# ddsPLS Exploration

Harpeth Lee

2/4/2022

```
library(ddsPLS2)
library(MASS)
library(spls)
```

```
ddsPLS2::ddsPLS2_App()
```

This code chunk opens an applet that can be used to build models using ddsPLS. Note that it requires the $X$ and $Y$ variables as separate csv files.

Code copied from the simulation_ssdpls2 repository created by Hadrien Lorenzo.

The `get_toy_example()` function simulates a data frame with `n` observations of where 50 of `p` predictors are associated with the single response variable.

```
# Creates a toy data set for the ddsPLS function
toy_ex <- get_toy_example()
```

```
# Creates model from the toy data
toy_mod <- ddsPLS(toy_ex$X, toy_ex$Y)

toy_results <- toy_mod$results
```

## Recreate Toy Example

This is a recreation of the toy example created by Hadrien Lorenzo, the original example can be found here.

```
# Creates toy data set to be used
simu_toy <- get_toy_example(n=50,sqrt_1_minus_sig2 = 0.9025,p = 1000)

# Creates vector of lambda values to be used
lambdas <- seq(0,1,length.out = 30)

# Sets number of bootstrap samples to run
n_B <- 100

# Creates model using ddsPLS algorithm
model_toy <- ddsPLS(simu_toy$X,simu_toy$Y,
                    doBoot = FALSE,
                    lambdas = lambdas,
                    n_B = n_B,
                    verbose = T # whether trace during process
                    )
```

```
model_toy_2 <- ddsPLS(simu_toy$X,simu_toy$Y,
                      doBoot = FALSE,
```

```
              criterion = "Q2",
              lambdas = lambdas,
              n_B = n_B,
              verbose = T # whether trace during process
              )
```

## Design 1

Generates **n** samples of **p** observations with **q** response variables. Projects 5 latent variables onto **p** components.

```
simu_1 <- get_design_1(n=50,sqrt_1_minus_sig2 = 0.99,p = 1000,q = 3)
```

What does the `NCORES` argument do? Setting it to integers greater than 1 gives an error.

Is there a way to include more components in the model?

```
model_1 <- ddsPLS(simu_1$X,simu_1$Y,
                  lambdas = lambdas,
                  n_B=n_B,
                  verbose=T)
```
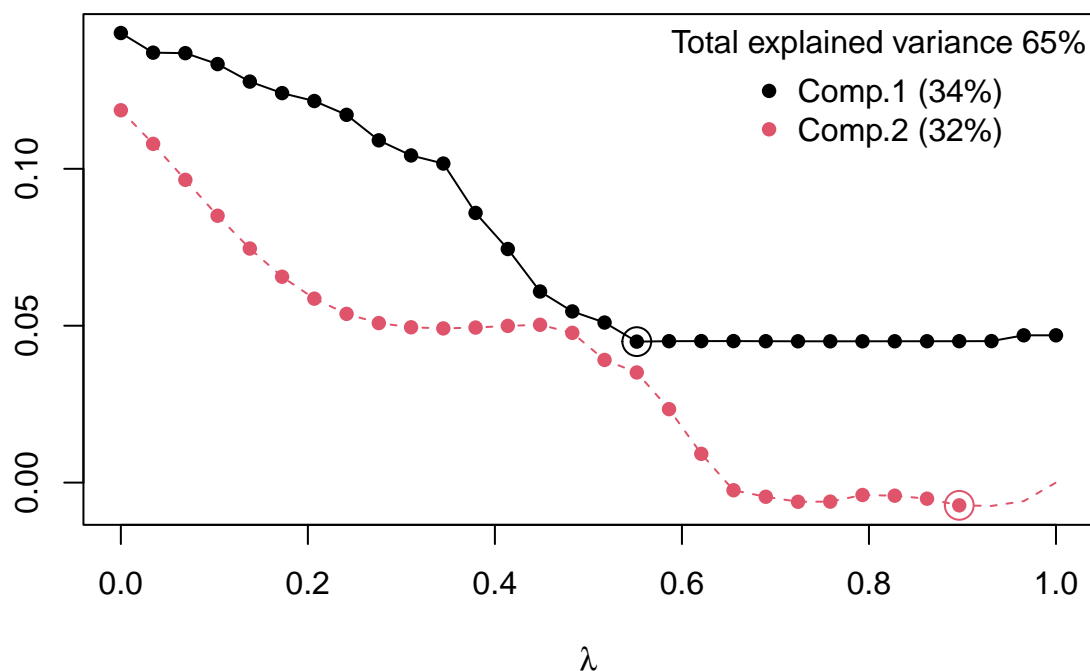
```
##                       --------------
##                      |    ddsPLS     |
## =====================--------------=====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2   R2h  Q2  Q2h VarExpl VarExpl.Tot
##        0.55 0.35 0.35 0.3 0.3     34%         34%
##                                   ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##      lambda  R2   R2h   Q2   Q2h VarExpl VarExpl.Tot
##         0.9 0.4 0.12 0.41 0.27     32%         65%
##                                   ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                  ...component 3 not built!
## =====================              =====================
##                       ===============
```
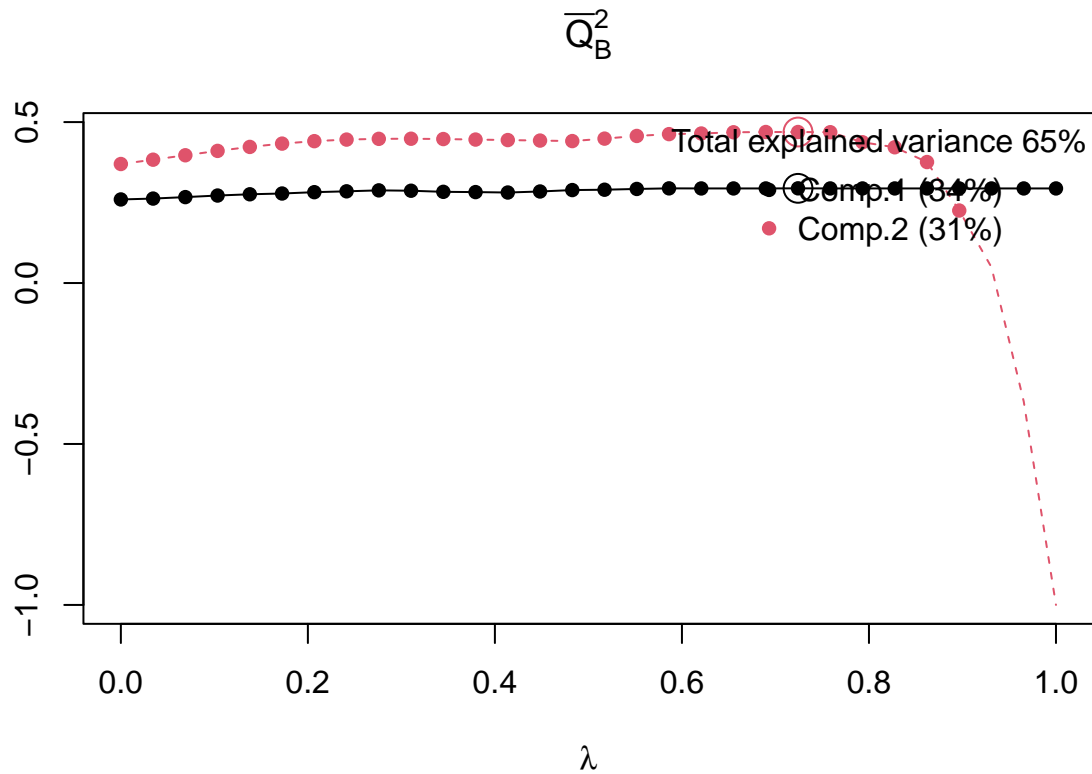
$$\overline{R}_B^2 - \overline{Q}_B^2$$

Total explained variance 65%
- Comp.1 (34%)
- Comp.2 (32%)

```r
model_1_Q2 <- ddsPLS(simu_1$X,simu_1$Y,
                     criterion = "Q2",
                     lambdas = lambdas,
                     n_B=n_B,
                     verbose=T)
```
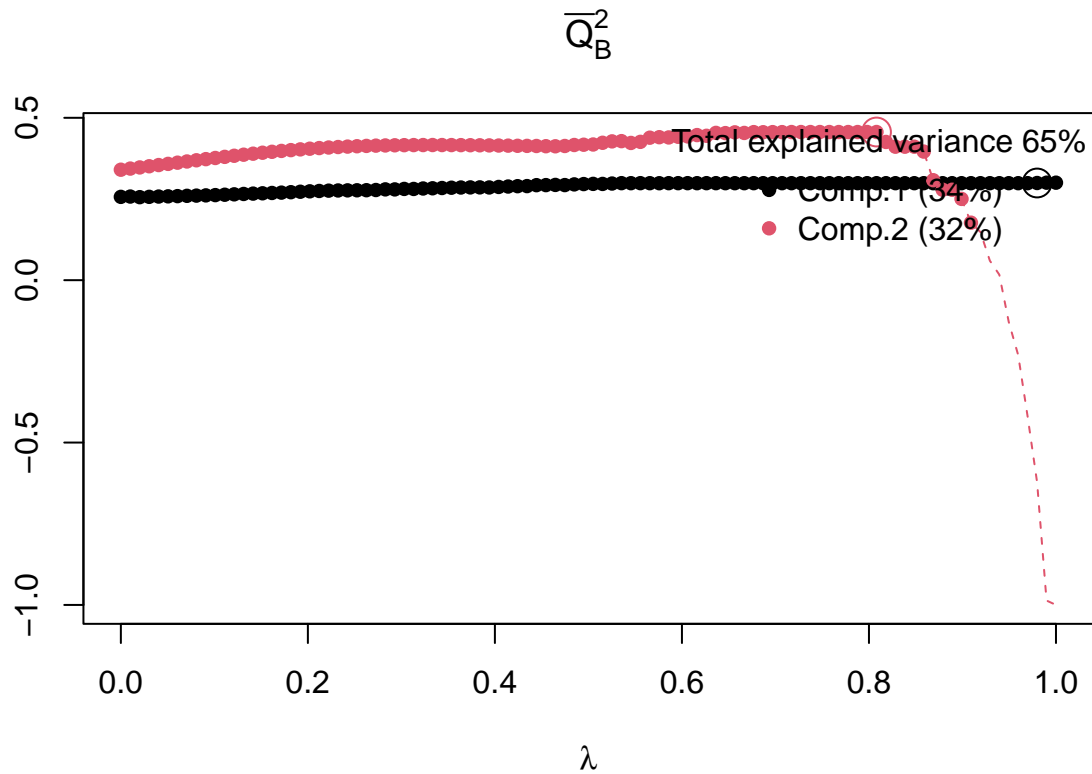
```
##                 ---------------
##                | |   ddsPLS    |
## ================---------------====================
## Should we build component 1 ? Bootstrap pending...
##     lambda   R2   R2h    Q2   Q2h VarExpl VarExpl.Tot
##       0.72 0.34  0.34  0.29  0.29    34%          34%
##                                   ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##     lambda   R2   R2h    Q2   Q2h VarExpl VarExpl.Tot
##       0.72 0.63  0.29  0.63  0.47    31%          65%
##                                   ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                   ...component 3 not built!
## ====================               ====================
##                     ===============
```

$$\overline{Q}_B^2$$

Total explained variance 65%

Comp.1 (34%)

Comp.2 (31%)

$\lambda$

```
model_1_lambda <- ddsPLS(simu_1$X, simu_1$Y,
                         criterion = "Q2",
                         lambdas = seq(0,1,length.out = 100),
                         n_B = n_B,
                         verbose = T)
```
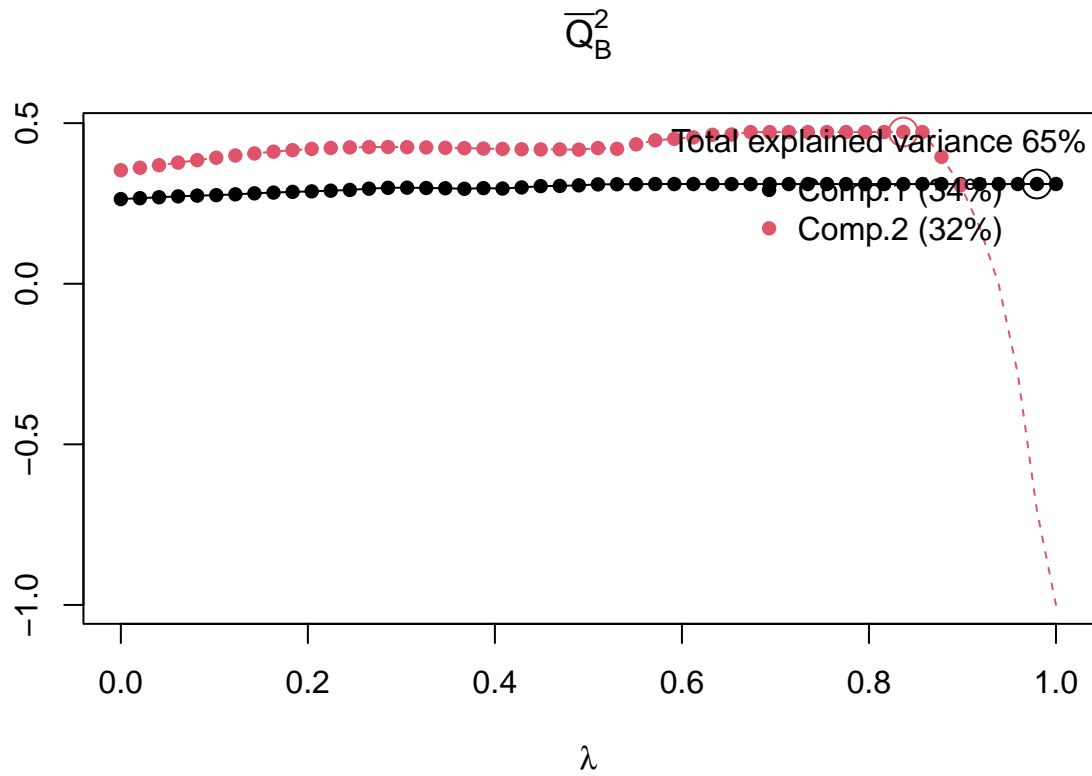
```
##                 ---------------
##                |    ddsPLS     |
## ===============                 ====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h  Q2 Q2h VarExpl VarExpl.Tot
##        0.98 0.34 0.34 0.3 0.3     34%         34%
##                                    ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##      lambda   R2 R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.81 0.64 0.3 0.62 0.46     32%         65%
##                                    ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                              ...component 3 not built!
## ====================                 ====================
##                   ================
```

