

# Surrogate Phenotype Regression Analysis

Zachary R. McCaw

2020-11-28

## Contents

- Setting
- Example Data
- Estimation
- Inference

## Setting

For each of  $n$  independent subjects, suppose two continuous outcomes are potentially observed. Let  $T_i$  denote the *target* outcome, and let  $S_i$  denote the *surrogate* outcome. Group the target and surrogate outcomes into a bivariate outcome vector  $Y_i = (T_i, S_i)'$ . For each subject, either the target or the surrogate is potentially missing. Suppose the target mean depends on a vector of covariates  $x_i$ , and the surrogate mean depends on a vector of covariates  $z_i$ :

$$\begin{aligned}\mu_{T,i} &= \mathbb{E}(T_i|x_i) = x_i'\beta \\ \mu_{S,i} &= \mathbb{E}(S_i|z_i) = z_i'\alpha\end{aligned}$$

Let  $\mu_i = (\mu_{T,i}, \mu_{S,i})'$  denote the mean vector. Consider the bivariate normal regression model:

$$\begin{pmatrix} T_i \\ S_i \end{pmatrix} \Big| (x_i, z_i) \sim N \left\{ \begin{pmatrix} x_i'\beta \\ z_i'\alpha \end{pmatrix}, \begin{pmatrix} \Sigma_{TT} & \Sigma_{TS} \\ \Sigma_{ST} & \Sigma_{SS} \end{pmatrix} \right\}$$

This package provides methods for estimation of the model parameters  $(\beta, \alpha, \Sigma)$ , and for inference on components of the target regression parameters  $\beta$ . In the case of bilateral (target, surrogate) missingness, estimation is performed via an expectation maximization (EM) procedure. In the case of unilateral target missingness, estimation is performed via an accelerated, generalized least squares (GLS) procedure.

## Example Data

Below, data are simulated for  $n = 10^3$  subjects. The target  $\mathbf{X}$  and surrogate  $\mathbf{Z}$  design matrices each contain an intercept and three standard normal covariates. The regression coefficient for the target outcome is  $\beta = (-1, 0.1, -0.1, 0)$ . The regression coefficient for the surrogate outcome is  $\alpha = (1, -0.1, 0.1, 0)$ . The target and surrogate outcome each have unit variance  $\Sigma_{TT} = \Sigma_{SS} = 1$ . The target-surrogate covariance, equivalently the correlation, is  $\Sigma_{TS} = \Sigma_{ST} = 0.5$ . An outcome matrix for which 10% of the target outcomes and 20% of the surrogate outcomes are missing completely at random is simulated using `rBNR`.

```

library(Spray)
set.seed(100)

# Observations.
n <- 1e3

# Target design.
X <- cbind(1, matrix(rnorm(3 * n), nrow = n))

# Surrogate design.
Z <- cbind(1, matrix(rnorm(3 * n), nrow = n))

# Target parameter.
b <- c(-1, 0.1, -0.1, 0)

# Surrogate parameter.
a <- c(1, -0.1, 0.1, 0)

# Covariance matrix.
sigma <- matrix(c(1, 0.5, 0.5, 1), nrow = 2)

# Generate data.
Y <- rBNR(X, Z, b, a, t_miss = 0.1, s_miss = 0.2, sigma = sigma);
t <- Y[, 1]
s <- Y[, 2]

```

## Formatting Assumptions

The target and surrogate outcome vectors ( $\mathbf{t}$ ,  $\mathbf{s}$ ) both have length  $n$ . The unobserved values of the target or surrogate outcome are set to NA. The target  $\mathbf{X}$  and surrogate  $\mathbf{Z}$  model matrices are numeric, with all factors and interactions expanded. The model matrices contain no missing values.

