# Additional documentation for GSG

Michael B. Morrissey[1] and Krzysztof Sakrejda[2]

December 22, 2013

[1]michael.morrissey@st-andrews.ac.uk
[2]sakrejda@cns.umass.edu

# Contents

## 1 Selection gradients and fitness functions for human birth weight and gestation length via variation in neonatal survival

The tensor product smooth-based generalized additive model in Morrissey and Sakrejda (2013) was fitted by:

```
library(mgcv)
data(humanNeonatal)
neonatalGam <- gam(nns~te(bw,gest), family='binomial', data=humanNeonatal)
```

We then used the function `gam.gradients()` to obtain selection gradients

```
> library(gsg)
> gradientsGam <- gam.gradients(neonatalGam, phenotype=c("bw","gest"),
+              n.boot=1000, standardize=TRUE)
Calculating bootstrap standard errors...

      ... estimated completion at  2012-06-10 16:19:03 ...done.
>
> round(gradientsGam,4)
          estimates     SE P.value
B-bw         0.0223 0.0034   0.000
B-gest       0.0037 0.0031   0.242
G-bw        -0.0350 0.0048   0.000
G-gest      -0.0087 0.0025   0.000
G-bw-gest   -0.0042 0.0037   0.300
```

The computation with 1000 bootstrap replicates took approximately 1.9 hours using a personal computer with an Intel Core 2 processor at 1.8 GHz. The same computation required approximately 7.5 minutes on an Intel i7 at 4.2 GHz using 4 cores.

## 2 Plotting a fitness landscape

The bivariate fitness landscape in Morrissey and Sakrejda (2013) was obtained by:

```
neonatal.fl<-fitness.landscape(mod= neonatalGam,
        phenotype=c("bw","gest"),plt.density=10,PI.method='n')
```

and the plot was made similarly to:

```
p<-matrix(neonatal.fl$Wbar,10,10,byrow=TRUE)
par(mar=c(5.5,6,1,1),oma=rep(1,4),las=1,cex.lab=1.2)
```

```
54  contour(t(p),xaxt='n',yaxt='n',xlab="Mean birth mass (kg)",ylab="")
55  axis(at=seq(0,1,length.out=10),
56          round(unique(neonatal.fl$points[,1]),2),side=1)
57  axis(at=seq(0,1,length.out=10),
58          round(unique(neonatal.fl$points[,2]),2),side=2)
59  par(las=0)
60  mtext(side=2,outer=TRUE,line=-1.5,
61          "Mean gestation length (days)",cex=1.2)
```

## 62  3   The Lande-Arnold selection analysis as a special case

63  A quadratic approximation of the bivariate human neonatal fitness function can be ob-

64  tained by:

```
65  neonatalQuadratic <- gam(nns~bw+gest+I(bw^2)+
66              I(gest^2)+I(bw*gest), family='gaussian',
67              data=humanNeonatal)
```

68      Obtaining the first and second order partial derivatives of this function is an implemen-

69  tation of the Lande and Arnold (1983) selection analysis as a special case of the general

70  formulation described in Morrissey and Sakrejda (2013):

```
71  > gradientsQuadratic <- gam.gradients(neonatalQuadratic,
72  +           phenotype=c("bw","gest"),
73  +           n.boot=1000, standardize=TRUE)
74  Calculating bootstrap standard errors...
75
76      ... estimated completion at  2012-06-10 17:00:13 ...done.
77  >
78  > round(gradientsQuadratic,4)
79          estimates     SE P.value
80  B-bw          0.0292 0.0040   0.000
81  B-gest        0.0045 0.0035   0.198
82  G-bw         -0.0599 0.0059   0.000
83  G-gest       -0.0171 0.0049   0.000
84  G-bw-gest    -0.0102 0.0042   0.012
```

85      Note that standardizations necessary for the Lande and Arnold (1983) analysis (mean

86  standardization of traits and analysis of fitness on the relative scale, scaling of 0.5 for the

87  diagonal quadratic coefficients; Stinchcombe et al. 2008) are intrinsic to the calculations

