

1 **WHEN ARE SEX-SPECIFIC EFFECTS REALLY SEX-SPECIFIC?**

2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27

Eunice H. Chin\* and Julian K. Christians

*Department of Biological Sciences, Simon Fraser University, 8888 University Dr.,  
Burnaby, B.C., V5A 1S6, Canada*

Word Count: 2210

\* To whom all correspondence should be addressed:

Eunice H. Chin,  
Department of Biological Sciences,  
Simon Fraser University,  
8888 University Dr.,  
Burnaby, B.C,  
V5A 1S6 Canada  
Phone: (778) 782-7398  
Email: [echin@sfu.ca](mailto:echin@sfu.ca)

28 **ABSTRACT**

29 We examined developmental programming studies that reported sex-specific effects  
30 published between 2012 and 2014, and examined whether the authors reported a  
31 statistical approach to explicitly test whether the effect of treatment differed between the  
32 sexes, e.g., a sex by treatment interaction term. Less than half of the studies that reported  
33 sex-specific effects described explicitly testing whether effects were indeed sex-specific;  
34 in most cases, an effect was considered “sex-specific” if it was significant in one sex but  
35 not the other. This is not a robust approach, since significance in one sex and lack of  
36 significance in the other sex does not imply a significant difference between the sexes.  
37 However, sample size often limits statistical power to detect interactions. We suggest  
38 that if the effect is significant in only one sex, but the interaction term is not significant,  
39 alternative solutions would be to present the confidence intervals for the effect size for  
40 each sex, or to use Bayesian approaches to calculate the probability that the effect sizes  
41 differ between the sexes. We present a simple example of a Bayesian analysis to  
42 illustrate that this approach is reasonably easy to implement and interpret.

43

44 Key words: sex-specific effects, developmental programming

45 **INTRODUCTION**

46 Recently, there has been increasing interest in sex-specific effects of  
47 developmental programming<sup>1,2</sup>. Literature searches on PubMed and Web of Science  
48 from 2000 to 2014 show that there has been a substantial increase in studies examining  
49 sex-specific effects of developmental programming – searches for the terms “sex-specific  
50 or sex-dependent” and “fetal or development” and “programming” yielded only 14  
51 studies published in 2000 and 84 studies published in 2014. Differential susceptibility to  
52 developmental programming between the sexes has been demonstrated widely in animal  
53 models<sup>2</sup>. Moreover, funding agencies such as National Institutes of Health (NIH) and  
54 Canadian Institutes for Health Research (CIHR) are increasingly encouraging the  
55 incorporation of gender and sex into research designs where appropriate.

56 In assessing whether effects are sex-specific, it is necessary to explicitly test  
57 whether an effect differs between the sexes; performing analyses separately for each sex  
58 is not sufficient. For example, if the effect of treatment is significantly different from  
59 zero in males, but not in females, it does not necessarily follow that the effect size in  
60 males is significantly different from the effect size in females. A non-significant effect in  
61 females is not evidence that the effect size in females is actually zero. Conversely, if the  
62 effect of treatment is significant and in the same direction in both sexes, it does not  
63 necessarily follow that the magnitude of the effect of treatment is equal in both sexes .  
64 We argue that, in order to report an effect as “sex-specific” or “sex-dependent”, it is  
65 necessary to perform an explicit test of the difference in effect size between the sexes,  
66 e.g., using a sex by treatment interaction term in the statistical model. We examined the  
67 reporting of statistical models in experiments that reported sex-specific effects of

68 developmental programming, and whether an explicit statistical test of differences in  
69 effect between the sexes was described.

70

## 71 **METHODS**

### 72 *Papers*

73 We searched the ISI Web of Knowledge (<http://www.isiwebofknowledge.com>)  
74 and PubMed (<http://www.ncbi.nlm.nih.gov/pubmed>) databases for papers published in  
75 2012 through 2014, using a combination of the search terms “sex-specific or sex-  
76 dependent”, “fetal or development or gestational”, and “programming”. To be included  
77 in our survey, studies had to involve a maternal treatment during pregnancy, followed by  
78 a measurement of offspring physiology. We assessed whether the authors performed an  
79 explicit statistical test of differences in effect between the sexes, e.g., a sex by treatment  
80 interaction term.

