

SVA Recount

Christopher Lo

9/19/2022

Introduction

What happens to the number of Surrogate Variables as we increase the number of studies in recount3?

To illustrate this, we examined 4 studies that involved case/control prostate tumor and normal samples, with a minimal of 25 samples per study. For each study, we examined the number of SVs (using method = "be"), and we also aggregated them one by one and examined the number of SVs.

1. SRP118614: "Overall design: Matched high-grade (GS=7(4+3)) prostate tumor and adjacent normal specimens from 16 patients (8 AAM and 8 EAM) were subjected to two replicate runs of RNA-sequencing."
2. SRP002628: "Overall design: We sequenced the transcriptome (polyA+) of 20 prostate cancer tumors and 10 matched normal tissues using Illumina GAI platform. Then we used bioinformatic approaches to identify prostate cancer specific aberrations which include gene fusion, alternative splicing, somatic mutation."
3. SRP212704: "Overall design: Strand specific total RNA seq was performed using frozen patient matched prostate cancer tissue in biological duplicates. Purpose: The goal of present study is to compare transcript level changes between normal and tumor of same individuals"
4. SRP027258: "We utilized RNA sequencing to test the hypothesis that SFN modifies the expression of genes that are critical in prostate cancer progression. Normal prostate epithelial cells, and androgen-dependent and androgen-independent prostate cancer cells were treated with 15 μ M SFN and the transcriptome was determined at 6 and 24 hour time points."

We also have a second case/control set on TB samples.

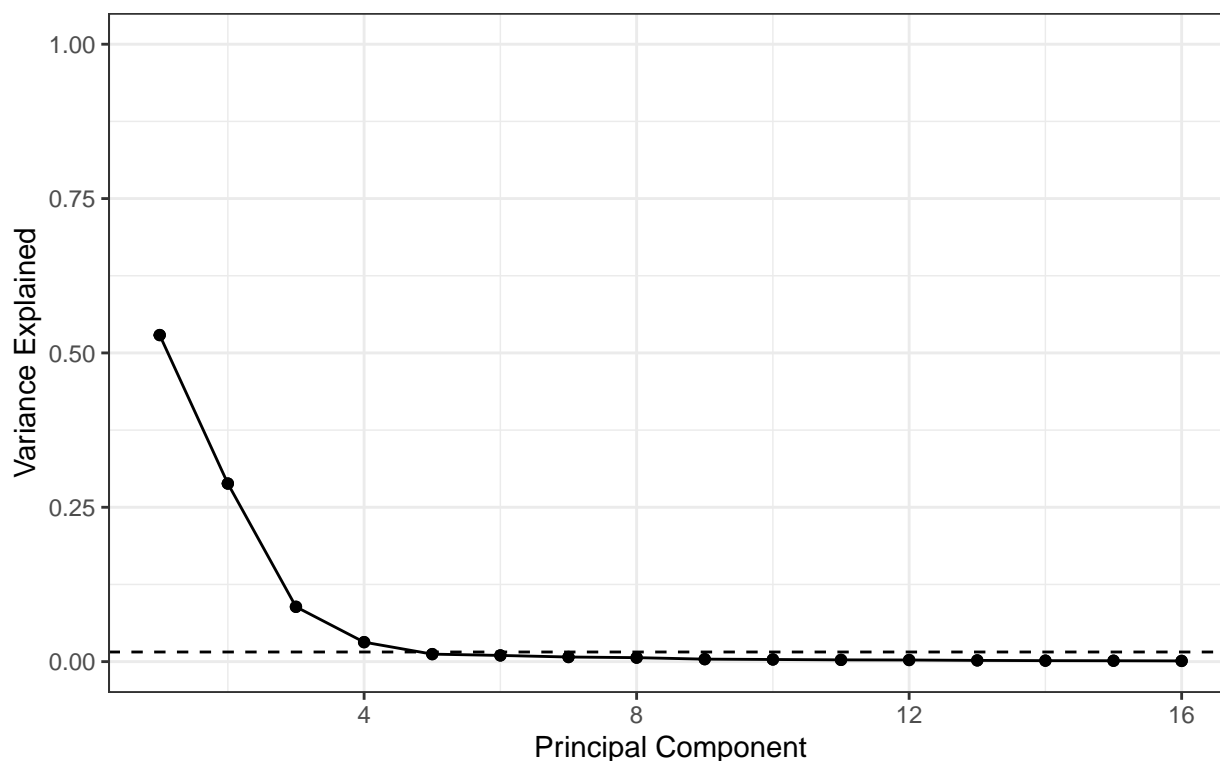
1. SRP098758 (n=428): "Samples are collected from subjects in a household contact study after a person comes back to the household, diagnosed with TB. Samples are collected every 6 months up to, 18 months. Some people go on to develop TB (cases) where as some others do not (controls). Here we are trying to establish a gene signature to predict the occurrence of TB."
2. ERP115010 (n=360): "Close contacts of active TB were defined as individuals with a cumulative duration of exposure of greater than eight hours in a confined space to the index case prior to initiation of treatment. Known human immunodeficiency virus (HIV)-positive patients were excluded. At enrollment, interferon gamma release assays (IGRAs) were done using the QuantiFERON-TB Plus assay (Qiagen, Germany), and peripheral blood was collected into Tempus tubes for whole genome transcriptional profiling by RNA sequencing. Participants who progressed to active TB were identified by linkage with the national electronic TB register. Local case notes were reviewed to identify individuals who had received preventative treatment. This submission contains data from n=360 adult participants, of which n=9 progressed to TB during a median follow-up time of 1.9 years. [note: we identify case base on IGRA assay result, which does not distinguish active vs. latent TB.]