88  implemented in `gam.gradients`:

```
89  humanNeonatal$st.bw <- (humanNeonatal$bw-mean(humanNeonatal$bw))/
90                         sd(humanNeonatal$bw)
91  humanNeonatal$st.gest <- (humanNeonatal$gest-mean(humanNeonatal$gest))/
92                         sd(humanNeonatal$gest)
93  humanNeonatal$w<-humanNeonatal$nns/mean(humanNeonatal$nns)
94  neonatalQuadraticStandardized <- gam(w~ st.bw + st.gest +I(0.5* st.bw^2)
95                         +I(0.5*st.gest^2)+I(st.bw*st.gest), family='gaussian',
96                         data=humanNeonatal)
97  gradientsQuadraticS <- gam.gradients(neonatalQuadraticStandardized,
98                         phenotype=c("st.bw","st.gest"),
99                         n.boot=1000, standardize=TRUE)
```

100  This produces the same selection gradients estimates. Differences in the standard errors

101  are due to MC error.

```
102  > round(gradientsQuadraticS,4)
103                    estimates     SE P.value
104  B-st.bw              0.0292 0.0038   0.000
105  B-st.gest            0.0045 0.0035   0.190
106  G-st.bw             -0.0599 0.0063   0.000
107  G-st.gest           -0.0171 0.0048   0.000
108  G-st.bw-st.gest     -0.0102 0.0042   0.018
```

## 109  4  Compromises between model flexibility and simplicity

110  As acknowledged in the discussion of Morrissey and Sakrejda (2013), it will not always be

111  sensible to fit fully flexible smooth terms for characterizing multivariate fitness functions.

112  The large neonatal survival databased allowed the bivariate tensor product smooth to be

113  fitted, but such data are often not available in evolutionary quantitative genetic studies of

114  wild populations. Slightly less flexible models may often be sensible, and can be handled

115  in the analytical framework supported by the R package GSG. A generally useful approach

116  may be to model fitness as semi-parametric smooth functions of each variable, while han-

117  dling interactions parametrically. This fitness function could be applied to the analysis of

118  the human neonatal data via:

119  `neonatalLessFlexible<-gam(nns~s(bw)+s(gest)+bw:gest,`

```
120                    family='binomial',data=humanNeonatal)
```

121    Analysis based on this somewhat less flexible characterization of the fitness function

122  proceeds similarly, and provides very similar results:

```
123 > gradientsLessFlexible<-gam.gradients(neonatalLessFlexible,
124 +                      phenotype=c("bw","gest"),
125 +                      n.boot=1000, standardize=TRUE)
126 Calculating bootstrap standard errors...
127
128      ... estimated completion at  2012-06-11 09:20:08 ...done.
129 > round(gradientsLessFlexible,4)
130          estimates     SE P.value
131 B-bw        0.0217 0.0038   0.000
132 B-gest      0.0033 0.0033   0.346
133 G-bw       -0.0339 0.0063   0.000
134 G-gest     -0.0184 0.0045   0.000
135 G-bw-gest  -0.0019 0.0034   0.542
```

136    This more constrained model may in fact have some interpretive benefits, for example,

137  the lack of statistical support for the interaction between birth weight and gestation length

138  in the fitness function compliments the estimate of the small (and also statistically unsup-

139  ported) off-diagonal element of the matrix of quadratic selection coefficients (see above and

140  Morrissey and Sakrejda 2013):

```
141 > summary(neonatalLessFlexible)
142
143 Family: binomial
144 Link function: logit
145
146 Formula:
147 nns ~ s(bw) + s(gest) + bw:gest
148
149 Parametric coefficients:
150              Estimate Std. Error z value Pr(>|z|)
151 (Intercept)  3.7033796  4.5862541    0.807     0.419
152 bw:gest     -0.0005008  0.0051294   -0.098     0.922
153
154 Approximate significance of smooth terms:
155          edf Ref.df Chi.sq  p-value
156 s(bw)   3.861  4.843 113.24  < 2e-16 ***
157 s(gest) 5.073  6.090  30.74 3.09e-05 ***
```

```
158  ---
159  Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1   1
160
161  R-sq.(adj) =  0.235   Deviance explained = 22.7%
162  UBRE score = -0.67517  Scale est. = 1         n = 7036
```