81

## 82 **RESULTS**

83 In total, there were 34 rodent studies that met the search criteria from 2012 to  
84 2014, and of these, 31 reported sex-specific effects. Only 11 studies reported a sex by  
85 treatment interaction, utilizing general linear models, while the remaining 20 studies  
86 analyzed the sexes separately. Studies that analyzed the sexes separately utilized t-tests,  
87 Mann-Whitney T-tests, Newman-Keul’s tests and Fisher’s tests. All 9 studies of other  
88 organisms, including sheep, pigs and baboons, reported sex-specific effects. However,  
89 only 3 studies reported a sex by treatment interaction term, utilizing general linear models

90 while the others analyzed the sexes separately using ANOVAs followed by Bonferroni or  
91 Dunnet's post-hoc tests.

92

### 93 **DISCUSSION**

94 Only a third of studies reporting sex differences in developmental programming  
95 effects included a sex by treatment interaction term in their statistical models. Instead the  
96 sexes were analyzed separately, and significant effects detected in only one sex were  
97 described as sex-specific. However, this is not a robust approach, and studies in  
98 developmental programming should include explicit statistical tests for sex-specific  
99 effects, e.g., by including interaction terms between sex and treatment.

100 While small sample sizes may limit the statistical power to detect interactions, this  
101 is not a justification for excluding explicit statistical tests for sex-specific effects. Where  
102 the sex by treatment interaction term is not significant, an alternative approach would be  
103 to report confidence intervals for the magnitude of effect in each sex. While confidence  
104 intervals would overlap between the sexes if the interaction was not significant, this  
105 approach would illustrate the potential magnitude of differences between the sexes.

106 Another alternative approach would be the use of Bayesian analysis. Traditional  
107 null hypothesis significant testing (NHST) uses P-values, which describe the probability  
108 of observing the test statistic or one more extreme if the null hypothesis is true (e.g., there  
109 is no effect of treatment). In contrast, Bayesian analysis quantify degree of belief or  
110 uncertainty in a parameter or hypothesis using probability distributions<sup>3</sup>. Thus, while  
111 NHST calculates describe the probability of observing the test statistic or one more  
112 extreme, given some null hypothesis, Bayesian approaches calculate the probability of

113 some hypothesis or parameter value, given the data<sup>4</sup>. Therefore, Bayesian approaches  
114 actually allow more intuitive statements to be made than does NHST. While Bayesian  
115 approaches provide the opportunity to incorporate prior knowledge into analyses, it is  
116 possible to use prior distributions that are uninformative, i.e., that have little impact on  
117 the results<sup>3</sup>. Bayesian analysis begins with a prior probability distribution (e.g., for  
118 parameter values) and uses the prior and sample data to produce a posterior distribution,  
119 from which one can determine the probability of some parameter value, given the data<sup>5</sup>  
120 (see Appendix 1 for more detail).

121

122 In the Appendix, we provide an example of a Bayesian approach to express the  
123 level of confidence that effect size differs between males and females, using both SAS  
124 and an open-source package available for R<sup>6</sup>. This approach focuses not on whether or  
125 not the difference in effect size between males and females is statistically significant (i.e.,  
126 different from zero), but on the level of confidence that this difference is sufficiently  
127 large to be biologically interesting, which is arguably a more important issue. Our  
128 example illustrates the flexibility of Bayesian analysis, which allows one to calculate the  
129 probability of a customized parameter value or hypothesis. Furthermore, while Bayesian  
130 approaches may currently be unfamiliar to many biologists and clinicians, the analyses  
131 that we have described are very easy to implement. This approach provides probabilities  
132 that are straightforward to interpret, but does not provide a clear-cut, “significant or not”  
133 answer in the way that a P-value might. However, NHST does not necessarily provide  
134 less ambiguous results, e.g., a non-significant P-value does not provide evidence that an  
135 effect does not exist, and conversely a significant P-value does not necessarily indicate

136 that an effect is biologically important. Furthermore, because NHST reduces the results to  
137 a dichotomy (significant or not), with small sample sizes one is more likely accept the  
138 null hypothesis even if it is false, i.e., commit a type II error. With small sample sizes,  
139 there will be more uncertainty in Bayesian estimates<sup>7</sup>, but there will not be a greater  
140 chance of error. Therefore, even if sample size limits the statistical power to detect a sex  
141 by treatment interaction term, Bayesian approaches allow researchers in developmental  
142 programming to assess the probability of sex-specific effects.

143

#### 144 **ACKNOWLEDGEMENTS**

145 We thank Richard Morey, Tim Swartz and Ian Bercovitz for advice regarding the  
146 interpretation of Bayesian analyses and William Tian for help with the literature search.