4

$$\overline{Q}_B^2$$

Total explained variance 65%
Comp.1 (34%)
Comp.2 (32%)

$\lambda$

```r
model_1_lambda_2 <- ddsPLS(simu_1$X, simu_1$Y,
                    criterion = "Q2",
                    lambdas = seq(0,1,length.out = 50),
                    n_B = n_B,
                    verbose = T)
```

```
##                  _____
##                 |    ddsPLS    |
## ==================--------------==================
## Should we build component 1 ? Bootstrap pending...
##      lambda  R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.98 0.33 0.33 0.31 0.31    34%         34%
##                                 ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##      lambda  R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.84 0.63 0.29 0.63 0.47    32%         65%
##                                 ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                 ...component 3 not built!
## ====================              ==================
##                      ===============
```

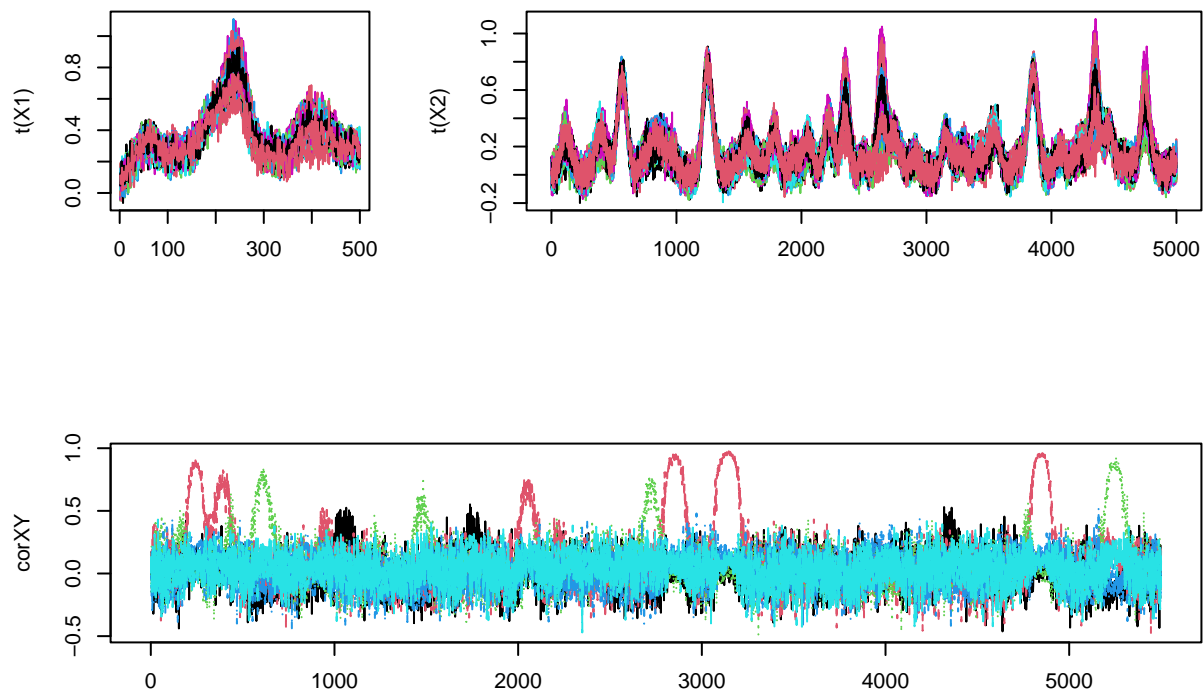$$\overline{Q}_B^2$$

## Different Simulations of Design 1 Data

There is a problem with `get_design_1`, `q` cannot take values other than 5.

```
data_1 <- get_design_1(n = 100, p = 1000, q = 5)

ddsPLS(data_1$X,data_1$Y,
                criterion = "Q2",
                lambdas = lambdas,
                n_B=n_B,
                verbose=T)
```

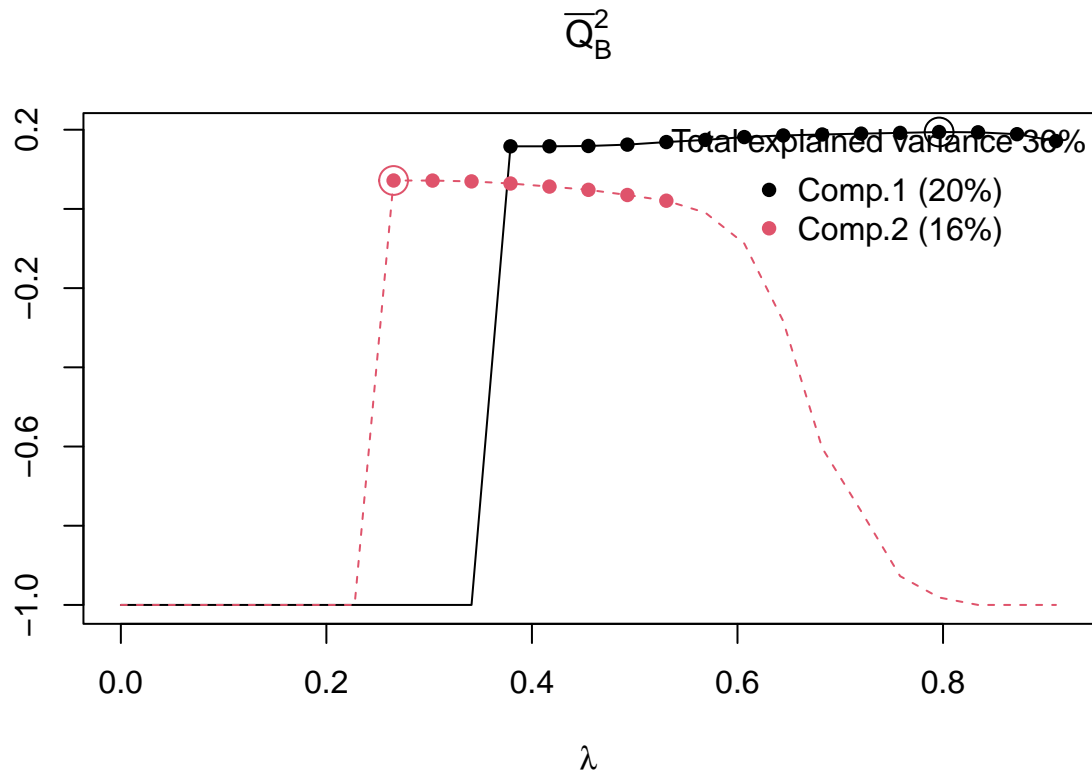## Design 2

```
simu_2.1 <- get_design_2(plot = T)
```

```
model_2.1 <- ddsPLS(simu_2.1$X$X1, simu_2.1$Y,
                    criterion = "Q2",
                    n_lambdas = 25,
                    verbose = TRUE)
```
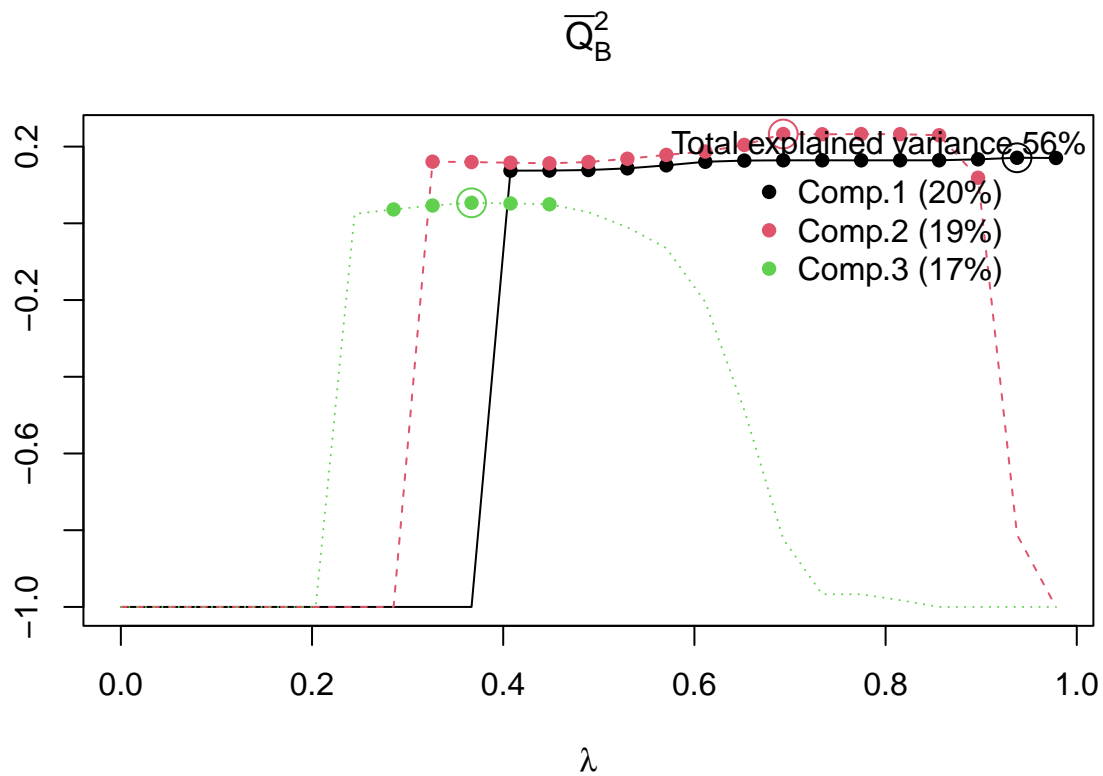
```
##                   --------------
##                  |    ddsPLS    |
## ========================----------------=====================
## Should we build component 1 ? Bootstrap pending...
##      lambda  R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.8 0.19 0.19 0.19 0.19    20%          20%
##                                       ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##      lambda  R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.27  0.4  0.2 0.24 0.07    16%          36%
##                                       ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                   ...component 3 not built!
## =====================              =====================
##                      ================
```
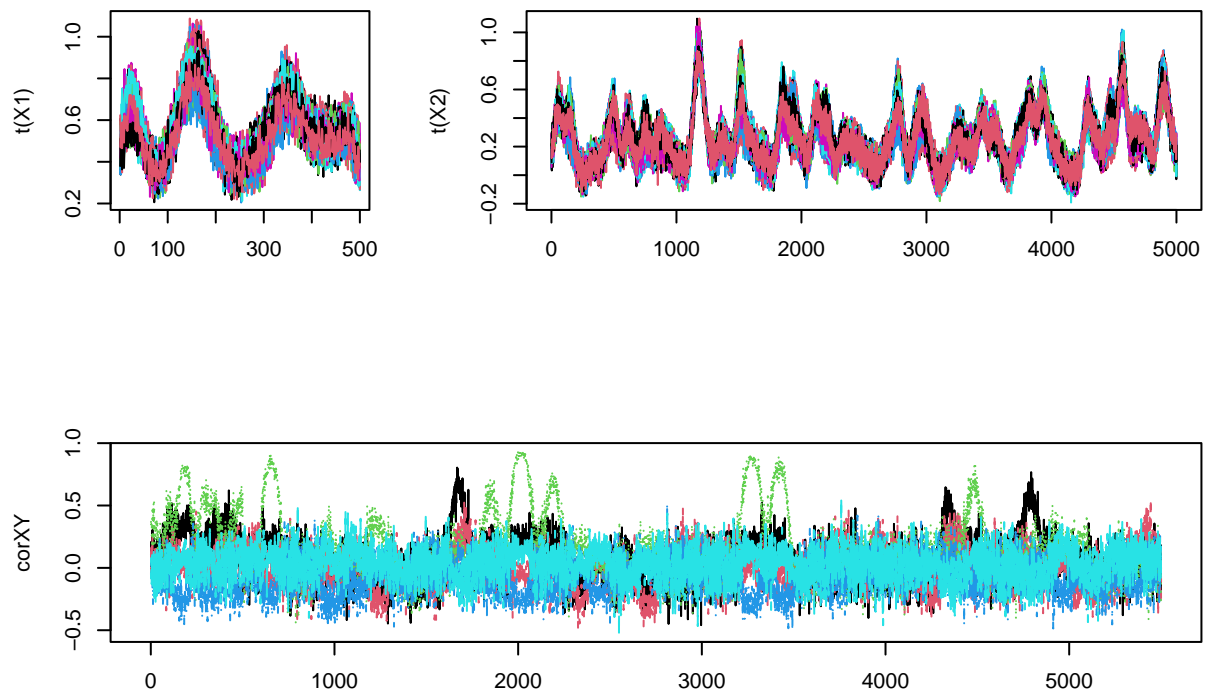
$$\overline{Q}_B^2$$

```r
model_2.2 <- ddsPLS(simu_2.1$X$X2, simu_2.1$Y,
                    criterion = "Q2",
                    n_lambdas = 25,
                    verbose = TRUE)
```

```
##                  --------------
##                 |    ddsPLS    |
## ======================--------------=====================
## Should we build component 1 ? Bootstrap pending...
##      lambda  R2 R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.94 0.2 0.2 0.17 0.17    20%          20%
##                                  ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.69 0.38 0.19 0.38 0.23    19%          39%
##                                  ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.37 0.57 0.19 0.39 0.05    17%          56%
##                                  ...component 3 built!
## Should we build component 4 ? Bootstrap pending...
##                                  ...component 4 not built!
## =====================              =====================
##                    ================
```

$$\overline{Q}_B^2$$

Total explained variance 56%
- Comp.1 (20%)
- Comp.2 (19%)
- Comp.3 (17%)