3. SRP126580 (n=54): “Overall design: We undertook RNA Sequencing (RNA-Seq) of our earlier Berry et al. 2010 (GSE19444 and GSE19442) cohorts and additionally set up a prospective cohort study at Leicester (UK) in subject groups of incident TB and recent TB contacts, respectively. In the Leicester cohort, we performed systematic longitudinal sampling and clinical characterisation first, to validate our TB signature using RNA-Seq in a new and independent cohort of individuals with active TB and LTBI, and secondly to provide longitudinal data in a low TB incidence setting. All samples in this series were re-analyzed from GSE19444. There are links on each sample page to the original sample.”
[QC note: 28% of genes have count of 0]

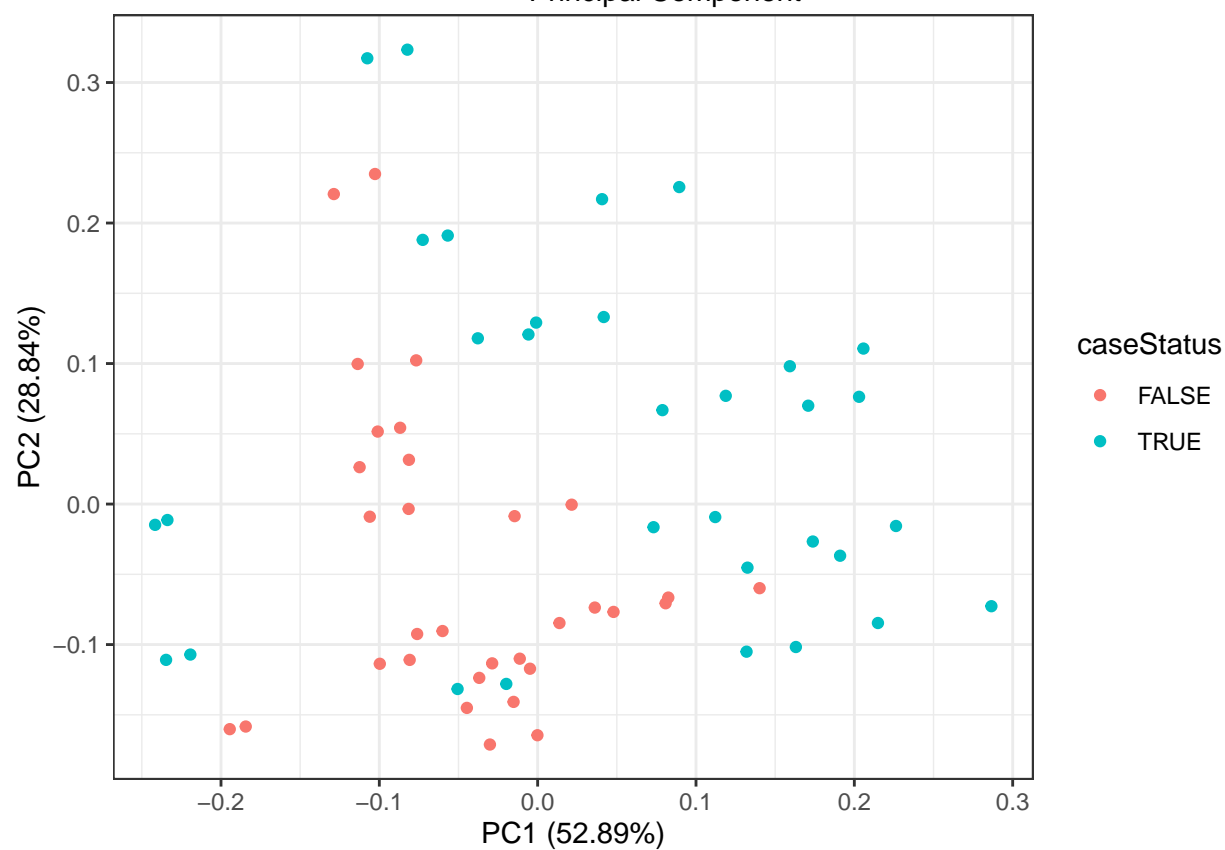
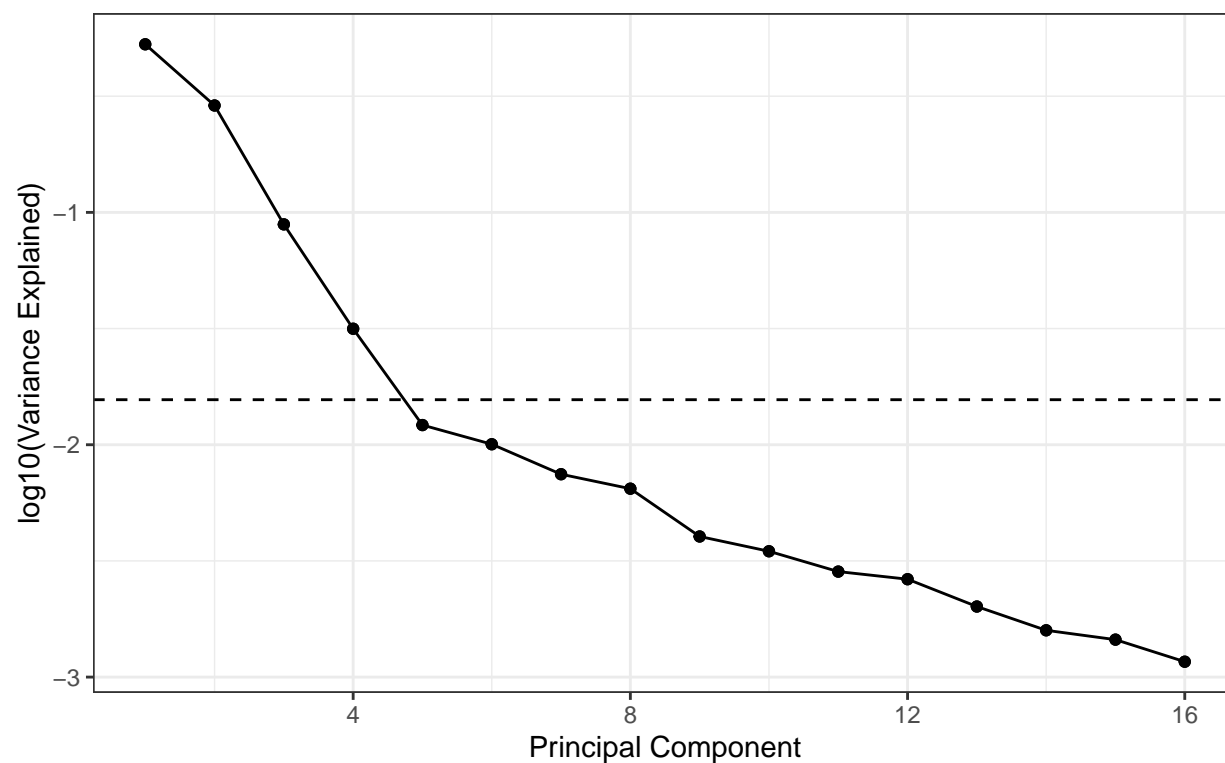
```
## 1
## [1] "SizeFactor distribution:"
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.4696 0.7496  1.0697  1.1054  1.4532  2.0903
```