## 163  5   Notes about algorithms for calculating standard errors and/or

## 164     p-values

165  The parametric bootstrap, as applied in Morrissey and Sakrejda (2013) is the default
166  method for obtaining coefficients of selection gradients and prediction intervals fitness
167  landscapes, in each function in GSG. Alternative algorithms include case bootstrapping,
168  simulation from an approximation to the posterior distribution of the gam model param-
169  eters, and a permutation test (P-values only). The two bootstrap algorithms, and the
170  posterior simulations, allow the smoothing parameters to be fixed across replicates, or
171  re-estimated. By default, they are fixed following Schluter (1988).

## 172  6   A brief example with a Poisson fitness response

173  Fitness data are often counts, and so reasonably modelled as Poisson variables. Implement-
174  ing the methods described in Morrissey and Sakrejda (2013) using GSG is straightforward
175  for Poisson or other fitness distributions is straightforward. The functions in GSG that
176  extract data from a fitted `gam` object rely on prediction on the data scale, and so analysis
177  based on different assumed distributions of fitness simply require fitting a model with a
178  different error structure.

179     The example code below simulates a Poisson fitness response as a function of a sin-
180  gle trait, and shows the implementation of an analysis to obtain the associated selection
181  gradient:

```
182  > n<-200
183  > z<-rnorm(n,0,1)
```

```
184 > W<-rpois(n,exp(1+z-0.5*z^2))
185 > simPoisData<-as.data.frame(list(z=z,W=W))
186 >
187 > simPoisGam<-gam(W~s(z),family='poisson',data=simPoisData)
188 >
189 > gradientsPoisSim<-gam.gradients(simPoisGam,phenotype="z")
190 Calculating bootstrap standard errors...[1] 100
191
192       ... estimated completion at  2012-06-11 09:30:52 ...done.
193 >
194 > round(gradientsPoisSim,4)
195     estimates     SE P.value
196 B-z    0.4423 0.0642   0.000
197 G-z   -0.2068 0.0852   0.034
```

## 198  7  Direct calculation of selection differentials

199 Selection differentials are defined most simply as the change in the central moments of the

200 phenotypic distribution due to selection (Endler, 1986; Lande and Arnold, 1983). Gen-

201 erally, these can be calculated as the difference between the means, variances, and co-

202 variances, weighted by fitness, and the unweighted moments. These are calculated using

203 `moments.differentials()` in the R package GSG

```
204 > humanDifferentials<-moments.differentials(
205 +         z=humanNeonatal[,c("bw","gest")],
206 +         W=humanNeonatal$nns,n.boot=1000,standardized=TRUE)
207 >
208 > round(humanDifferentials,4)
209       Coefficient      SE P-value
210 S 1         0.0667 0.0055       0
211 S 2         0.0612 0.0056       0
212 C 1        -0.2057 0.0153       0
213 C 2        -0.2160 0.0183       0
214 C 1,2      -0.1919 0.0157       0
```

## 215  8  Lasso and ridge regression selection analysis

216 Selection gradients were obtained from the regularised regression analyses in Morrissey

217 (2013) by tricking `gam.gradients()` into doing the analysis. First the regression analyses

218   were fitted; using the lasso as an example:

```
219  library(glmnet)
220  data(SoayLambs)
221
222  phen<-c("WEIGHT","HINDLEG","HORNLEN","lnKeds")
223  covars<-SoayLambs[, phen]
224  for(i in 1:4){
225    for(j in 1:i){
226      covars<-cbind(covars,covars[,phen[i]]*covars[,phen[j]])
227      names(covars)[length(names(covars))]<-paste(phen[i],phen[j],sep="")
228    }
229  }
230
231
232  lamb.lasso<-cv.glmnet(x=as.matrix(covars), y=
233      SoayLambs$W, family='binomial',alpha=1)
```