147

#### 148 **REFERENCES**

- 149 1. Dunn GA, Morgan CP, Bale TL. Sex-specificity in transgenerational epigenetic  
150 programming. *Hormones and Behavior*. 2011;59(3), 290-295.
- 151 2. Aiken CE, Ozanne SE. Sex differences in developmental programming models.  
152 *Reproduction*. 2013;145, R1-R13.
- 153 3. Wagenmakers EJ. A practical solution to the pervasive problems of p values.  
154 *Psychonomic Bulletin and Review*. 2007;14(5), 779-804.
- 155 4. Stephens PA, Buskirk SW, Martínez del Rio C. Inference in ecology and  
156 evolution. 22. 2007;4(192-197).
- 157 5. Eddy SR. What is Bayesian statistics? *Nature Biotechnology*. 2004;22(9), 1177-  
158 1178.
- 159 6. R-Core, Team. R: A language and environment for statistical computing. Vienna,  
160 Austria2015; Available from: <http://www.R-project.org/>.
- 161 7. Kruschke JK. Bayesian estimation supersedes the t Test. *Journal of Experimental*  
162 *Psychology: General*. 2013;142(2), 573-603.

163

164

165 Appendix 1. Example of a Bayesian approach to express the level of confidence that  
166 effect size differs between males and females.

167 To illustrate a Bayesian approach, we have simulated a data set where a trait has  
168 been measured in male and female individuals subjected to one of two treatments (Table  
169 A1). When analyzing the sexes separately using a one-way ANOVA, there is a  
170 significant effect of treatment on the trait in males ( $F_{1,18} = 5.19$ ;  $P = 0.04$ ), but not in  
171 females ( $F_{1,18} = 0.03$ ;  $P = 0.86$ ). However, when analyzing the sexes together using a  
172 two-way ANOVA and including a sex by treatment interaction term, the interaction is  
173 marginally non-significant ( $F_{1,36} = 3.46$ ;  $P = 0.07$ ). This is a situation in which many  
174 authors would be inclined to report sex-specific effects, and yet we have argued that sex-  
175 specific effects should not be reported unless there is a significant sex by treatment  
176 interaction term. A Bayesian analysis provides an alternative way to express the level of  
177 confidence that effect size differs between males and females.

178 Before considering a Bayesian approach, it is necessary to consider what is tested  
179 by including a sex by treatment interaction term in a model, i.e., whether the effect of  
180 treatment differs between the sexes. If there are only two treatments, the sex by treatment  
181 interaction can be quantified by a single number: the difference in effect size between  
182 males and females, and the P-value for the interaction tests whether this number is  
183 significantly different from zero. In our example, the effect size in males is  $(11.41 -$   
184  $10.13) = 1.28$ , whereas it is  $(10.16 - 10.25) = -0.09$  in females. Thus, the difference in  
185 effect sizes is  $-0.09 - 1.28 = -1.37$  (note that this value could also be calculated as 1.37,  
186 depending on which means are subtracted from which). While such estimates of the  
187 magnitude of an interaction are often not reported, it is useful to consider this estimate

188 and whether a difference in effect size is biologically important, rather than simply  
189 whether or not a difference is statistically significant from zero.

190 We performed Bayesian analyses using a BAYES statement in proc GENMOD in  
191 SAS ver. 9.3 (see Appendix 2 for code), and using the BayesFactor package version  
192 0.9.11-1 in R (see Appendix 3 for code). To interpret the results of a Bayesian analysis, it  
193 is necessary to understand the posterior distribution. The analysis uses random  
194 simulation to approximate the posterior distribution by generating many (e.g., 10000)  
195 combinations of parameter values that are consistent with the observed data and the prior  
196 distribution (although an uninformative prior can be selected which will have little  
197 influence on the results)<sup>7</sup>. In this case, the parameters include the effects of sex, treatment  
198 and the sex by treatment interaction, and each combination of parameters generated by  
199 the analysis is a posterior sample.

200 One of the results of the SAS analysis is an estimate of the sex by treatment  
201 interaction, i.e., 1.37, as calculated above; this parameter is not calculated automatically  
202 by BayesFactor, but can be obtained from its output (see Appendix 3). Another result of  
203 the SAS analysis is the 75th percentile of the sex by treatment interaction among  
204 posterior samples, which is -0.86 (Table A2). In other words, the sex by treatment  
205 interaction is less than -0.86 in 75% of posterior samples. Thus, one could report that  
206 there is a probability of 0.75 that the difference in effect size between males and females  
207 is 0.86 or greater. If 0.86 was considered a biologically important difference, one could  
208 conclude that there was a substantial probability (0.75) that the difference in effect  
209 between males and females was biologically important. Note that the Bayesian approach  
210 allows a more intuitive statement than that one would obtain from a traditional

211 confidence interval (e.g., “were we to repeat the experiment many times, the 95%  
212 confidence interval would include the true value in 95% of repetitions”).