SRP118614

Number of SVs: 3

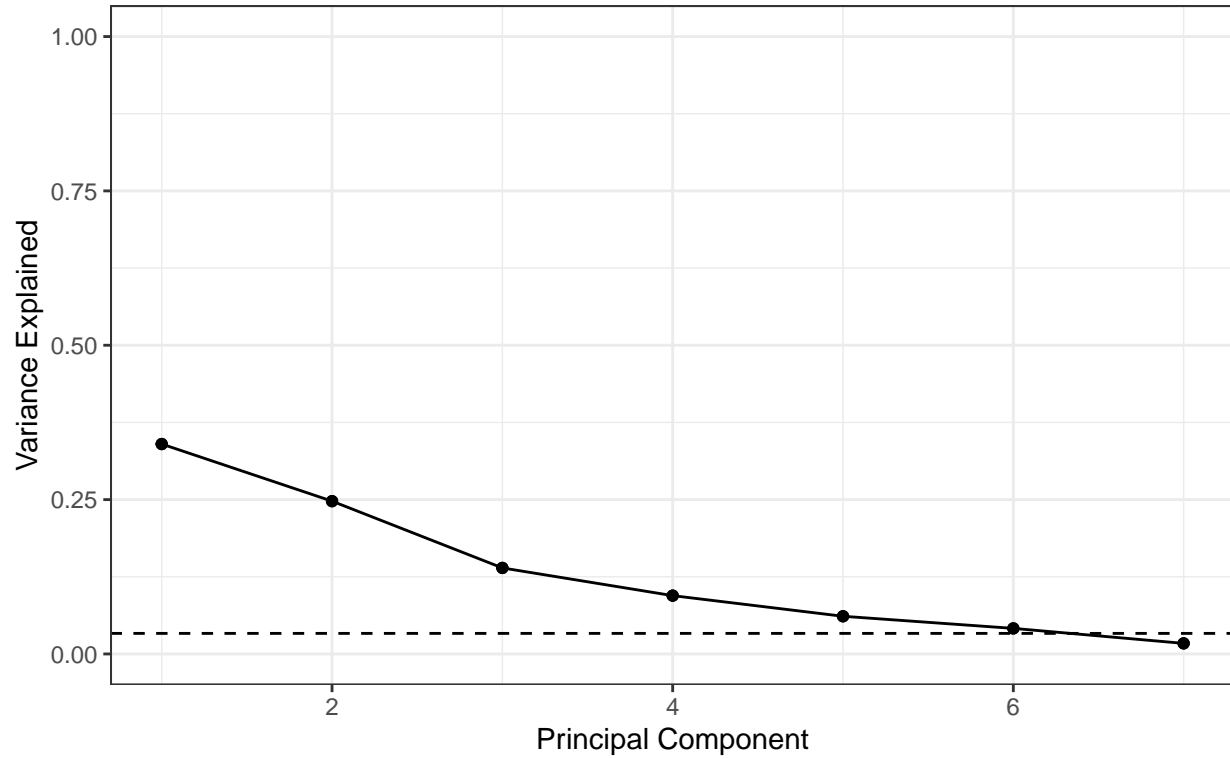


SRP118614
Number of SVs: 3

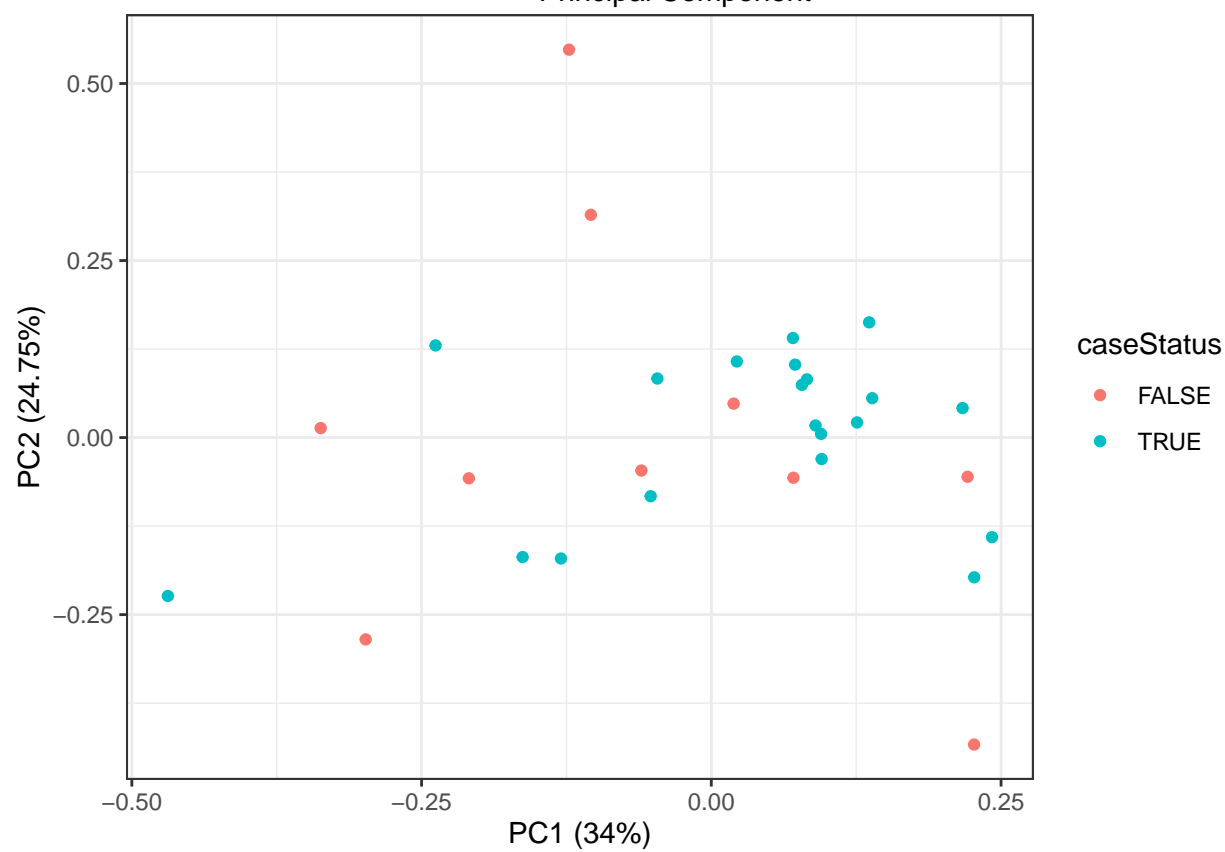
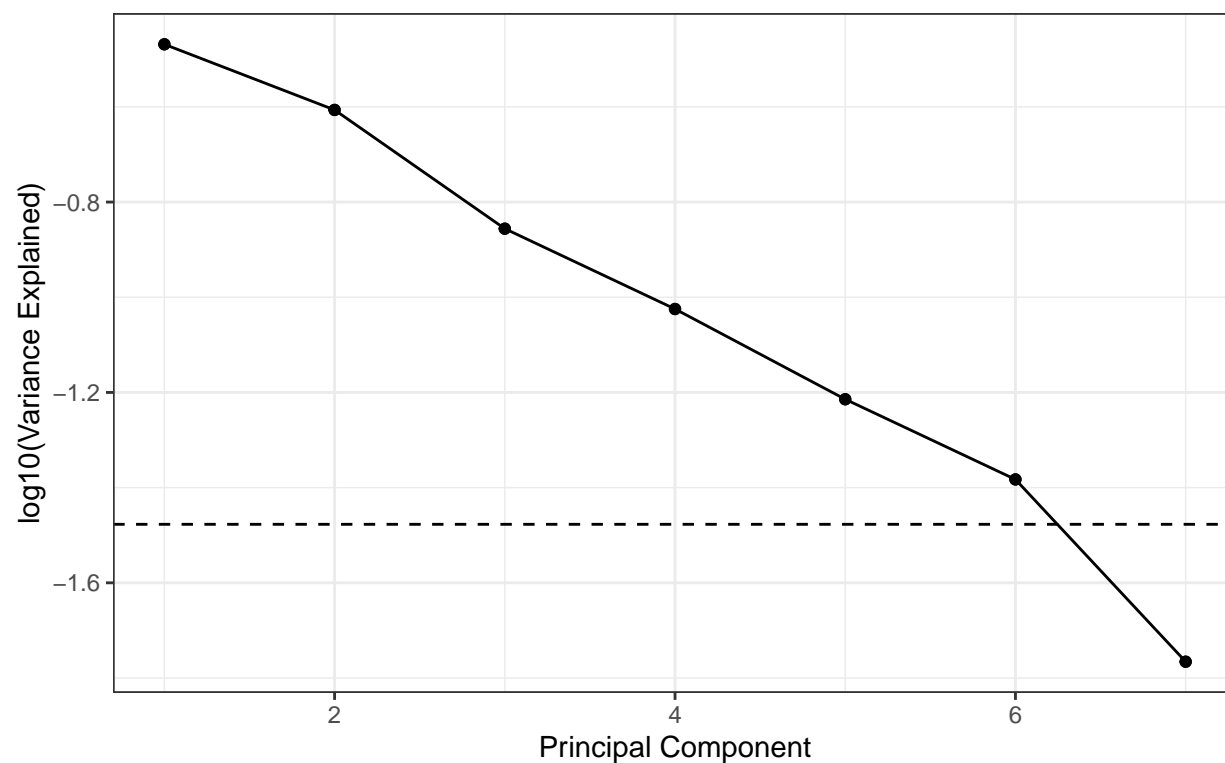


```
## Number of significant surrogate variables is: 3
## Iteration (out of 5 ):1 2 3 4 5 2
## [1] "SizeFactor distribution:"
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.2691 0.4948  1.4120  1.2100  1.8021  1.9976
```

