \mathbf{0} \mathbf{b}_s)]\right) \\
&\propto \exp\left(\frac{-1}{2}[(\mathbf{b}_s^t (\mathbf{V}^{-1}) \mathbf{b}_s)] - (\mathbf{b}_s^t (\mathbf{V}^{-1} \hat{\beta}_s) + (\hat{\beta}_s^t \mathbf{V}^{-1} \mathbf{b}_s))\right) (10)
\end{aligned}$$

Analogous to the previous situation, we have omitted an additional term contained in the likelihood of $\Pr(\hat{\beta}_s | \mathbf{X}_s, \mathbf{E})$ because it didn't contain \mathbf{b}_s , $\hat{\beta}_s^t \mathbf{V}^{-1} \hat{\beta}_s$ and we have added an additional term to only the alternative model:

$$\mathbf{b}_s^t (\mathbf{V}^{-1} \hat{\beta}_s) (\mathbf{V}^{-1} + \mathbf{W}^{-1}) (\hat{\beta}_s^t \mathbf{V}^{-1} \mathbf{b}_s) \quad (11)$$

Crucially, in computing the Bayes Factor

$$\frac{\Pr(\hat{\beta}_s | \mathbf{X}_s, \mathbf{E}, M1)}{\Pr(\hat{\beta}_s | \mathbf{X}_s, \mathbf{E}, M0)} \quad (12)$$

Only the added term is present in the numerator and not the denominator, and so we can see that the Bayes Factor will be proportional to:

$$\exp(\mathbf{b}_s^t (\mathbf{V}^{-1} \hat{\beta}_s) (\mathbf{V}^{-1} + \mathbf{W}^{-1}) (\hat{\beta}_s^t \mathbf{V}^{-1} \mathbf{b}_s)) \quad (13)$$

which, intuitively, augments the likelihood of the data by the prior covariance matrix of the effects, incorporating prior belief about shared effect and thus increasing our power to detect such homogeneity. This results from the fact that \mathbf{b}_s , $\hat{\beta}_s^t \mathbf{V}^{-1} \hat{\beta}_s$ is present in both the alternative and null models, but $\exp(\mathbf{b}_s^t (\mathbf{V}^{-1} \hat{\beta}_s) (\mathbf{V}^{-1} + \mathbf{W}^{-1}) (\hat{\beta}_s^t \mathbf{V}^{-1} \mathbf{b}_s))$ is added only in the alternative models, where there exists a distribution of non-zero values for \mathbf{b}_s . Furthermore, if we imagine that the likelihood function will have most of its density at $\mathbf{b}_s = \hat{\beta}_s$ and thus the the marginal probability $\Pr(\hat{\beta}_s | \mathbf{X}_s, \mathbf{E}, M1)$ will really depend on the prior probability assigned to \mathbf{b}_s at $\mathbf{b}_s = \hat{\beta}_s$, the use of a prior which assumes and thus upweights values of \mathbf{b}_s near a shared 'mean' effect size will be of benefit if such a shared effect size exists.

- 1. In the simplest case, assuming the τ is known, if we consider as the product of all marginal likelihoods, then how are we really exploiting the covariance in effects? I still don't see it – the \mathbf{Y} vectors are $n \times 1$, and so the covariance in effects in the off diagonal terms is not in the ABF_{random}

This is only because we integrate over \mathbf{b}_s first - if we had integrated over $\bar{\mathbf{b}}$ first, it would not have been so. We would have brought everything back into 1, and then it becomes factored again because the residuals are uncorrelated.

- 2. We can use a product because we assume the vectors of residual errors are independent across populations, but the effects are not (in any case by the max H case). So is the product because of the conditional exchangeability?

Yes!

- 3. I see that we are pooling the Bayes Factors by the product, which I guess increases our sample size (by sheer more terms in our product) but even if the residuals are independent, we miss the covariance in effects inherent in the data

False, this is accounted for by the prior covariance in effects which is captured in the shared mean — thus though the residuals are independent, the covariance in effects implies that the effects are not.

- 4. Perhaps the product refers to conditional on the mean $\bar{\mathbf{b}}$, they are exchangeable (i.e., independent) but then we are replying only on the BF(fixed) to reinforce our covariance in effects.

Yes – this is true only in the met analysis case. The exchangeability allows us to factor the likelihood with \mathbf{b}_s conditional on $\bar{\mathbf{b}}$ because each \mathbf{b}_s is i.i.d. conditional on $\bar{\mathbf{b}}$ (see DeFinetti's theorem); however, when we return to integrate $\bar{\mathbf{b}}$ out, the product no longer factorizes. In the met analysis case where the τ s are independent, integrating out τ returns the factorization; however, if we make no assumption about the independence of residual variances, then we must consider the distribution of the Bayes Factor of a multivariate vectors (i.e., the product of $n \times S$ vectors rather than $n \times s$ vectors).

From William's paper:

The off diagonal of matrix \mathbf{W}_g defines context-dependent prior correlation between non-zero regression coefficients. Incorporating this information enables borrowing strength across correlated components in β_g , thereby improving the efficiency of model selection.

But if we use a product over all \mathbf{b}_s , then how are we incorporating the off-diagonal elements? Is it just through the fixed effect term on the left, or does the term on the right incorporate the greater probability that \mathbf{b}_s arises from a shared effects mode? I don't think it can because it doesn't pay any attention to other data. Maybe if I could see where the prior on shared effect came into the Bayes factor ...

Also: Kass LaPlace Approximation: I see that the majority of the Mass of the BF will fall at the MLE, but don't we still need to weight by the prior probability of that parameter?

3 Answers!

- 2. We can use a product because we assume the vectors of residual errors are independent across populations, but the effects are not (in any case by the max H case). So is the product because of the conditional exchangeability?

Yes! See above explanations!