$\lambda$
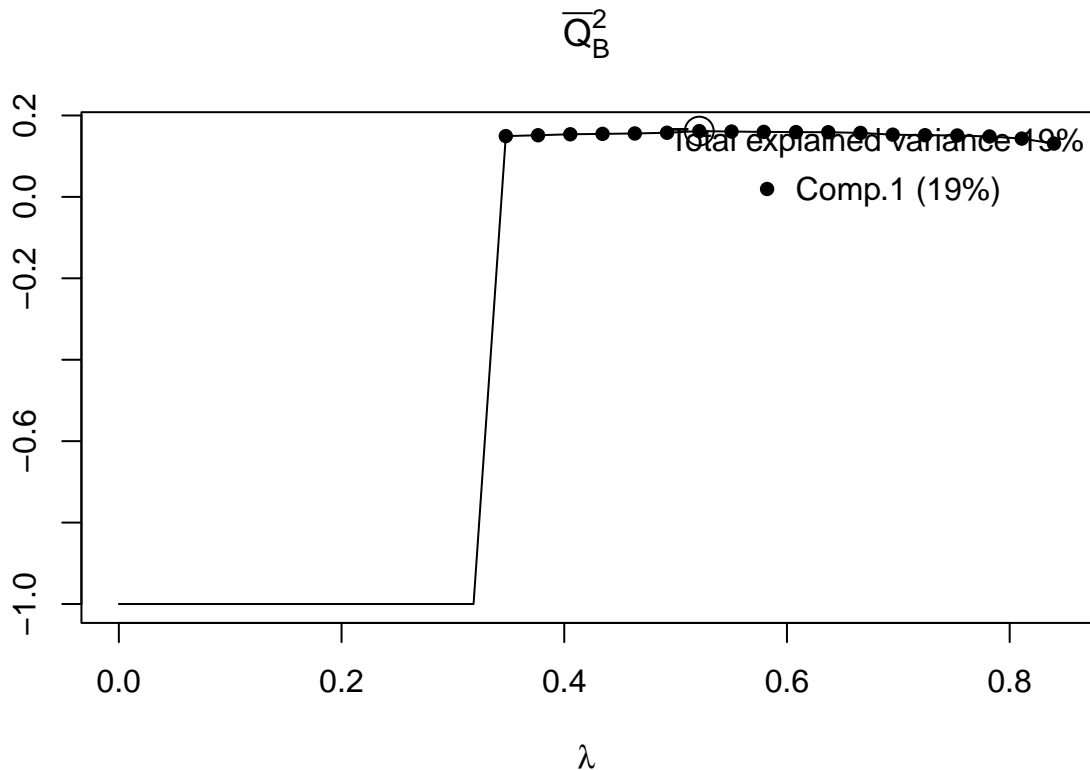
```
simu_2.2 <- get_design_2(seed = 2, ncpX = 20, plot = T)
```





```
model_2.3 <- ddsPLS(simu_2.2$Xs$X1, simu_2.2$Y,
                    criterion = "Q2",
                    n_lambdas = 30,
                    verbose = TRUE)
```
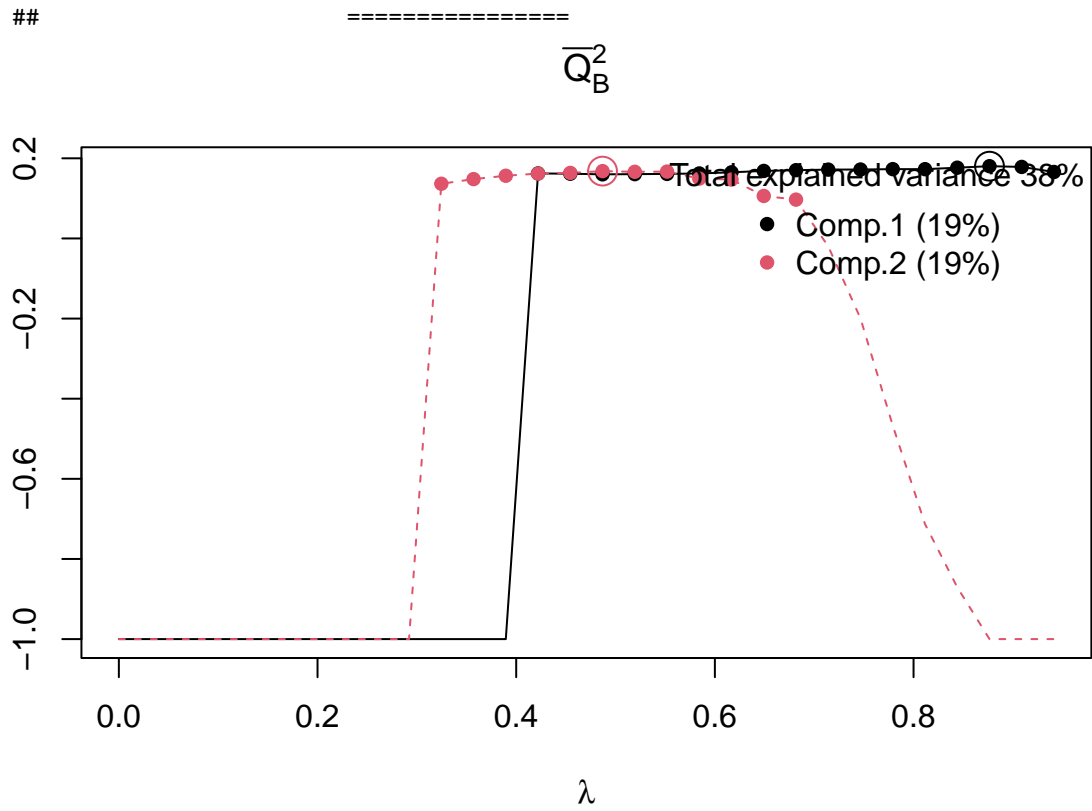
```
##                  --------------
```

```
##                         |   ddsPLS   |
## ====================---------------====================
## Should we build component 1 ? Bootstrap pending...
##       lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.52 0.22 0.22 0.16 0.16     19%          19%
##                                      ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                      ...component 2 not built!
## ====================          ====================
##               ================
```

$$\overline{Q}^2_B$$
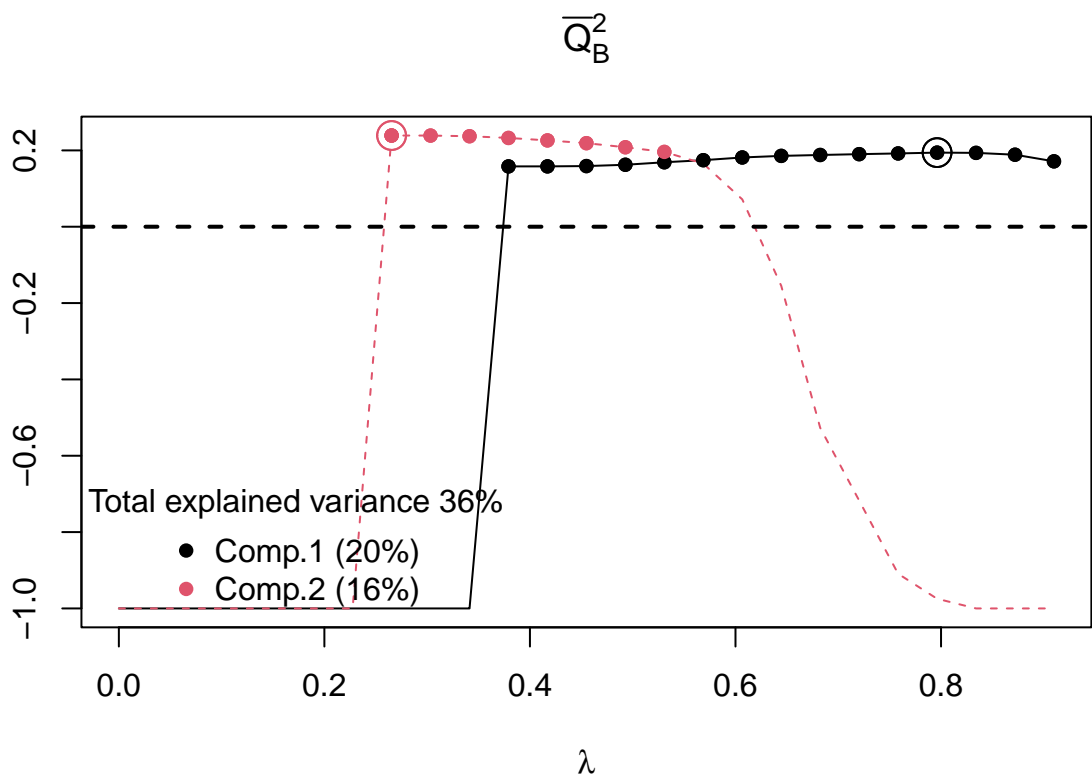


```
model_2.4 <- ddsPLS(simu_2.2$Xs$X2, simu_2.2$Y,
                    criterion = "Q2",
                    n_lambdas = 30,
                    verbose = TRUE)
```

```
##                         --------------
##                         |   ddsPLS   |
## ====================---------------====================
## Should we build component 1 ? Bootstrap pending...
##       lambda  R2 R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.88 0.2 0.2 0.18 0.18     19%          19%
##                                      ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##       lambda   R2 R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.49 0.39 0.2 0.32 0.17     19%          38%
##                                      ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                      ...component 3 not built!
## ====================          ====================
```

`##                       =================`

$$\overline{Q}^2_B$$



Model results can also be plotted using the `plot` function.

```
plot(model_2.1,type="Q2",legend.position = "bottomleft")
```