SRP002628
Number of SVs: 4



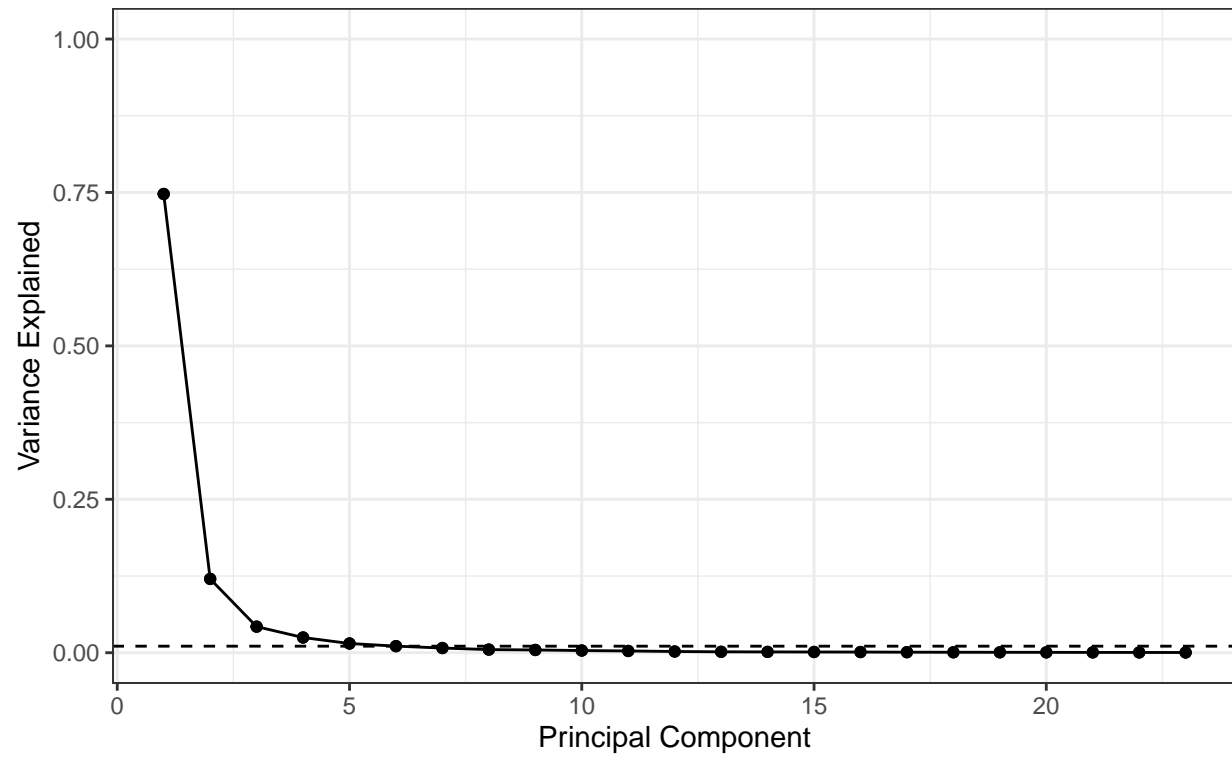
SRP002628
Number of SVs: 4



```
## Number of significant surrogate variables is: 4
## Iteration (out of 5 ):1 2 3 4 5
```