## Estimation

Estimation of the bivariate normal regression model is performed using `Fit.BNR`. If the surrogate outcome vector  $\mathbf{s}$  contains missing values, or if the surrogate design matrix  $\mathbf{Z}$  differs from the target design matrix  $\mathbf{X}$ , then the EM algorithm is applied. Otherwise, estimation is performed via GLS, which is significantly faster.

```

# Fit bivariate normal regression model.
fit <- Fit.BNR(
  t = t,
  s = s,
  X = X,
  Z = Z
)
show(fit)

## Objective increment: 1.76
## Objective increment: 0.00795
## Objective increment: 0.000344
## Objective increment: 3.17e-05
## Objective increment: 3.33e-06

```

```
## Objective increment: 3.58e-07
## 5 update(s) performed before tolerance limit.
##
##      Outcome Coefficient   Point      SE      L      U      p
## 1      Target           x1 -1.0500 0.0322 -1.1200 -0.9890 2.45e-234
## 2      Target           x2  0.1190 0.0276  0.0653  0.1740 1.52e-05
## 3      Target           x3 -0.0839 0.0298 -0.1420 -0.0255 4.87e-03
## 4      Target           x4  0.0160 0.0275 -0.0379  0.0698 5.61e-01
## 5 Surrogate           z1  0.9610 0.0335  0.8950  1.0300 6.80e-181
## 6 Surrogate           z2 -0.0765 0.0310 -0.1370 -0.0157 1.37e-02
## 7 Surrogate           z3  0.1240 0.0308  0.0637  0.1840 5.54e-05
## 8 Surrogate           z4 -0.0482 0.0300 -0.1070  0.0107 1.09e-01
##
##      Covariance Point      SE      L      U
## 1      Target 0.958 0.0450 0.891 1.030
## 2 Target-Surrogate 0.502 0.0377 0.464 0.539
## 3      Surrogate 0.948 0.0471 0.879 1.020
```

The output is an object of class `bnr` with these slots:

- `@Covariance` containing the target-surrogate covariance matrix.

```
round(fit@Covariance, digits = 3)
```

```
##      Target Surrogate
## Target    0.958    0.502
## Surrogate 0.502    0.948
```

- `@Covariance.info` containing the information matrix for  $(\Sigma_{TT}, \Sigma_{TS}, \Sigma_{SS})$ .

```
round(fit@Covariance.info, digits = 3)
```

```
##      Target-Target Target-Surrogate Surrogate-Surrogate
## Target-Target      837.789      -771.178      203.979
## Target-Surrogate    -771.178      1881.356     -779.437
## Surrogate-Surrogate  203.979     -779.437      800.209
```

- `@Covariance.tab` containing the estimated covariance parameters in tabular format.

```
fit@Covariance.tab
```

```
##      Covariance   Point      SE      L      U
## 1      Target 0.9582826 0.04500485 0.8908536 1.0308153
## 2 Target-Surrogate 0.5015663 0.03766063 0.4639057 0.5392269
## 3      Surrogate 0.9481284 0.04705099 0.8792191 1.0224384
```

- `@Regression.info` containing the information matrix for  $(\beta, \alpha)$ .

```
round(fit@Regression.info, digits = 3)
```

```
##      x1      x2      x3      x4      z1      z2      z3      z4
## x1 1218.880    6.001   -5.489   -8.886 -534.388  -17.155  -1.010  -17.337
## x2    6.001 1316.182   29.125  -74.870  -5.337   17.600   19.296   -4.338
## x3   -5.489   29.125 1141.283 -115.861   1.255   65.024   12.191   43.722
## x4   -8.886  -74.870 -115.861 1341.487   0.442    6.980  -41.193  -14.785
## z1 -534.388  -5.337   1.255   0.442 1126.463   21.000    7.194   47.448
## z2  -17.155   17.600   65.024    6.980   21.000 1047.657  -52.538   28.147
## z3   -1.010   19.296   12.191  -41.193    7.194  -52.538 1061.920   30.751
## z4  -17.337   -4.338   43.722  -14.785   47.448   28.147   30.751 1115.257
```