234   The coefficients of the fitted lasso model are thus:

```
235  > predict(lamb.lasso,type="coefficients",s="lambda.min")
236  15 x 1 sparse Matrix of class "dgCMatrix"
237                           1
238  (Intercept)     1.6051944
239  WEIGHT          1.0970974
240  HINDLEG         0.2661741
241  HORNLEN        -0.4738859
242  lnKeds         -0.2270427
243  WEIGHTWEIGHT    0.1396266
244  HINDLEGWEIGHT    .
245  HINDLEGHINDLEG   .
246  HORNLENWEIGHT   0.0687889
247  HORNLENHINDLEG   .
248  HORNLENHORNLEN   .
249  lnKedsWEIGHT     .
250  lnKedsHINDLEG  -0.2111347
251  lnKedsHORNLEN    .
252  lnKedslnKeds     .
253  >
```

254   These can be forced into a gam object, and then the gradients are obtained using

255   gam.gradients():

```
256  dummy.gam<-gam(W~WEIGHT+HINDLEG+HORNLEN+lnKeds
257     +I(WEIGHT^2)
```

```
258    +I(WEIGHT*HINDLEG)   +I(HINDLEG^2)
259    +I(WEIGHT*HORNLEN)   +I(HINDLEG*HORNLEN)   +I(HORNLEN^2)
260    +I(WEIGHT*lnKeds)    +I(HINDLEG*lnKeds)    +I(HORNLEN*lnKeds) + I(lnKeds^2),
261          family='binomial',data= SoayLambs)
262
263 predict(lamb.lasso,type="coefficients",s="lambda.min")
264
265 lasso.coefs<-as.numeric(predict(lamb.lasso,type="coefficients",s="lambda.min"))
266 dummy.gam$coefficients<-lasso.coefs
267
268 lasso.grads<-gam.gradients(mod=dummy.gam,phenotype=phen,se.method='n')
```

269   The lasso-based selection gradients are thus:

```
270 > lasso.grads
271                       estimates SE P.value
272 B-WEIGHT            0.161052875 NA      NA
273 B-HINDLEG           0.039538491 NA      NA
274 B-HORNLEN          -0.086004182 NA      NA
275 B-lnKeds           -0.022464858 NA      NA
276 G-WEIGHT           -0.046488274 NA      NA
277 G-HINDLEG          -0.007482505 NA      NA
278 G-HORNLEN          -0.017072158 NA      NA
279 G-lnKeds           -0.005697433 NA      NA
280 G-WEIGHT-HINDLEG   -0.020150751 NA      NA
281 G-WEIGHT-HORNLEN    0.049823337 NA      NA
282 G-HINDLEG-HORNLEN   0.008351592 NA      NA
283 G-WEIGHT-lnKeds     0.020535502 NA      NA
284 G-HINDLEG-lnKeds   -0.031363207 NA      NA
285 G-HORNLEN-lnKeds   -0.007713071 NA      NA
286
```

287   Obtaining standard errors and P-values for such an analysis does not seem meaningful,

288 as the shrinkage and variable selection inherent in the lasso (or elastic net regression,

289 generally) to some extent generates parameters that reflect both the pattern in the data

290 and the extent to which it is statistically supported.


291 **9   Generalised projection-pursuit regression and selection gradi-**

292     **ents**


293 Characterisation of a fitness landscape might proceed as above:

```
294  data(SoayLambs)
295  phen<-c("WEIGHT","HINDLEG","HORNLEN","lnKeds")
296  fit.land<-gppr(y="W",xterms=phen,
297                 data=SoayLambs,family='binomial')
298
299  grads<-gppr.gradients(mod=fit.land,
300                 phenotype=phen,
301                 family='binomial')
```

302    In which case the gradients produced are

```
303  > grads$ests
304                          estimates          SE P.value
305  B-WEIGHT              1.942585e-01 0.063735725   0.002
306  B-HINDLEG            3.779547e-02 0.059054385   0.494
307  B-HORNLEN           -1.212769e-01 0.044498700   0.006
308  B-lnKeds            -3.235531e-02 0.031340850   0.284
309  G-WEIGHT            -8.071137e-02 0.054259581   0.012
310  G-HINDLEG            4.556036e-05 0.015213039   0.660
311  G-HORNLEN           -2.682024e-02 0.025293852   0.018
312  G-lnKeds             5.079316e-04 0.004820774   0.930
313  G-WEIGHT-HINDLEG    -1.956484e-02 0.024748309   0.496
314  G-WEIGHT-HORNLEN     5.039030e-02 0.028296653   0.008
315  G-HINDLEG-HORNLEN    8.371374e-03 0.018093772   0.502
316  G-WEIGHT-lnKeds      1.344184e-02 0.014885079   0.292
317  G-HINDLEG-lnKeds    -4.058068e-05 0.006674169   0.866
318  G-HORNLEN-lnKeds    -1.045944e-02 0.010646091   0.276
319  >
```