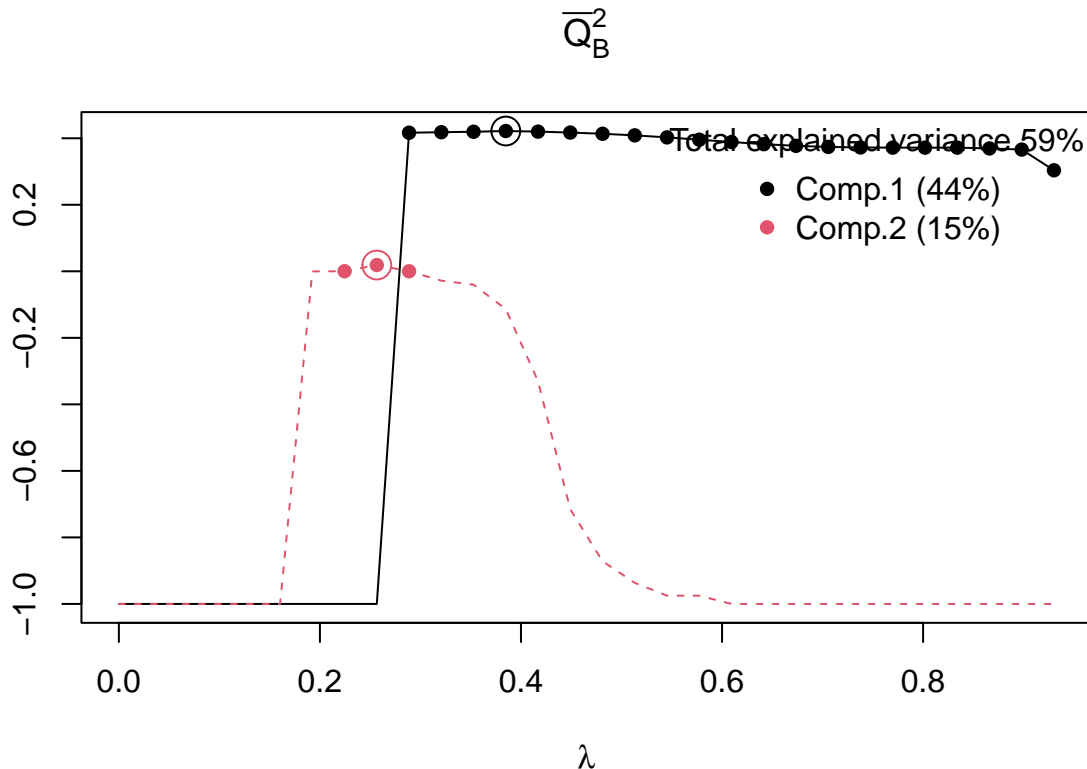$$\overline{Q}^2_B$$

## Get Data Simulation

The following `get_data` function is from the vignette for the `ddsPLS` package

The variable `eps` seems to relate the predictors and response, as well as `phi`. The dimension of `phi` specifies the number of latent variables.

```
data_3.1 <- get_data()

model_3.1 <- ddsPLS(data_3.1$X, data_3.1$Y,
                    criterion = "Q2",
                    n_lambdas = 30,
                    verbose = TRUE)
```
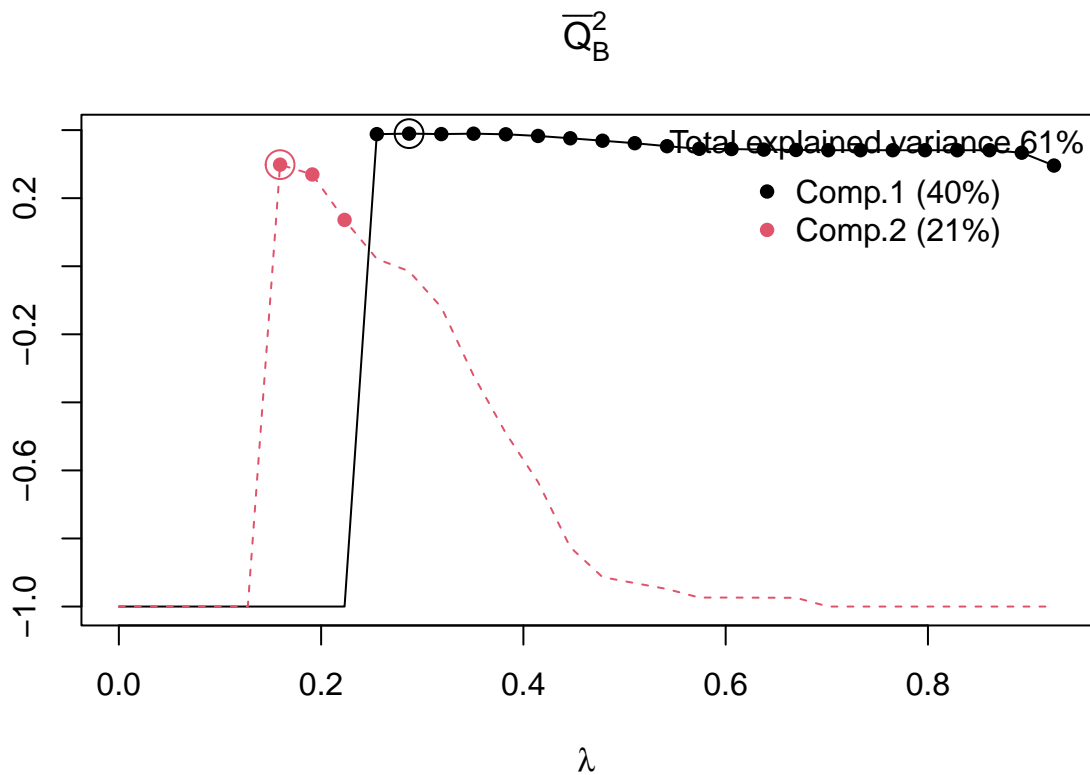
```
##                      --------------
##                     |    ddsPLS    |
## =====================--------------=====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2   R2h   Q2   Q2h VarExpl VarExpl.Tot
##        0.38 0.43 0.43 0.42 0.42     44%         44%
##                                  ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##      lambda   R2 R2h   Q2   Q2h VarExpl VarExpl.Tot
##        0.26 0.63 0.2 0.43 0.02     15%         59%
##                                  ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                  ...component 3 not built!
## =====================            =====================
##                      ================
```

$$\overline{Q}_B^2$$



```
data_3.2 <- get_data(p1 = 50, p2 = 50, p3 = 50, p = 250)
```

```
model_3.2 <- ddsPLS(data_3.2$X, data_3.2$Y,
                    criterion = "Q2",
                    n_lambdas = 30,
                    verbose = TRUE)
```
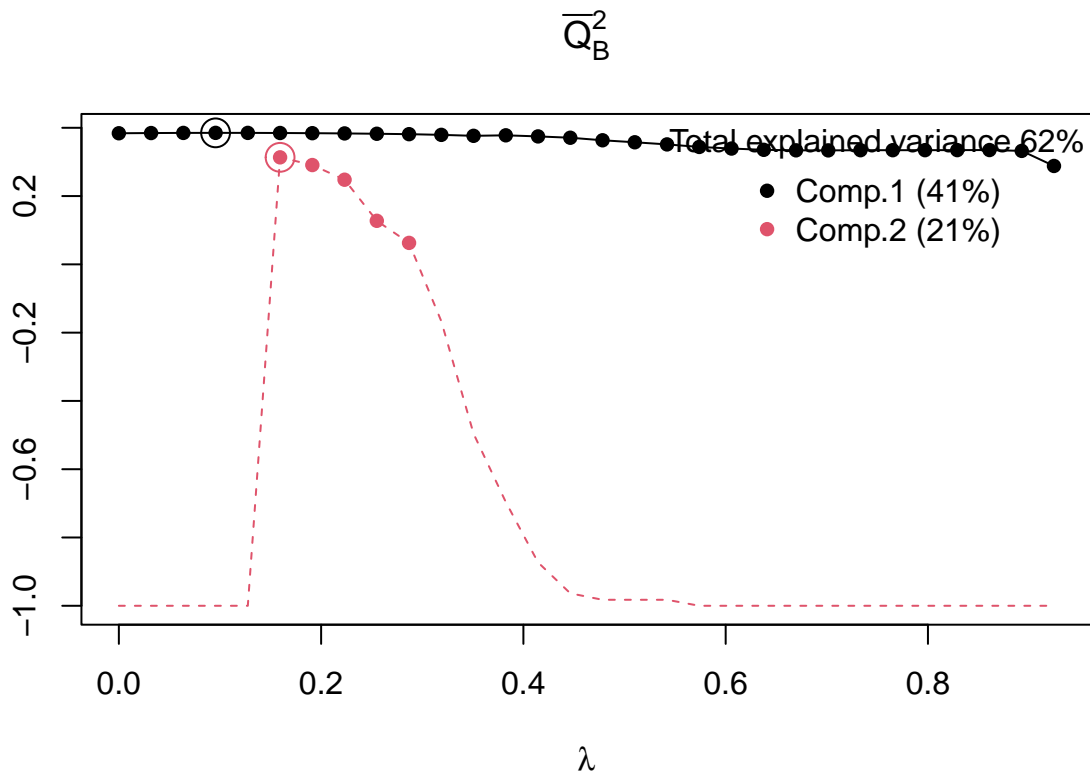
```
##                    _____
##                   |     ddsPLS     |
## ===================----------------====================
## Should we build component 1 ? Bootstrap pending...
##       lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.29 0.41 0.41 0.39 0.39     40%          40%
##                                      ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##       lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.16 0.62  0.2 0.55  0.3     21%          61%
##                                      ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                ...component 3 not built!
## =====================              ====================
##                    ================
```

$$\overline{Q}_B^2$$



```
model_3.3 <- ddsPLS(data_3.2$X, data_3.2$Y,
                    criterion = "Q2",
                    n_lambdas = 30,
                    verbose = TRUE,
                    LD = TRUE)
```

```
##                    _____
##                   |     ddsPLS     |
## ===================----------------====================
```

```
## Should we build component 1 ? Bootstrap pending...
##       lambda R2 R2h   Q2  Q2h VarExpl VarExpl.Tot
##          0.1 0.4 0.4 0.39 0.39     41%          41%
##                                     ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##       lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.16 0.62 0.21 0.57 0.31     21%          62%
##                                     ...component 2 built!
## Should we build component 3 ? Bootstrap pending...
##                                     ...component 3 not built!
## =====================            =====================
##                     ================
```

$$\overline{Q}_B^2$$



```
model_3.4 <- ddsPLS(data_3.2$X, data_3.2$Y,
                    criterion = "diffR2Q2",
                    n_lambdas = 30,
                    verbose = TRUE,
                    LD = TRUE)
```
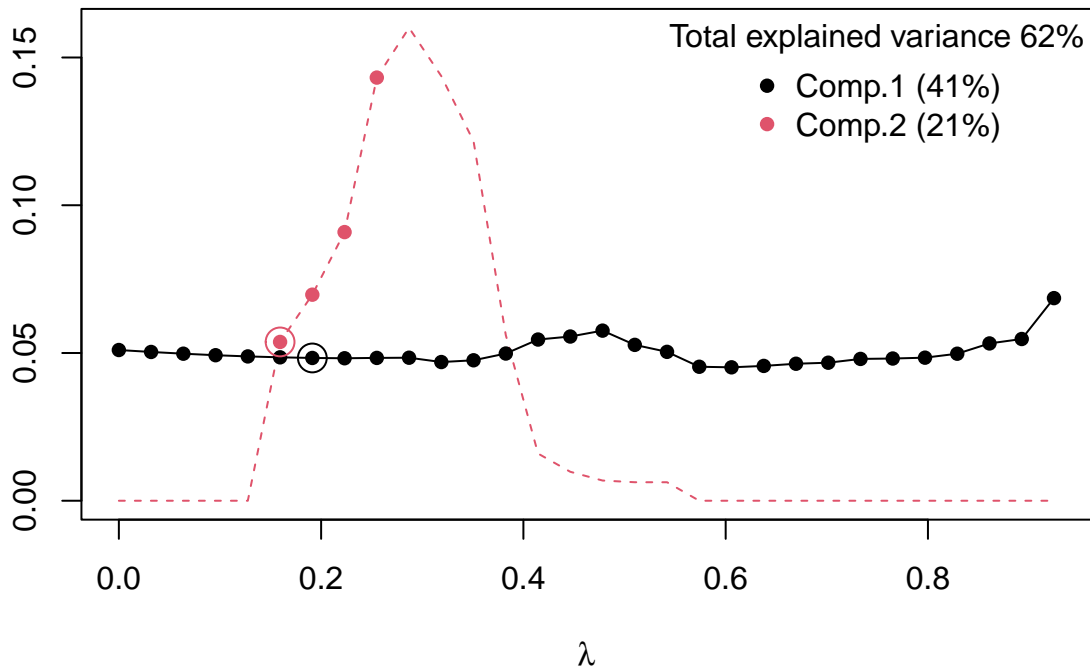
```
##                  --------------
##                 |    ddsPLS    |
## =====================--------------=====================
## Should we build component 1 ? Bootstrap pending...
##       lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.19 0.42 0.42 0.37 0.37     41%          41%
##                                     ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##       lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.16 0.62 0.21 0.57 0.29     21%          62%
##                                     ...component 2 built!
```

```
## Should we build component 3 ? Bootstrap pending...
##                                  ...component 3 not built!
## ====================              =====================
##                     ================
```

$$\overline{R}_B^2 - \overline{Q}_B^2$$



## Novel Simulations

The following code simulates a data set with 100 uncorrelated predictors and 1 response variable all sampled from a normal distribution.

```
Sigma <- diag(100)

sim_pred <- mvrnorm(n = 1000, mu = rep(0, 100), Sigma = Sigma)

sim_resp <- matrix(rnorm(1000), 1000, 1)

ddsPLS(sim_pred, sim_resp, verbose = TRUE)
```

```
##                     ---------------
##                    |    ddsPLS     |
## =====================---------------=====================
## Should we build component 1 ? Bootstrap pending...
##                                  ...component 1 not built!
##              ...no Q2r large enough for tested lambda.
## ====================              =====================
##                     ================

##
## Call:
## NULL
```

```
##
## No ddsPLS model built.
```

As expected no model is built as performance is awful. Interestingly message "no Q2r large enough for tested lambda" is given for justification, seems to suggest it checks just $Q^2$. Perhaps this just means that mean estimation performs better.

```
Sigma <- matrix(c(1,.75,.75,1),2,2)

n <- 20
p <- 5
p <- p - 2

sim_preds <- cbind(mvrnorm(n = n, rep(0, 2), Sigma), matrix(rep(0,n*p),n, p))

sim_resp <- as.matrix(apply(sim_preds,1,function(x) 5*x[1]+x[2]))

sim_preds <- sim_preds + matrix(rnorm(n*(p+2), sd = 0.6), n,(p+2))
sim_resp <- sim_resp + matrix(rnorm(n, sd = 0.8), n,1)
```
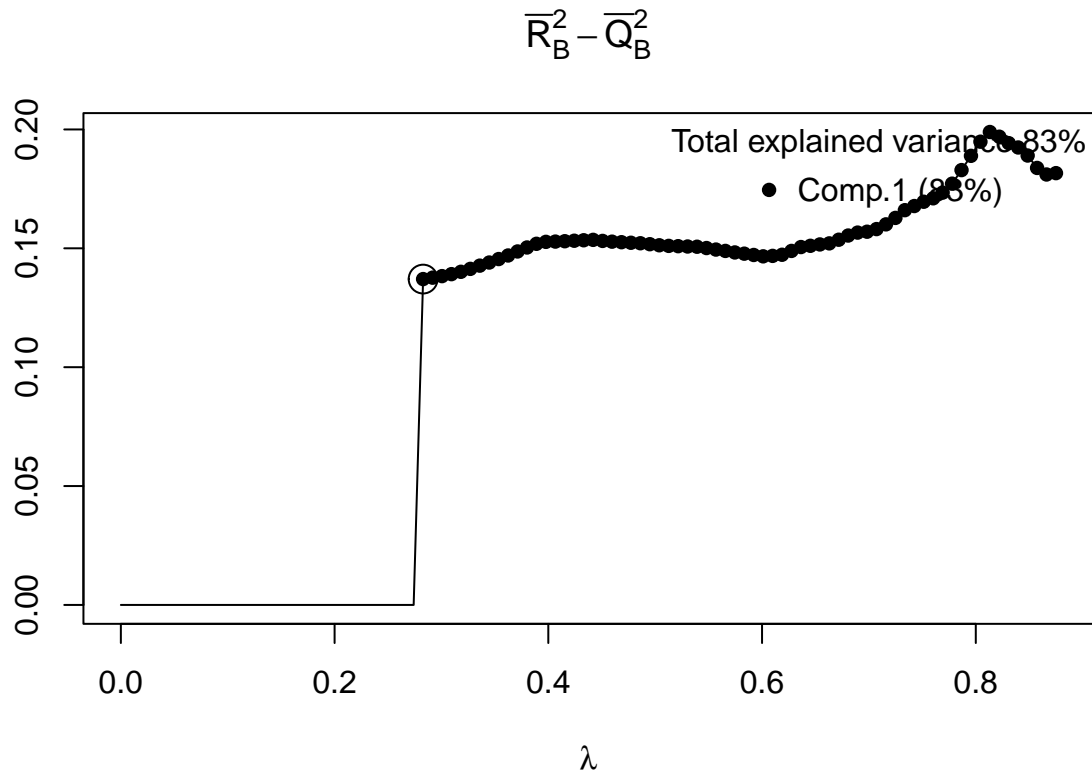
The above code simulates data with $n$ observations of $p$ predictors and 1 response variable. There are two predictors from which responses are linearly generated. Random noise is then added to the predictors and the response. As expected, the ddsPLS model performs very well.

```
pls_model <- ddsPLS(sim_preds, sim_resp, verbose = TRUE)
```

```
##                          ---------------
##                         |     ddsPLS     |
## =====================---------------=====================
## Should we build component 1 ? Bootstrap pending...
##       lambda   R2   R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.28 0.79 0.79 0.65 0.65     83%          83%
##                                      ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                      ...component 2 not built!
## =====================              =====================
##                       ===============
```