- @Regression.tab containing the estimated regression parameters in tabular format.

```
fit@Regression.tab
```

```
##      Outcome Coefficient      Point      SE      L      U
## 1      Target      x1 -1.05209739 0.03218776 -1.11518425 -0.98901053
## 2      Target      x2  0.11947911 0.02761939  0.06534610  0.17361211
## 3      Target      x3 -0.08392919 0.02981181 -0.14235926 -0.02549912
## 4      Target      x4  0.01597697 0.02748197 -0.03788669  0.06984064
## 5 Surrogate      z1  0.96106173 0.03350996  0.89538342  1.02674004
## 6 Surrogate      z2 -0.07646045 0.03101070 -0.13724031 -0.01568059
## 7 Surrogate      z3  0.12401459 0.03076198  0.06372222  0.18430697
## 8 Surrogate      z4 -0.04817376 0.03001576 -0.10700358  0.01065606
##
##      P
## 1 2.449017e-234
## 2 1.519002e-05
## 3 4.873169e-03
## 4 5.609965e-01
## 5 6.799248e-181
## 6 1.367784e-02
## 7 5.543993e-05
## 8 1.085051e-01
```

- @Residuals containing the target and surrogate residuals.

```
round(head(fit@Residuals), digits = 3)
```

```
##      Target Surrogate
## 1  0.869      0.889
## 2 -0.793     -0.803
## 3 -0.850     -1.646
## 4  0.103      0.464
## 5  1.868      NA
## 6 -0.151      0.014
```

## Inference

Wald and Score tests on  $\beta$  are specified using a logical vector `is_zero`, with length equal to the number of columns in the target model matrix `X`, and indicating which regression coefficients are zero under the *null hypothesis*. At least one element of `is_zero` must be `TRUE` (i.e. a test must be specified) and at least one element of `is_zero` must be `FALSE` (i.e. a null model must be estimable).

Below, various hypotheses are tested on the example data. The first is an overall test of  $H_0 : \beta_1 = \beta_2 = \beta_3 = 0$ , which is false. The second assesses  $H_0 : \beta_1 = \beta_2 = 0$ , which is again false, leaving  $\beta_3$  unconstrained. The final considers  $H_0 : \beta_3 = 0$ , which is true, leaving  $\beta_1$  and  $\beta_2$  unconstrained. All models include an intercept  $\beta_0$  under the null.

```
cat("Joint score test of b1 = b2 = b3 = 0", "\n")
test_spec <- c(FALSE, TRUE, TRUE, TRUE)
signif(Test.BNR(t, s, X, Z, is_zero = test_spec, report = FALSE, test = "Wald"), digits = 2)
signif(Test.BNR(t, s, X, Z, is_zero = test_spec, report = FALSE, test = "Score"), digits = 2)

cat("\n", "Joint score test of b1 = b2 = 0, treating b3 as a nuisance", "\n")
test_spec <- c(FALSE, TRUE, TRUE, FALSE)
signif(Test.BNR(t, s, X, Z, is_zero = test_spec, report = FALSE, test = "Wald"), digits = 2)
```

```

signif(Test.BNR(t, s, X, Z, is_zero = test_spec, report = FALSE, test = "Score"), digits = 2)

cat("\n", "Individual score test of b3 = 0, treating b2 and b3 as nuisances", "\n")
test_spec <- c(FALSE, FALSE, FALSE, TRUE)
signif(Test.BNR(t, s, X, Z, is_zero = test_spec, report = FALSE, test = "Wald"), digits = 2)
signif(Test.BNR(t, s, X, Z, is_zero = test_spec, report = FALSE, test = "Score"), digits = 2)

## Joint score test of b1 = b2 = b3 = 0
##      Wald      df      p
## 2.7e+01 3.0e+00 7.2e-06
##      Score      df      p
## 2.6e+01 3.0e+00 1.2e-05
##
## Joint score test of b1 = b2 = 0, treating b3 as a nuisance
##      Wald      df      p
## 2.6e+01 2.0e+00 2.0e-06
##      Score      df      p
## 2.5e+01 2.0e+00 3.3e-06
##
## Individual score test of b3 = 0, treating b2 and b3 as nuisances
##      Wald      df      p
## 0.34 1.00 0.56
##      Score      df      p
## 0.34 1.00 0.56

```