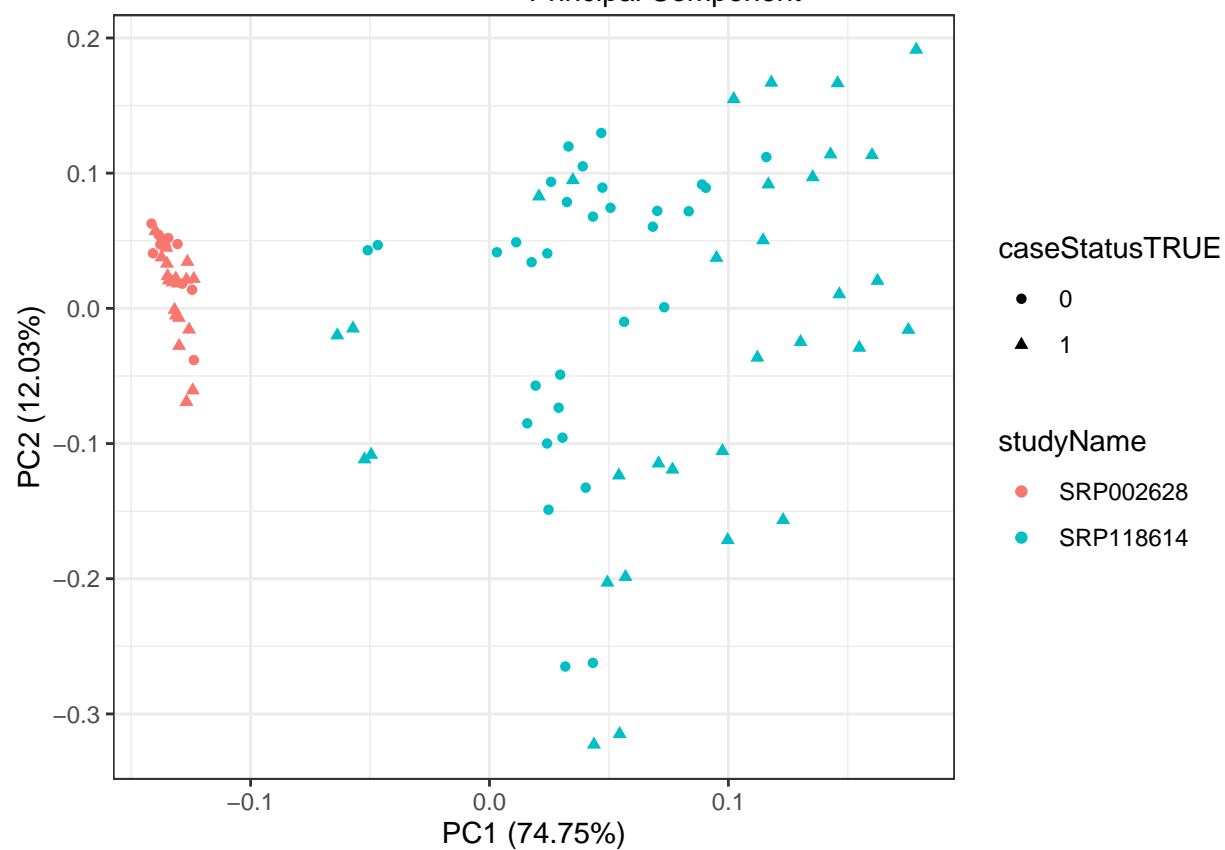
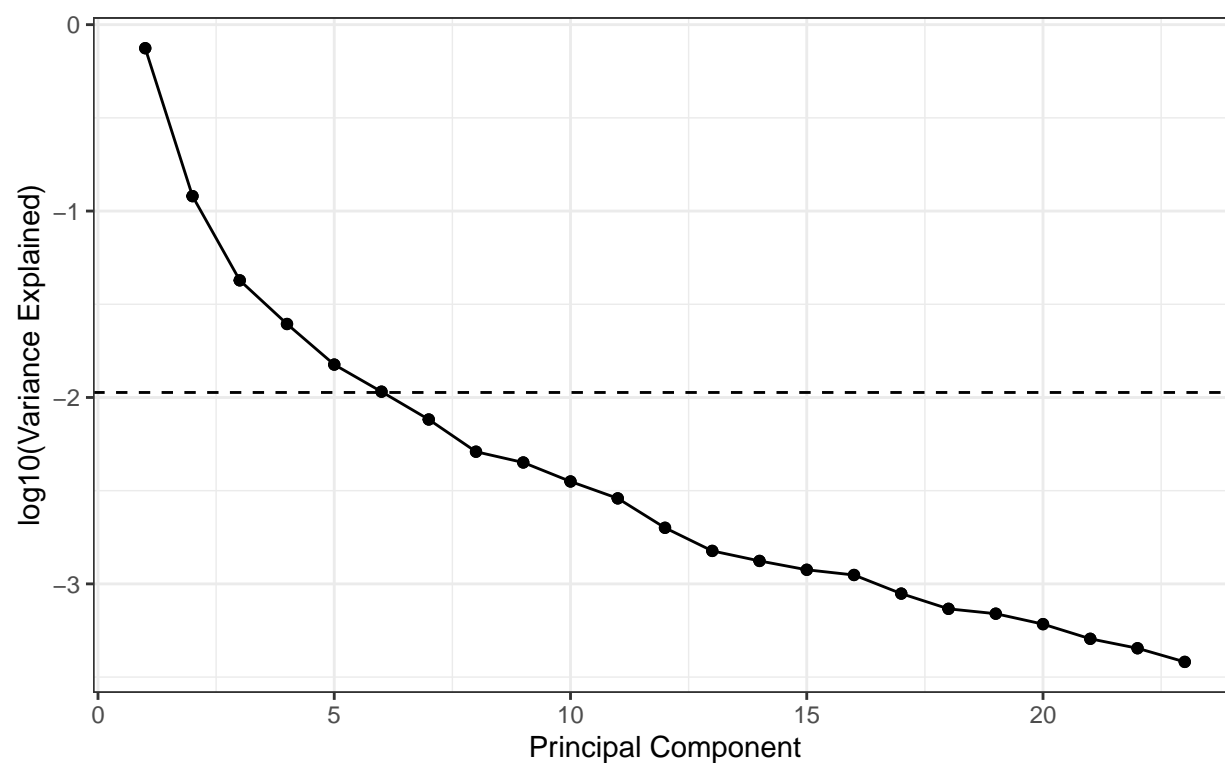
SRP118614 + SRP002628