213           What if 0.86 was not considered a biologically important difference? One could  
214 alternatively identify a different value for the sex by treatment interaction, and determine  
215 its probability. This can be achieved quite easily using the distribution of the sex by  
216 treatment interaction among posterior samples (which can be obtained using the  
217 OUTPOST option in the BAYES statement in SAS, or the POSTERIOR function in  
218 BayesFactor). For instance, if a difference in effect size of 1.2 would be considered  
219 biologically important, one can determine where 1.2 occurs among the set of posterior  
220 samples. In this example, among 10000 posterior samples, the 5836th lowest value is -  
221 1.20054 and the 5837th lowest value is 1.19970 (values will vary slightly from run to run  
222 because the algorithm used by the analysis involves random simulation). Thus, one could  
223 conclude that “there is a probability of 0.58 (5836/10000) that the effect size in males is  
224 at least 1.2 units greater than that in females” (Table A2). Such a statement might not  
225 provide very convincing support for the existence of a sex-specific effect. However, this  
226 statement has incorporated consideration of what is or is not a biologically important  
227 difference, in contrast to the observation that the effect is significant in one sex, but not in  
228 the other. The selection of a “biologically important” effect size could be achieved  
229 objectively, e.g., by using an effect size observed in a seminal paper in the field, or an  
230 effect that would be considered clinically important in humans.

231           Discussing the difference in effect size in absolute terms will not be very intuitive  
232 in many cases. There might be a variety of other thresholds to determine whether the  
233 difference in effect size between males and females is biologically important. For

234 instance, it might be meaningful to express the difference in effect size as a percentage,  
235 e.g., what is the probability that the difference in effect size between males and females is  
236 at least 10% of the value in control males? This is possible using the posterior samples  
237 from a Bayesian analysis by (1) calculating the difference in effect size between males  
238 and females for each posterior sample, (2) calculating the mean value for control males in  
239 each posterior sample, (3) assessing whether the difference in effect size is greater than  
240 10% of the mean of control males in each posterior sample, and (4) counting the  
241 proportion of posterior samples for which this condition is true, i.e., the posterior  
242 probability that the difference in effect size between males and females is at least 10% of  
243 the mean value in control males, which in this example is 0.62 (Table A2). Again, this is  
244 not convincing evidence of a sex-specific effect, but unlike traditional NHST, this  
245 approach has assessed the evidence of a biologically important sex-specific effect. Thus,  
246 this approach provides a more meaningful interpretation of the data than the traditional  
247 approach described at the beginning of this example, i.e., one way ANOVAs performed  
248 separately for each sex, without a sex by treatment interaction term.  
249  
250