$$\overline{R}_B^2 - \overline{Q}_B^2$$

## Complex Simulated Data

The general structure of simulated data is $\mathbf{X} = \mathbf{A}^T\phi + \epsilon_X$ and $\mathbf{Y} = \mathbf{D}^T\phi + \epsilon_Y$. Note that $\phi$ provides the structure between the two. Code structures it as $\mathbf{X} = \phi\mathbf{A} + \epsilon_X$ and similarly for $\mathbf{Y}$. $\epsilon$ is added random error. $\mathrm{Cov}(\mathbf{X}, \mathbf{Y}) = \mathbf{D}^T\mathbf{A}$.

```r
sim_data <- function(n = 5, p = 10, q = 2, R = 5, x = 3, noise_weight = 1, D_method = "new") {

  # Ensures x<=R, if x>R the dimension of A is incompatible with phi
  if(x > R){
    x = R
  }

  # Creates A and D matrices
  A <- matrix(c(rep(rep(1,p),x), rep(rep(0,p),R-x)), ncol = p)

  if(D_method == "new") {
     D <- matrix(rep(1, R*q), nrow = R)
  } else {
    D <- diag(max(q, R))[1:R, 1:q]
  }

  d <- ncol(A)+nrow(A)+ncol(D)
  psi <- MASS::mvrnorm(n = n,mu = rep(0,d),Sigma = diag(d))
  phi <- psi[,1:nrow(A)]

  epsilon_X <- mvrnorm(n = dim(phi)[1],
                       rep(0, dim(A)[2]),
                       Sigma = noise_weight*diag(dim(A)[2]))
```

```
  epsilon_Y <- mvrnorm(n = dim(phi)[1],
                       rep(0, dim(D)[2]),
                       Sigma = noise_weight*diag(dim(D)[2]))

  X <- phi %*% A + epsilon_X
  Y <- phi %*% D + epsilon_Y



  list(X=X, Y=Y)
}
```

```
# This chunk isn't running, my guess is that setting a seed causes
# problems for the bootstrapping done to run ddsPLS

set.seed(1999)
sim_1 <- sim_data()

sim1_pls <- ddsPLS(sim_1$X, sim_1$Y,
                   n_B = 30,
                   n_lambdas = 25,
                   verbose = TRUE)
```

ddsPLS is not an appropriate method when working with a small sample size as

```
set.seed(1999)
sim_1 <- sim_data()
```

```
sim1_pls <- ddsPLS(sim_1$X, sim_1$Y,
                   verbose = TRUE)
```

```
set.seed(1894)
```

```
sim_2 <- sim_data(n = 100, p = 150, q = 5, R = 5, x = 4, D_method = "old")
```
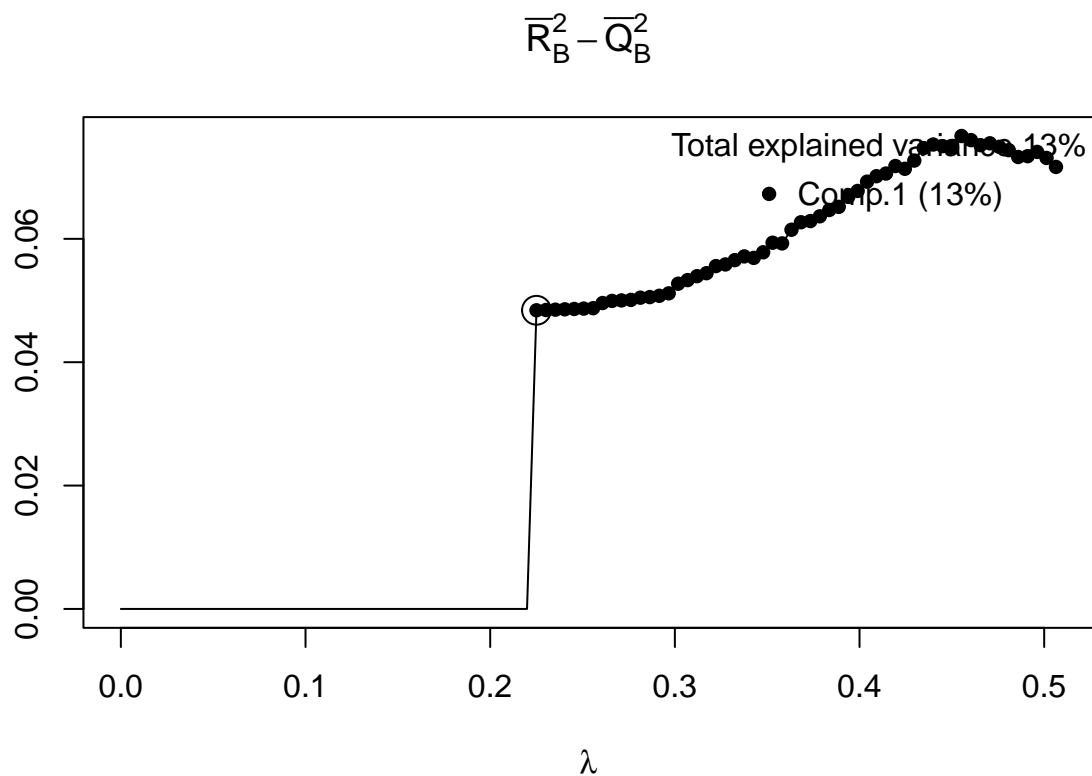
```
sim2_pls <- ddsPLS(sim_2$X, sim_2$Y,
                   verbose = TRUE)
```
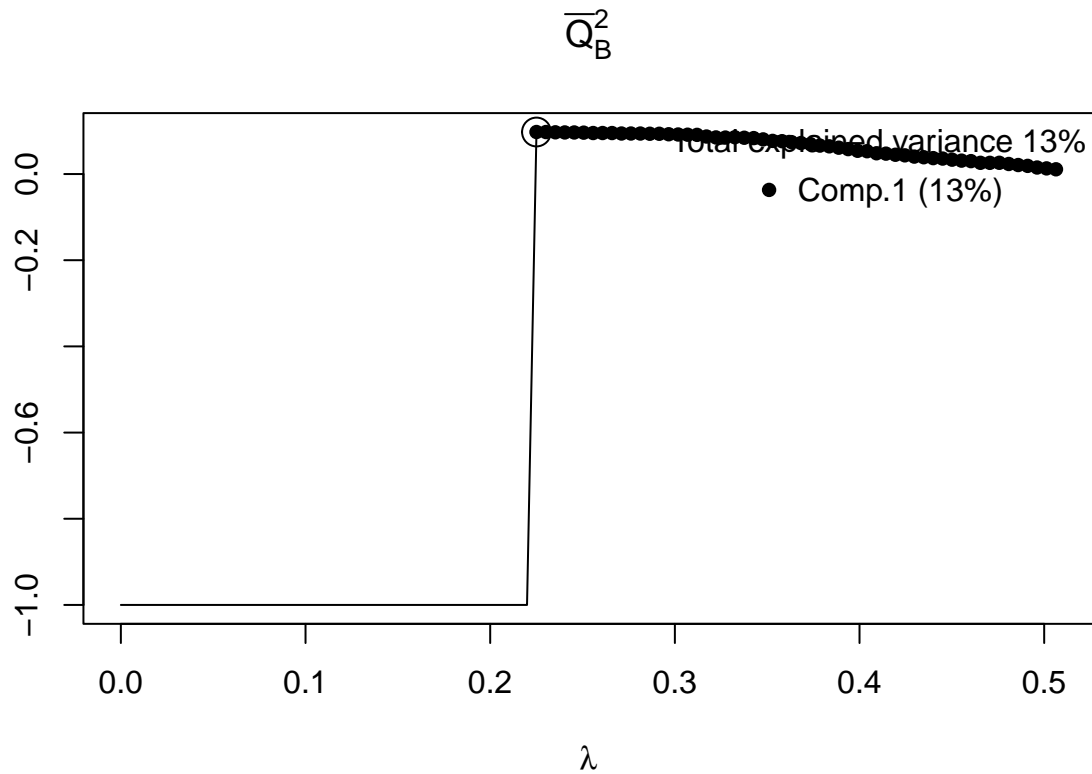
```
##                      --------------
##                     |    ddsPLS    |
## ====================----------------=====================
## Should we build component 1 ? Bootstrap pending...
##      lambda  R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.23 0.14 0.14 0.09 0.09    13%         13%
##                                      ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                      ...component 2 not built!
## ====================              =====================
##                      ================
```

$$\overline{R}_B^2 - \overline{Q}_B^2$$

```
sim2_pls <- ddsPLS(sim_2$X, sim_2$Y,
                   criterion = "Q2",
                   verbose = TRUE)
```

```
##                   --------------
##                  |    ddsPLS    |
## ===================--------------====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h  Q2 Q2h VarExpl VarExpl.Tot
##        0.23 0.14 0.14 0.1 0.1     13%         13%
##                                 ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                 ...component 2 not built!
## ====================              ====================
##                   ===============
```