Number of SVs: 2



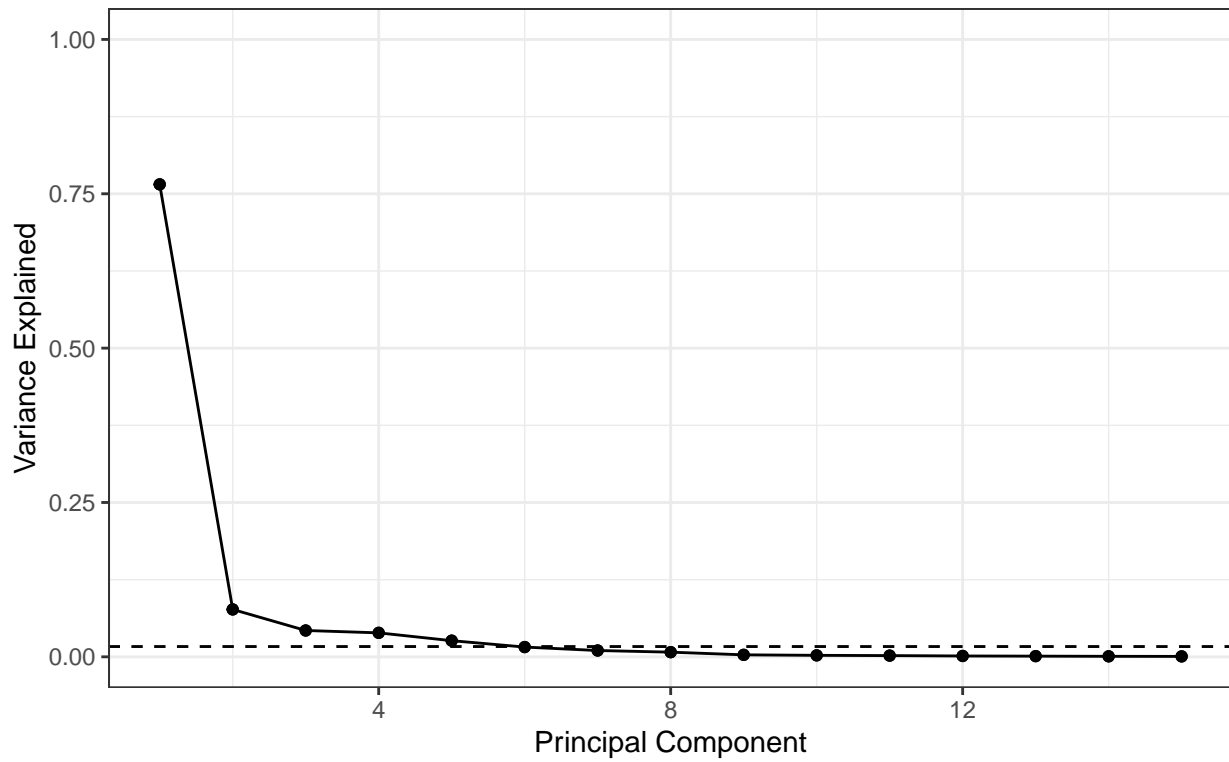
SRP118614 + SRP002628

Number of SVs: 2

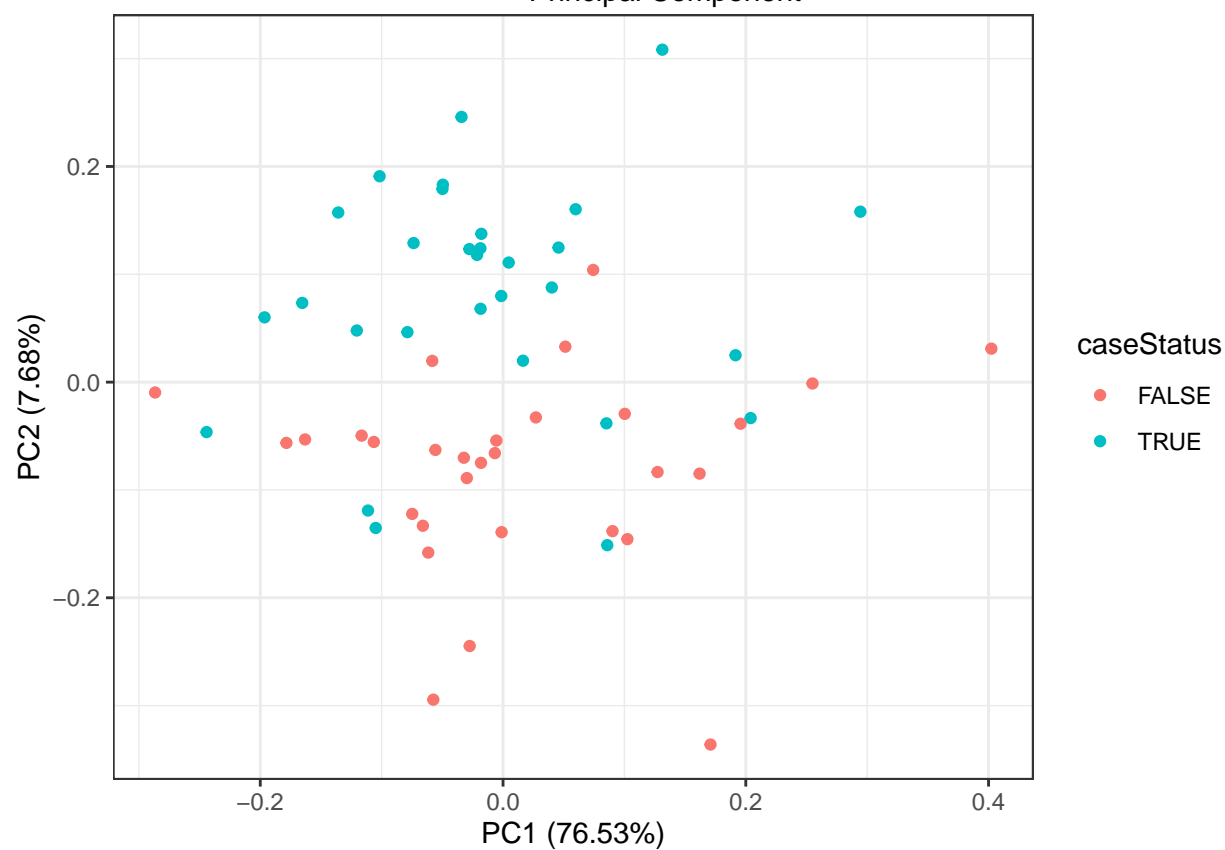
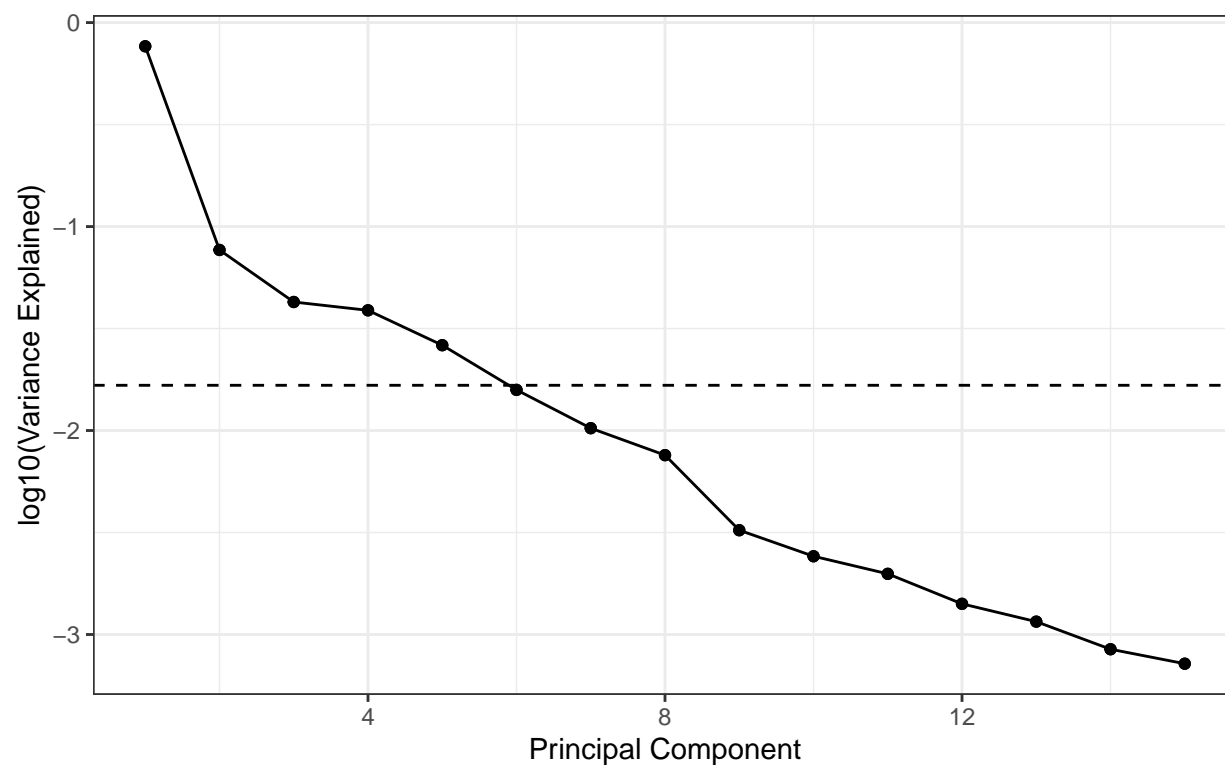


```
## Number of significant surrogate variables is: 2
## Iteration (out of 5 ):1 2 3 4 5 Number of significant surrogate variables is: 8
## Iteration (out of 5 ):1 2 3 4 5 3
## [1] "SizeFactor distribution:"
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.5444 0.8273  1.0219  1.0575  1.2505  1.8502
```