251 Appendix 2. SAS code.

```
252 * treat refers to treatment, which can be 1 or 2;
253 * sex is coded as M or F;
254
255 proc genmod;
256     class treat sex;
257     model trait = treat|sex/ type3;
258     lsmeans treat*sex;
259
260 * The code below is the default analysis, which uses a noninformative
261 uniform prior;
262 * The "outpost" option creates a dataset (named "post") containing the
263 posterior samples;
264     bayes outpost = post;
265
266 * this is an alternative, which specifies the use of a Jeffreys'
267 noninformative prior;
268     bayes COEFFPRIOR=JEFFREYS outpost = post;
269
270 * In the dataset containing the posterior samples, SAS creates a
271 variable "treatlsexf", which is the estimate of the difference in effect
272 size between males and females;
273 * The dataset is sorted by treatlsexf to determine, e.g., where a
274 difference of 1.2 ranks among posterior samples;
275 proc sort data = post;
276     by treatlsexf;
277 proc print;
278
279 * For each posterior sample, the estimate for each of the four groups
280 (i.e., treatment 1 males, treatment 1 females, etc.) can be calculated;
281 * The way SAS estimates the parameters is that one group mean is set as
282 the intercept, and then effects of treatment, sex and treatment*sex are
283 estimated as deviations from that reference group;
284 * In this example, the parameter "treat1" describes the effect of
285 treatment 1 compared to treatment 2, the parameter "sexf" describes the
286 effect of being female compared to being male, and the parameter
287 "treatlsexf" describes the deviation of the mean of the treatment 1
288 females from what would be expected given the effects of treatment and
289 sex;
290 data post;
291     set post;
292     trtlsexf = intercept + treat1 + sexf + treatlsexf;
293     trtlsexm = intercept + treat1;
294     trt2sexf = intercept + sexf;
295     trt2sexm = intercept;
296
297 * Once the means of each of the four groups are calculated, it is
298 possible to calculate customized parameters, e.g., whether the
299 difference in effect size is greater than 10% of the mean of treatment 1
300 males in each posterior sample;
301     ratio = 0.1*trtlsexm;
302
303 * abs(treatlsexf) returns the absolute value of the difference in effect
304 size between the sexes for each posterior sample;
305     if abs(treatlsexf) > ratio then check = 1; else check = 0;
```

```
306
307 * The variable "check" allows the calculation of the proportion of
308 posterior samples in which the difference in effect size between males
309 and females is at least 10% of the value in treatment 1 males;
310
311 proc freq;
312     tables check;
313
```

```

314 Appendix 3. R code.
315 # load the "BayesFactor" package
316 library(BayesFactor)
317
318 # The following command computes the Bayes Factor for a linear model
319 (hence "lmBF")
320 # The model is the full model, i.e., treatment, sex, and interaction
321 term
322 # "full" is the name of the object containing the results - another name
323 could be used instead
324 # "rscaleFixed" is the prior scale, which is set to a large number to
325 avoid shrinkage
326 # Shrinkage isn't necessarily undesirable, however, avoiding shrinkage
327 is necessary to make the results comparable with those from the SAS
328 analysis described in Appendix 2
329 # For a discussion of shrinkage, see Kruschke \(2013\)
330 full <- lmBF(trait ~ treat + sex + treat:sex, data=sextreatint,
331 rscaleFixed = 100)
332
333 # The "posterior" function samples from the posterior distribution
334 # BayesFactor calculates the parameters differently than SAS (Appendix
335 2)
336 # In this example, parameters include
337 # the grand mean
338 # a female effect and a male effect of male (equal in magnitude to the
339 female effect, but in the opposite direction)
340 # an effect of treatment 1 and an effect of treatment 2 (equal in
341 magnitude to the effect of treatment 1, but in the opposite direction)
342 # four estimates treatment by sex interaction effects (i.e., one for
343 each group)
344 chainsfull <- posterior(full, iterations = 10000)
345
346 # Those adept at the R language could perform calculations using the
347 posterior distribution directly in R
348 # Alternatively, the posterior samples can be written to an Excel file
349 for analysis using different software
350 # Load the "xlsx" package - this is needed to write the posterior
351 samples to an Excel file
352 library(xlsx)
353
354 # Write the posterior samples to an Excel file (be patient)
355 write.xlsx(chainsfull, "BayesFactor_posteriors.xlsx")
356
357
358

```

359 Table A1. Simulated data set.

Treatment	Sex	Trait	Group average
1	M	9.8034	
1	M	13.5123	
1	M	10.4177	
1	M	12.4585	
1	M	9.9865	
1	M	11.6487	
1	M	12.2631	
1	M	10.5206	
1	M	10.9397	
1	M	12.5582	11.41
1	F	10.8819	
1	F	11.0768	
1	F	10.1114	
1	F	9.4873	
1	F	9.9139	
1	F	11.2894	
1	F	9.5791	
1	F	9.4602	
1	F	9.7483	
1	F	10.0784	10.16
2	M	8.6520	
2	M	8.3239	
2	M	12.3612	
2	M	10.9094	
2	M	10.9274	
2	M	10.3007	
2	M	9.5637	
2	M	9.7066	
2	M	9.2029	
2	M	11.3043	10.13
2	F	10.0371	
2	F	8.6253	
2	F	9.1643	
2	F	11.7190	
2	F	11.8492	
2	F	9.7477	
2	F	12.2489	
2	F	8.4438	
2	F	10.3420	
2	F	10.3152	10.25

360

361

362 Table A2. Comparison of results between SAS and R (see Appendices 2 and 3 for details)

Software	75th percentile of the difference in effect size between males and females	Probability that the effect size in males is at least 1.2 units greater than that in females	Probability that the difference in effect size between males and females is at least 10% of the value in control males
SAS, uniform prior	-0.86	0.58	0.62
SAS, Jeffreys' prior	-0.84	0.58	0.61
R, BayesFactor	-0.89	0.60	0.63

363