$$\overline{Q}_B^2$$

Total explained variance 13%

● Comp.1 (13%)

```
set.seed(108)

sim_3 <- sim_data(n = 50, p = 200, q = 5, R = 25, x = 10, D_method = "old")

sim3_pls <- ddsPLS(sim_3$X, sim_3$Y,
                    verbose = TRUE)

##                      --------------
##                      |    ddsPLS    |
## =====================--------------=====================
## Should we build component 1 ? Bootstrap pending...
##                                  ...component 1 not built!
##          ...no Q2r large enough for tested lambda.
## =====================              =====================
##                      ===============

sim3_pls <- ddsPLS(sim_3$X, sim_3$Y,
                    criterion = "Q2",
                    verbose = TRUE)

##                      --------------
##                      |    ddsPLS    |
## =====================--------------=====================
## Should we build component 1 ? Bootstrap pending...
##                                  ...component 1 not built!
##          ...no Q2r large enough for tested lambda.
## =====================              =====================
##                      ===============

set.seed(43)
```

```r
sim_4 <- sim_data(n = 250, p = 1000, q = 10, R = 100, x = 95, D_method = "old")

sim4_pls <- ddsPLS(sim_4$X, sim_4$Y,
                   verbose = TRUE)
```

```
##                   --------------
##                  |    ddsPLS    |
## ====================--------------====================
## Should we build component 1 ? Bootstrap pending...
##                                  ...component 1 not built!
##          ...no Q2r large enough for tested lambda.
## ====================            ====================
##                  ==============
```

```r
sim4_pls <- ddsPLS(sim_4$X, sim_4$Y,
                   criterion = "Q2",
                   verbose = TRUE)
```

```
##                   --------------
##                  |    ddsPLS    |
## ====================--------------====================
## Should we build component 1 ? Bootstrap pending...
##                                  ...component 1 not built!
##          ...no Q2r large enough for tested lambda.
## ====================            ====================
##                  ==============
```

```r
set.seed(99)

sim_5 <- sim_data(n = 500, p = 500, q = 10, R = 100, x = 95, D_method = "old")

sim5_pls <- ddsPLS(sim_5$X, sim_5$Y,
                   verbose = TRUE)
```
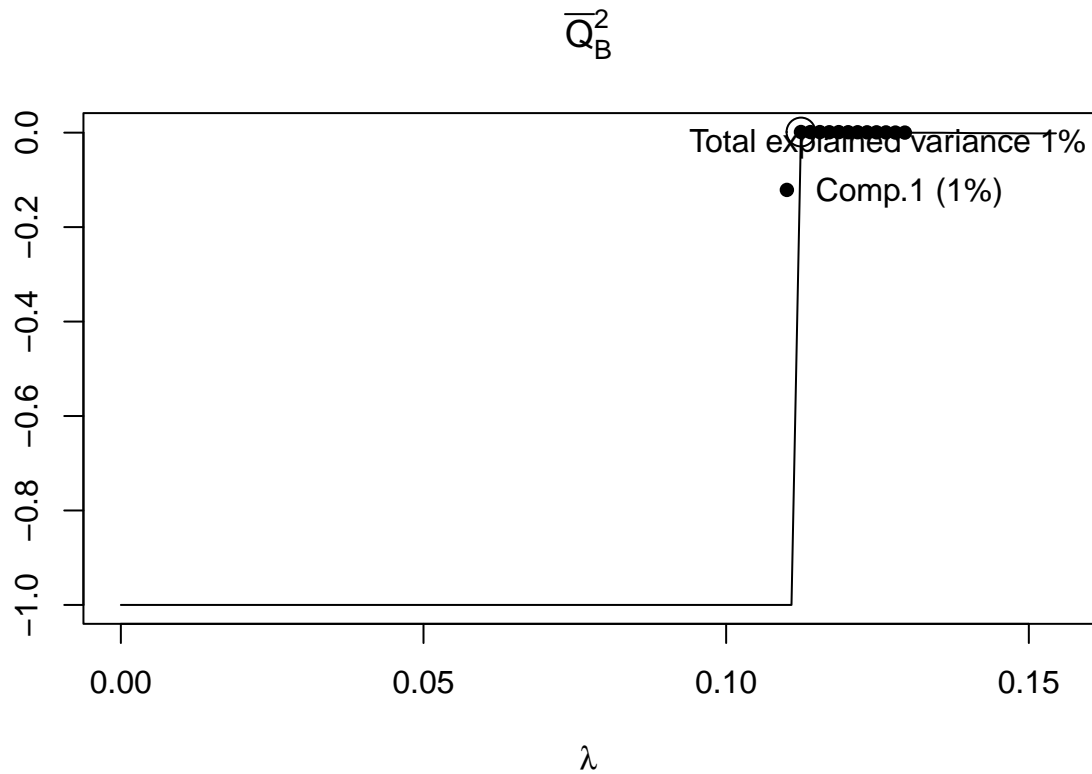
```
##                   --------------
##                  |    ddsPLS    |
## ====================--------------====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h Q2 Q2h VarExpl VarExpl.Tot
##        0.11 0.01 0.01  0   0      1%           1%
##                                  ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                  ...component 2 not built!
## ====================            ====================
##                  ==============
```
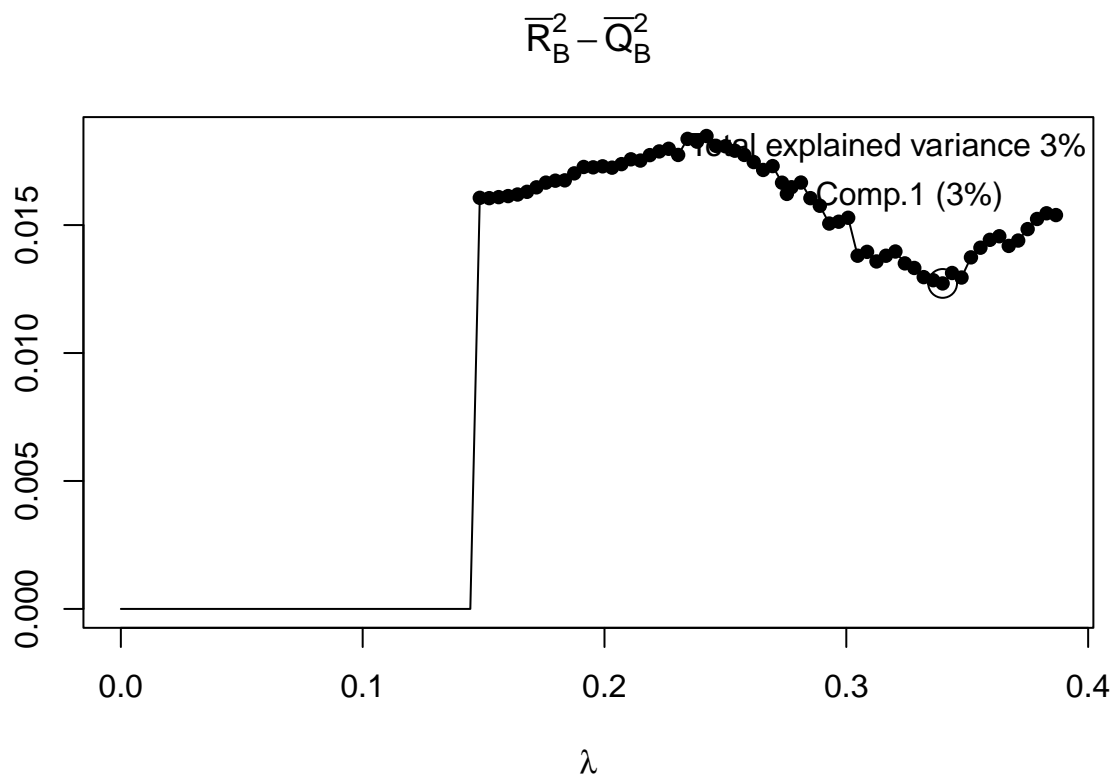
$$\overline{R}_B^2 - \overline{Q}_B^2$$

```
sim5_pls <- ddsPLS(sim_5$X, sim_5$Y,
                   criterion = "Q2",
                   verbose = TRUE)
```

```
##                  --------------
##                 |    ddsPLS    |
## ================---------------====================
## Should we build component 1 ? Bootstrap pending...
##     lambda   R2  R2h Q2 Q2h VarExpl VarExpl.Tot
##       0.11 0.01 0.01  0   0      1%          1%
##                                 ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                 ...component 2 not built!
## ====================            ====================
##                   ===============
```

$$\overline{Q}^2_B$$

```r
set.seed(99)

sim_6 <- sim_data(n = 250, p = 200, q = 10, R = 8, x = 7, D_method = "old")
```

```r
sim6_pls <- ddsPLS(sim_6$X, sim_6$Y,
                   verbose = TRUE)
```
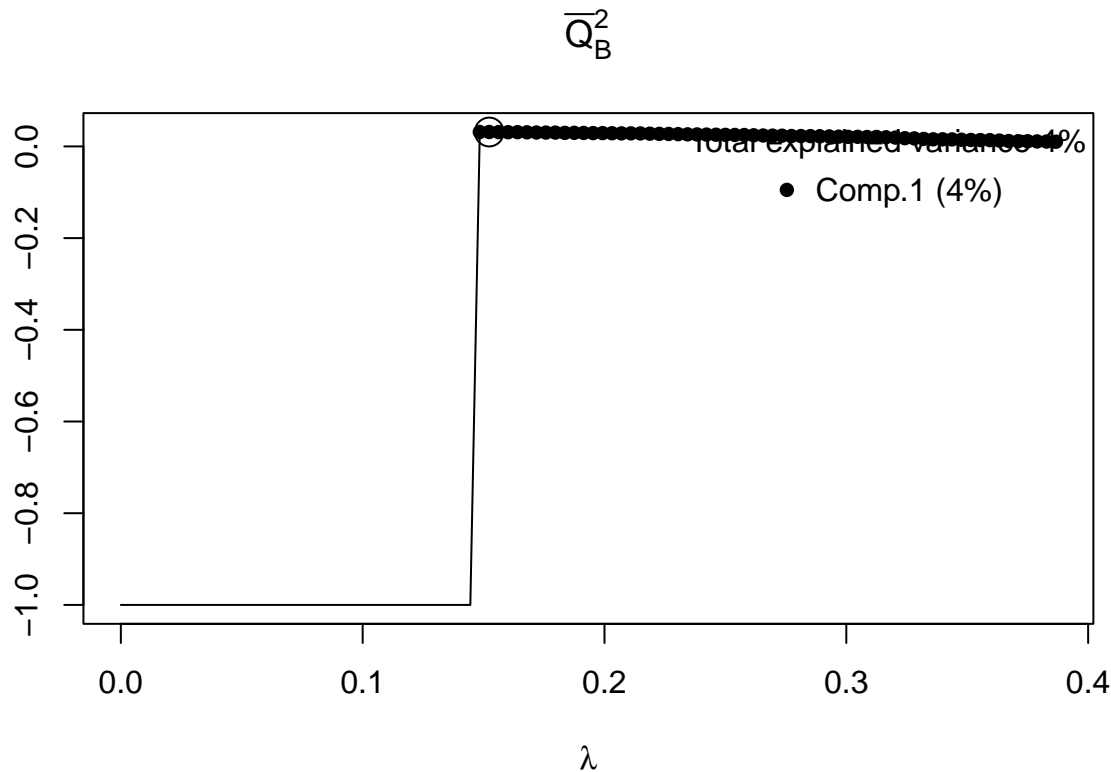
```
##                    ---------------
##                   |    ddsPLS    |
## =====================---------------====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.34 0.03 0.03 0.02 0.02     3%          3%
##                                    ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                    ...component 2 not built!
## ====================            ====================
##                    ================
```

$$\overline{R}_B^2 - \overline{Q}_B^2$$

Total explained variance 3%

Comp.1 (3%)

$\lambda$

```
sim6_pls <- ddsPLS(sim_6$X, sim_6$Y,
                   criterion = "Q2",
                   verbose = TRUE)
```

```
##                   --------------
##                  |   ddsPLS     |
## ==================--------------====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.15 0.05 0.05 0.03 0.03    4%            4%
##                                    ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                    ...component 2 not built!
## ====================              ====================
##                   ===============
```

$\overline{Q}^2_B$

It looks like there is a problem with the way I am simulating the data. All of the simulations have only explained a low amount of the variation in the data, perhaps due to the relation between `R` and `x`.

I changed `A` from a subset of a diagonal matrix to a matrix of all 1s and the function starts to work well. Running the code with different values of `q` and `noise_weight = 0` returns a model that perfectly explains the data.
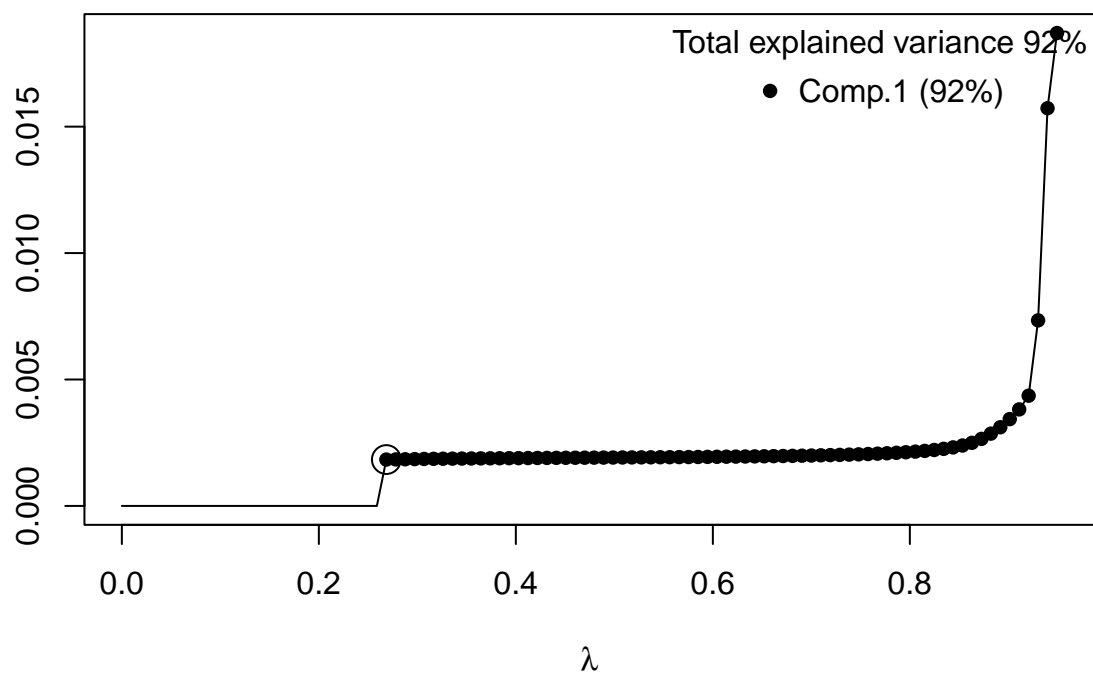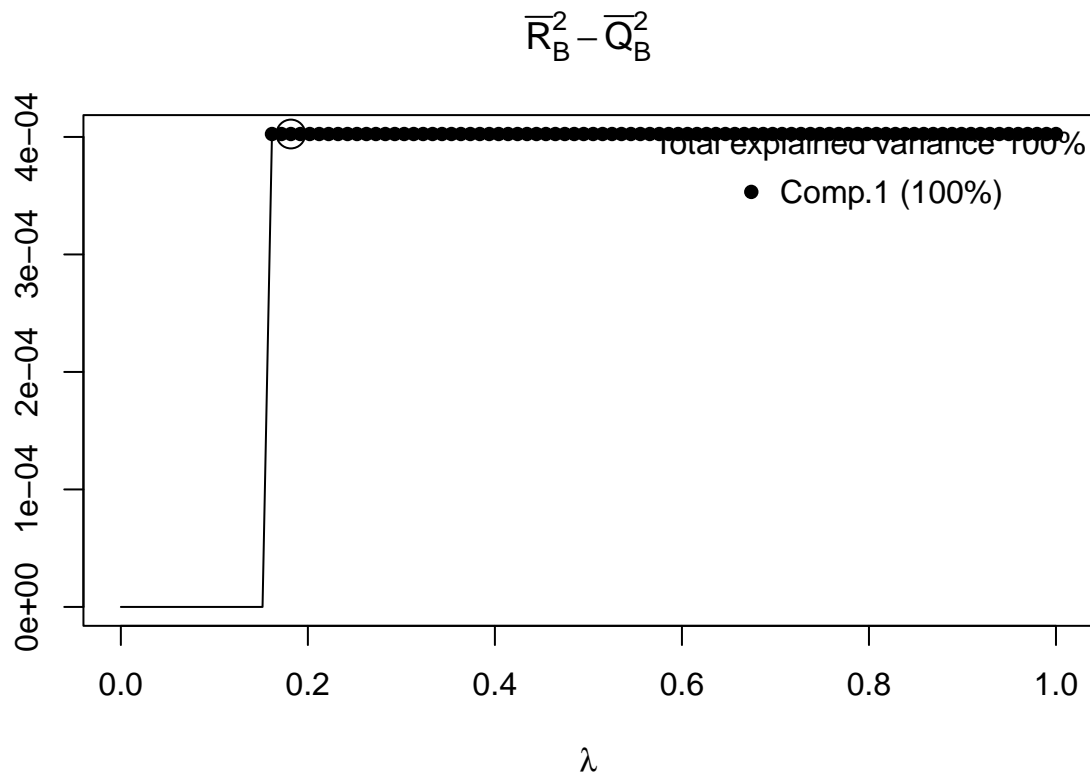
```
set.seed(12)

sim_7 <- sim_data(n = 100, p = 200, q = 5, noise_weight = 0.5)

sim7_pls <- ddsPLS(sim_7$X, sim_7$Y,
                   verbose = TRUE)
```

```
##                       ---------------
##                      |    ddsPLS     |
## =====================---------------=====================
## Should we build component 1 ? Bootstrap pending...
##       lambda  R2 R2h  Q2 Q2h VarExpl VarExpl.Tot
##         0.27 0.9 0.9 0.9 0.9     92%         92%
##                                     ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                     ...component 2 not built!
## =====================              =====================
##                       ===============
```
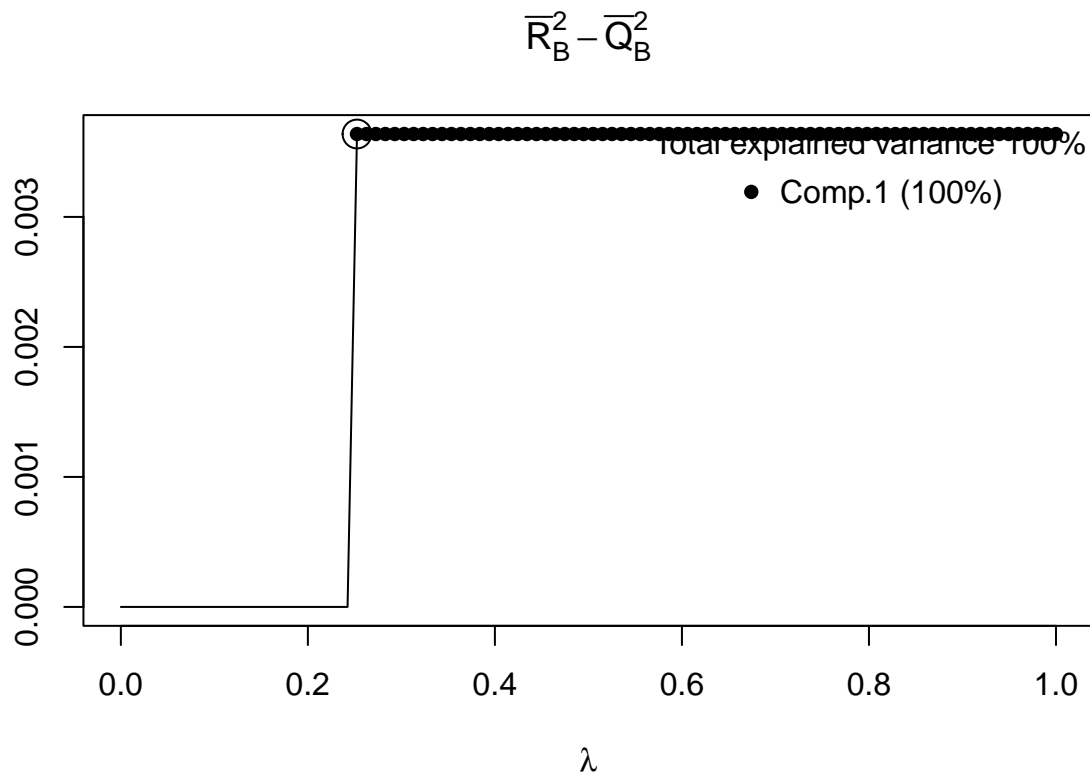
$$\overline{R}_B^2 - \overline{Q}_B^2$$