SRP212704
Number of SVs: 1



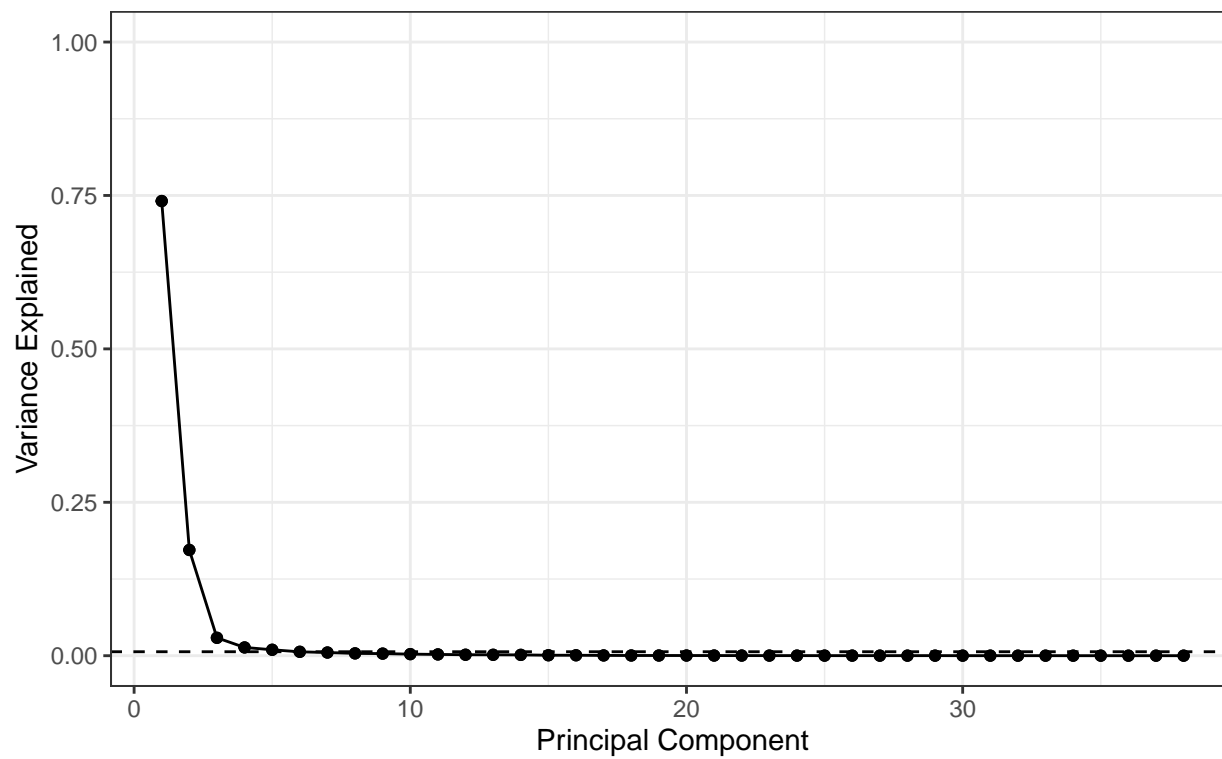
SRP212704
Number of SVs: 1



```
## Number of significant surrogate variables is: 1
## Iteration (out of 5 ):1 2 3 4 5
```