320    One might wish to obtain the selection gradients associated with the axes of phenotype
321  of the gppr analysis. This could be done by re-fitting a gam with the same type of regression
322  function to rotated data:

```
323  SoayLambs$SelTerm<-as.matrix(SoayLambs[,phen]) %*% as.matrix(fit.land$alpha)
324
325  new.mod<-gam(W~s(SelTerm,bs="cr"),data=SoayLambs,family='binomial')
326
327  grads2<-gam.gradients(mod=new.mod,phenotype="SelTerm",standardized=TRUE)
```

328  This yields the gradient estimates of selection along the axis defined by the gppr as

```
329  > grads2$ests
330              estimates          SE P.value
331  B-SelTerm  0.18616108 0.03970303   0.000
332  G-SelTerm -0.05821734 0.10574565   0.218
333  >
```

334   The SEs and P-values should be taken with a grain of salt. Since the gppr analysis has

335   specifically sought to find an axis that explains fitness variation, statistical inference of

336   selection focusing only on that direction, and not accounting for all the other directions

337   that were not chosen, is inappropriate. The P-values should thus be thought of as requiring

338   correction for multiple testing, although just how many tests (i.e., of phenotypic directions)

339   one should think of the gppr analysis as having conducted, I don't know.

340   It seems that it should be instructive, at least, to consider what variance in expected

341   fitness would have been apparently explained under an hypothesis of no selection. Although

342   I used this approach in Morrissey (2013), I do not specifically want to promote it at present

343   as a "canned solution", so it is not specifically implemented in any function in gsg. This

344   approach is pretty easily implemented, though:

```
345   n.perm<-1000
346   varWperm<-array(dim=1000)
347   for(x in 1:n.perm){
348      SoayLambs$permutedW<-SoayLambs$W[sample(1:length(SoayLambs$W),
349                       length(SoayLambs$W),replace=FALSE)]
350      perm.mod<-gppr(y="permutedW",xterms=phen,
351               data=SoayLambs,family='binomial')
352      varWperm[x]<-var(inv.logit(predict(perm.mod,type="raw")))
353   }
```

354   The variance in expected absolute fitness (survival probability) from the fitted model is

```
355   > var(inv.logit(predict(fit.land,type="raw")))
356   [1] 0.02917478
357   >
```

358   which is very much in the tail of our null distribution

```
359   > table(var(inv.logit(predict(fit.land,type="raw")))>varWperm)/n.perm
360
361   TRUE
362      1
363   >
364   > quantile(varWperm,probs=c(0.025,0.25,0.5,0.75,0.975))
365          2.5%           25%           50%           75%         97.5%
366   0.0004020845 0.0016562070 0.0027748300 0.0045031314 0.0097702922
367   >
```

# 368 **References**

369 Endler, J. A., 1986. Natural selection in the wild. Princeton University Press.

370 Lande, R. and S. J. Arnold, 1983. The measurement of selection on correlated characters.
371     Evolution 37:1210–1226.

372 Morrissey, M. B., 2013. In search of the best methods for multivariate selection analysis.
373     in preparation for submission to Methods in Ecology and Evolution .

374 Morrissey, M. B. and K. Sakrejda, 2013. Unification of regression-based approaches to the
375     analysis of natural selection. Evolution 67:2094–2100.

376 Schluter, D., 1988. Estimating the form of natural selection on a quantitative trait. Evo-
377     lution 42:849–861.

378 Stinchcombe, J. R., A. F. Agrawal, P. A. Hohenlohe, S. J. Arnold, and M. W. Blows, 2008.
379     Estimating nonlinear selection gradients using quadratic regression coefficients: double
380     or nothing? Evolution 62:2435–2440.