Total explained variance 92%
● Comp.1 (92%)

```
set.seed(12)

sim_8 <- sim_data(n = 100, p = 200, q = 150, noise_weight = 0)

sim8_pls <- ddsPLS(sim_8$X, sim_8$Y,
                   verbose = TRUE)
```

```
##                    ---------------
##                    |    ddsPLS    |
## ======================---------------======================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.18 0.99 0.99 0.99 0.99    100%        100%
##                                      ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                      ...component 2 not built!
## ====================            ====================
##                    ================
```

$$\overline{R}^2_B - \overline{Q}^2_B$$

If $q > n$, then the model is still able to perfectly explain the variance.
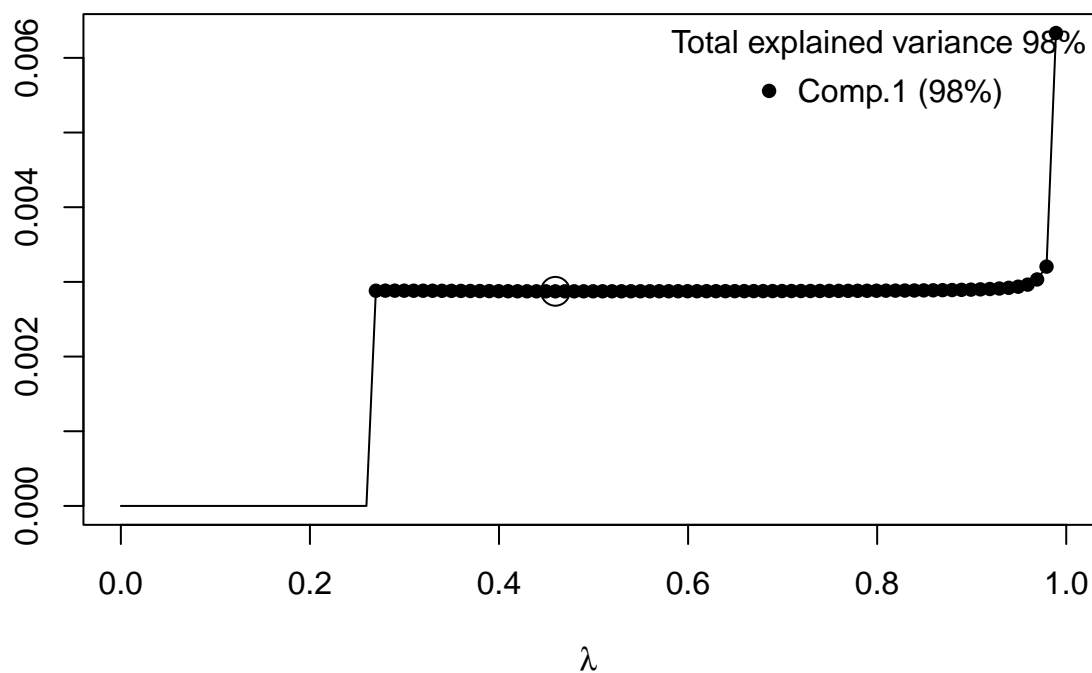
```
set.seed(12)

sim_9 <- sim_data(n = 50, p = 50, q = 150, noise_weight = 0)

sim9_pls <- ddsPLS(sim_9$X, sim_9$Y,
                    verbose = TRUE)
```

```
##                    ----------------
##                   |     ddsPLS      |
## ====================--------------====================
## Should we build component 1 ? Bootstrap pending...
##       lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##         0.25 0.97 0.97 0.97 0.97    100%        100%
##                                     ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                     ...component 2 not built!
## ====================              ====================
##                    ================
```

$$\overline{R}^2_B - \overline{Q}^2_B$$

Still works when $n < q$ and $p < q$.
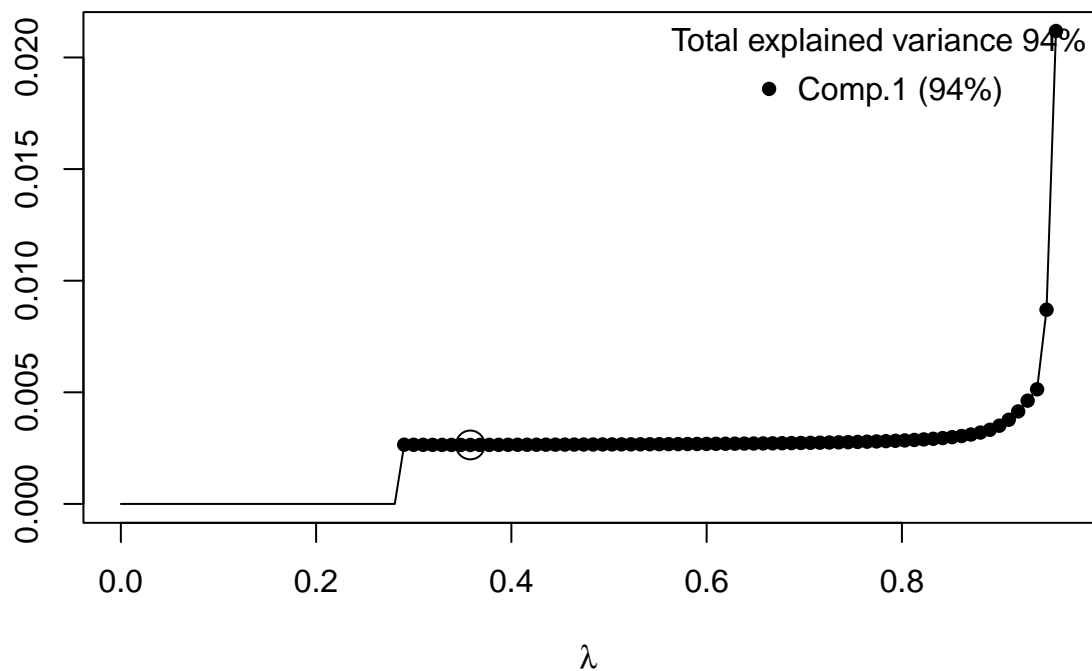
```
set.seed(12)
```

```
sim_9 <- sim_data(n = 100, p = 200, q = 5, noise_weight = 0.1)
```

```
sim9_pls <- ddsPLS(sim_9$X, sim_9$Y,
                   verbose = TRUE)
```

```
##                 --------------
##                |   ddsPLS    |
## ===================--------------===================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.46 0.97 0.97 0.97 0.97     98%          98%
##                                  ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                 ...component 2 not built!
## ===================              ===================
##                 ===============
```
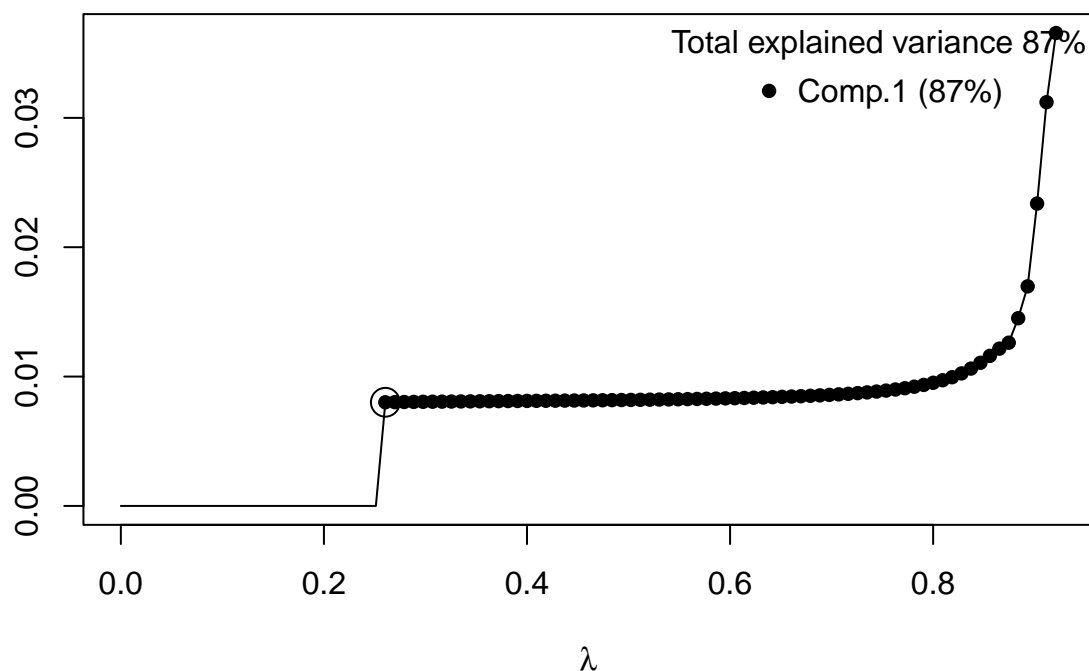
$$\overline{R}^2_B - \overline{Q}^2_B$$

```
set.seed(167)

sim_10 <- sim_data(n = 100, p = 200, q = 5, noise_weight = 0.25)

sim10_pls <- ddsPLS(sim_10$X, sim_10$Y,
                    verbose = TRUE)
```

```
##                      ---------------
##                      |    ddsPLS    |
## ======================---------------=====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.36 0.93 0.93 0.92 0.92    94%          94%
##                                      ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                      ...component 2 not built!
## =====================            =====================
##                      ===============
```

$$\overline{R}_B^2 - \overline{Q}_B^2$$

```r
set.seed(902)

sim_11 <- sim_data(n = 100, p = 200, q = 5, noise_weight = 0.5)

sim11_pls <- ddsPLS(sim_11$X, sim_11$Y,
                    verbose = TRUE)
```

```
##                      --------------
##                     |   ddsPLS     |
## =====================----------------=====================
## Should we build component 1 ? Bootstrap pending...
##      lambda   R2  R2h   Q2  Q2h VarExpl VarExpl.Tot
##        0.26 0.86 0.86 0.85 0.85     87%         87%
##                                   ...component 1 built!
## Should we build component 2 ? Bootstrap pending...
##                                   ...component 2 not built!
## =====================               =====================
##                      ================
```

$$\overline{R}_B^2 - \overline{Q}_B^2$$

As noise increases, the model performance predictably increases.

```r
var_func <- function(noise_weight){
    sim <- sim_data(n = 100, p = 200, q = 5, noise_weight = noise_weight)
    mod <- ddsPLS(sim$X, sim$Y)
    if(!is.null(tail(mod$varExplained$Cumu, n=1))) {
        return(c(noise_weight, tail(mod$varExplained$Cumu, n=1)))
    }
}

apply(matrix(c(1:10/10), nrow = 1), MARGIN = 2, var_func)
```
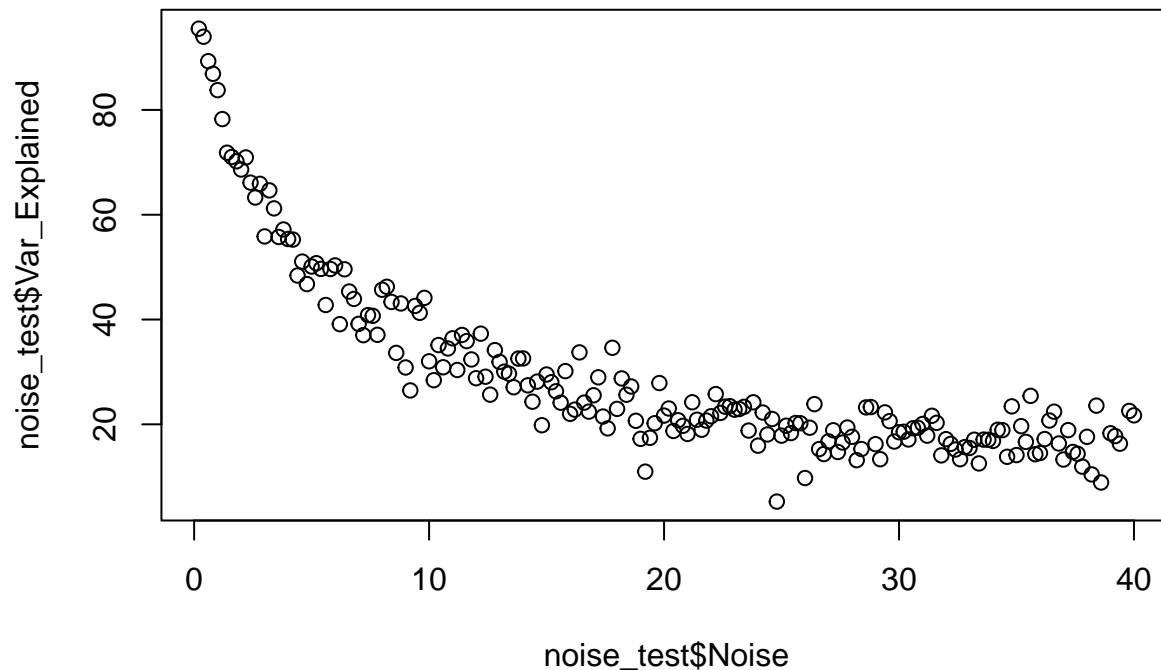
```
##          [,1]    [,2]     [,3]     [,4]     [,5]     [,6]     [,7]     [,8]
## [1,]   0.1000  0.20000  0.30000  0.40000  0.50000  0.60000  0.70000  0.80000
## [2,]  98.3625 96.46001 93.08608 89.61578 91.20199 90.06957 85.96129 85.71387
##          [,9]    [,10]
## [1,]   0.90000  1.00000
## [2,]  83.33498 83.20425
```

```r
noise_test <- apply(matrix(c(1:200/5), nrow = 1), MARGIN = 2, var_func)
```

```r
noise_test <- as.data.frame(do.call(rbind, noise_test))
colnames(noise_test) <- c("Noise", "Var_Explained")

plot(noise_test$Noise, noise_test$Var_Explained)
```

As we would predict, model performance decreases as the amount of noise increases. Initially, model performance decreases at a fairly rapid rate before becoming more gradual. Eventually, we would expect the percent variance explained to go to 0.
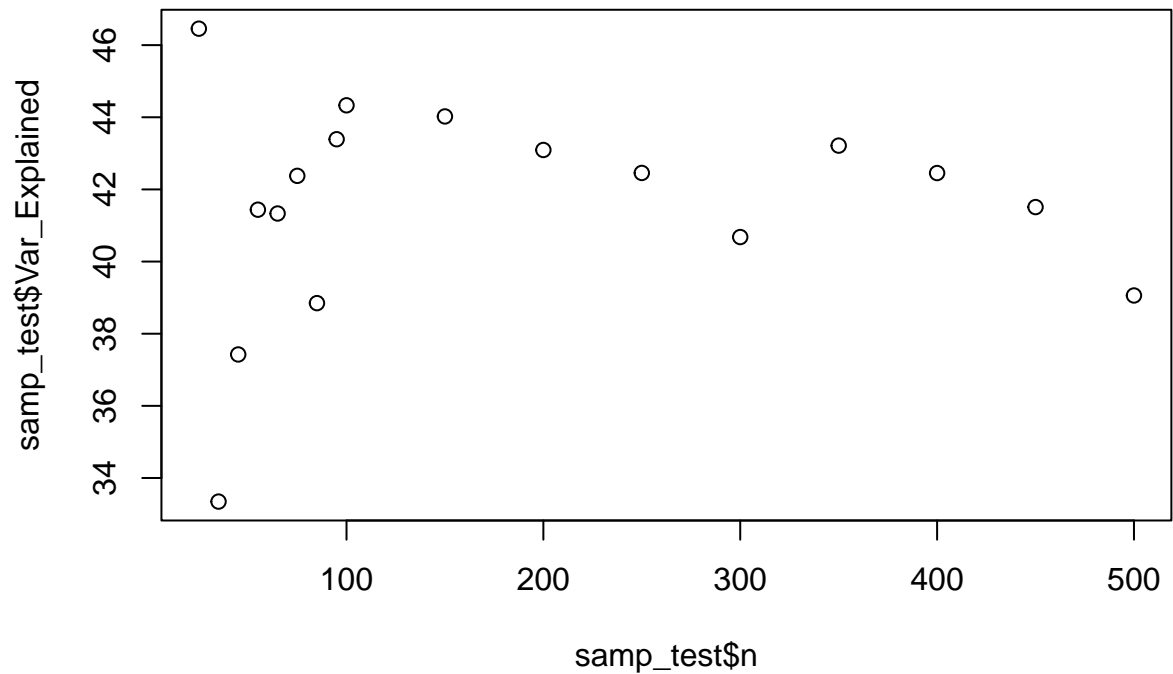
```
samp_func <- function(n, noise_weight){
  sim <- sim_data(n = n, p = 100, q = 5, noise_weight = noise_weight)
  mod <- ddsPLS(sim$X, sim$Y)
  if(!is.null(tail(mod$varExplained$Cumu, n=1))) {
    return(c(n, tail(mod$varExplained$Cumu, n=1)))
  }
}
```

```
samp_test <- apply(matrix(c(seq(from = 25, to = 95, by = 10),
                            seq(from = 100, to = 500, by = 50)),
                          nrow = 1),
                   MARGIN = 2,
                   samp_func,
                   noise_weight = 7)
```
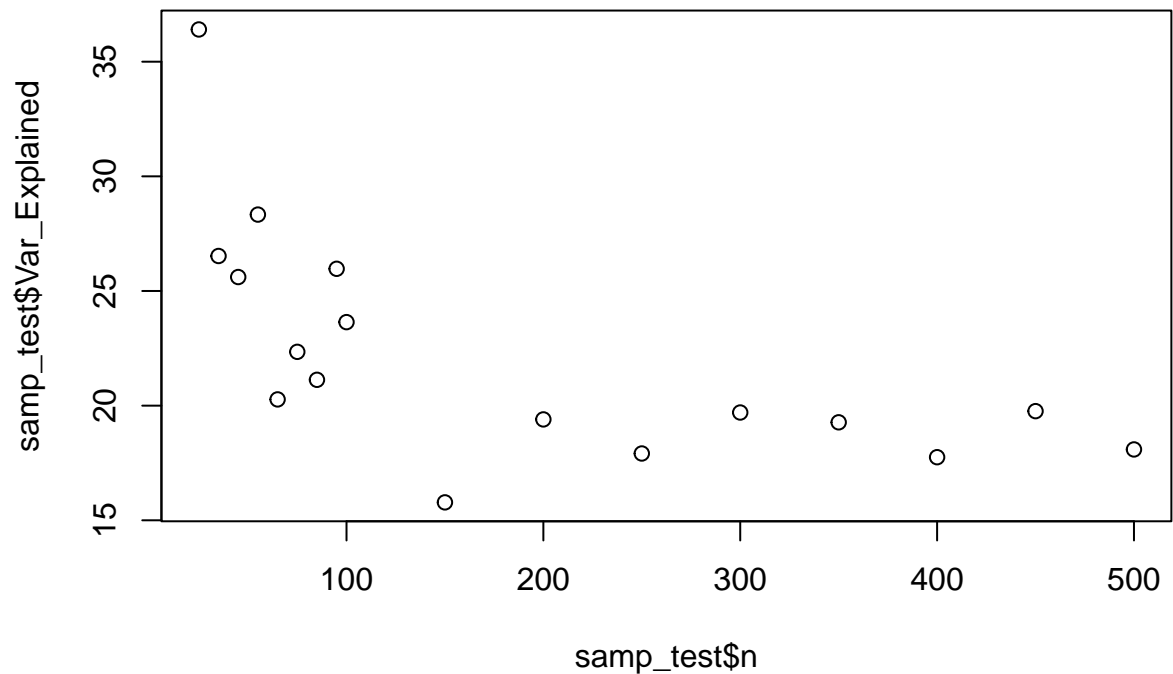
```
samp_test <- as.data.frame(t(samp_test))
colnames(samp_test) <- c("n", "Var_Explained")

plot(samp_test$n, samp_test$Var_Explained)
```

```
samp_test <- apply(matrix(c(seq(from = 25, to = 95, by = 10),
                            seq(from = 100, to = 500, by = 50)),
                          nrow = 1),
                   MARGIN = 2,
                   samp_func,
                   noise_weight = 20)
```

```
samp_test <- as.data.frame(t(samp_test))
colnames(samp_test) <- c("n", "Var_Explained")

plot(samp_test$n, samp_test$Var_Explained)
```
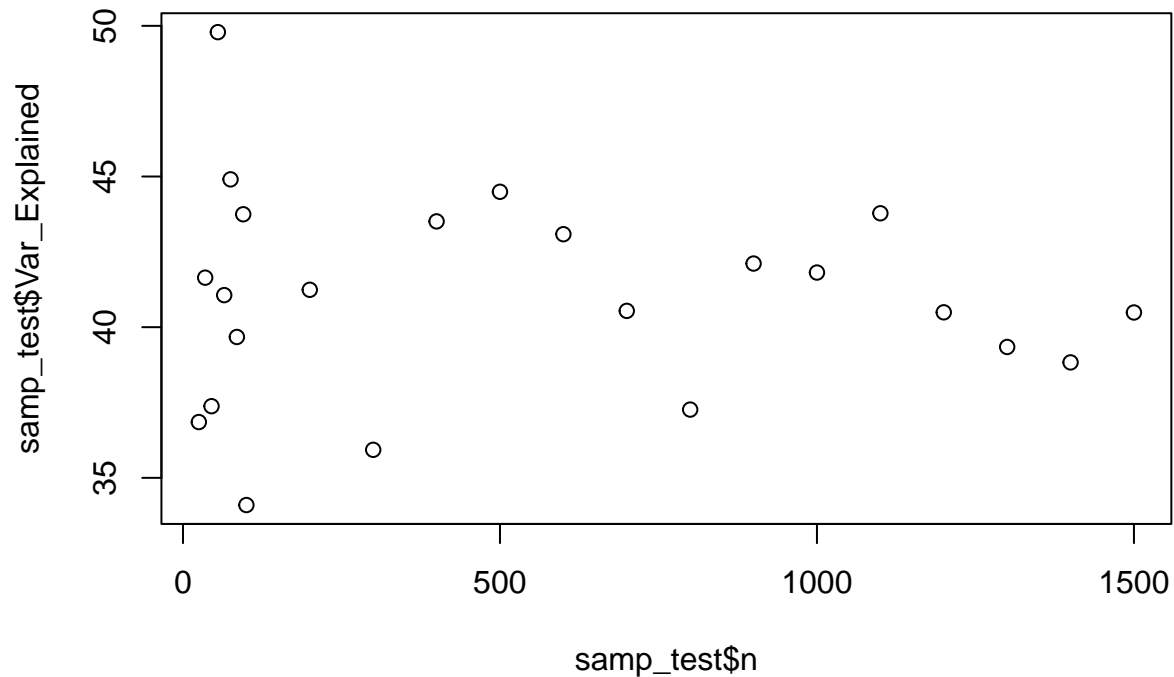
```
samp_test <- apply(matrix(c(seq(from = 25, to = 95, by = 10),
                            seq(from = 100, to = 1500, by = 100)),
                          nrow = 1),
                   MARGIN = 2,
                   samp_func,
                   noise_weight = 7)
```

```
samp_test <- as.data.frame(t(samp_test))
colnames(samp_test) <- c("n", "Var_Explained")

plot(samp_test$n, samp_test$Var_Explained)
```



Model performance seems to be much more variable at a low sample size before stabilizing. It looks like there may be a slight improvement as model size increases however this would need more inquiry. I am curious as to why models with small sample size can perform much better than those based on a larger sample size.