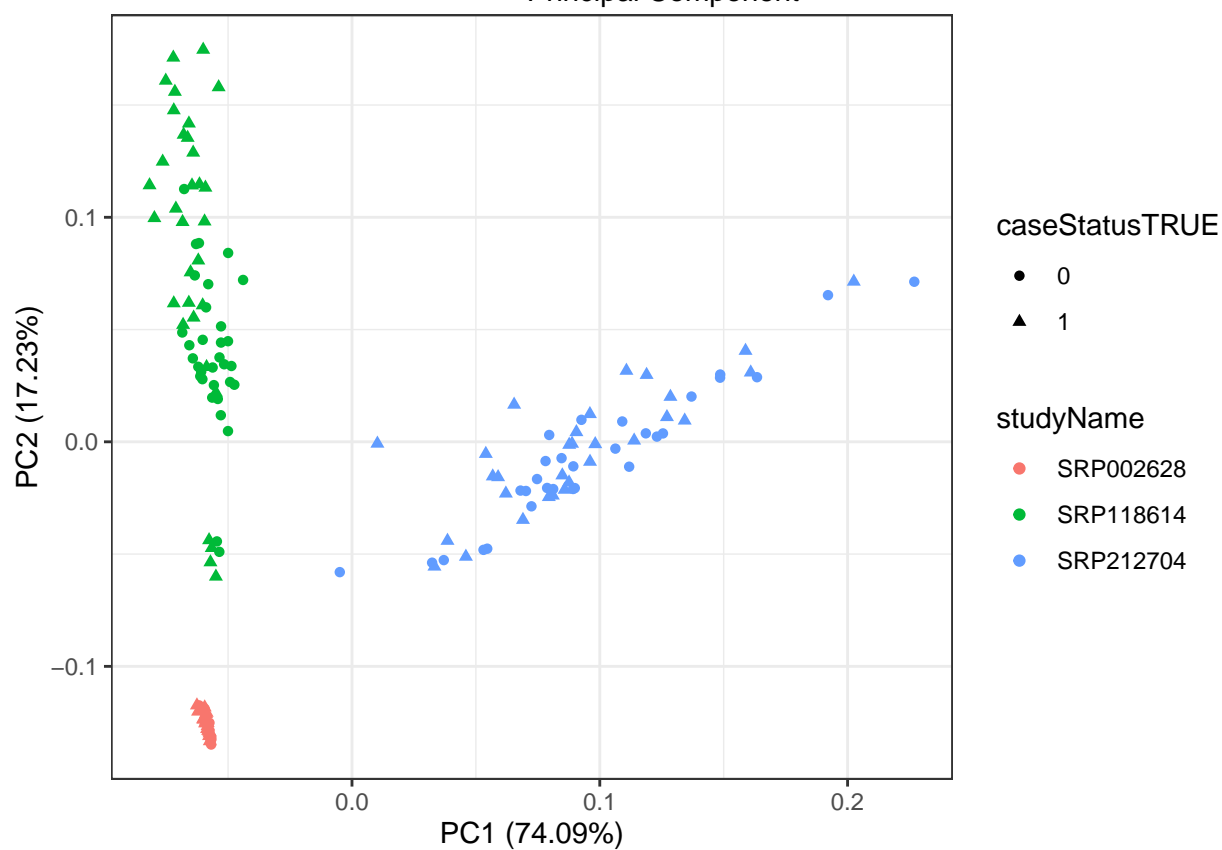
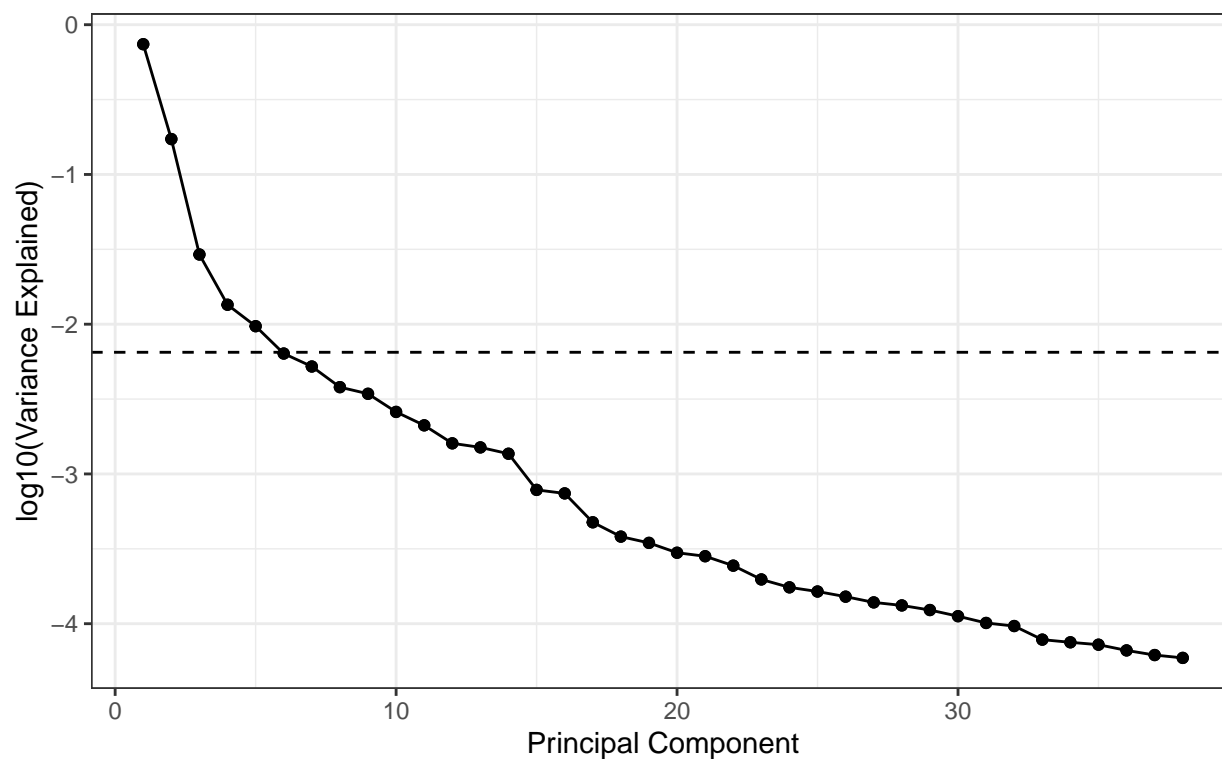
SRP118614 + SRP002628 + SRP212704

Number of SVs: 1



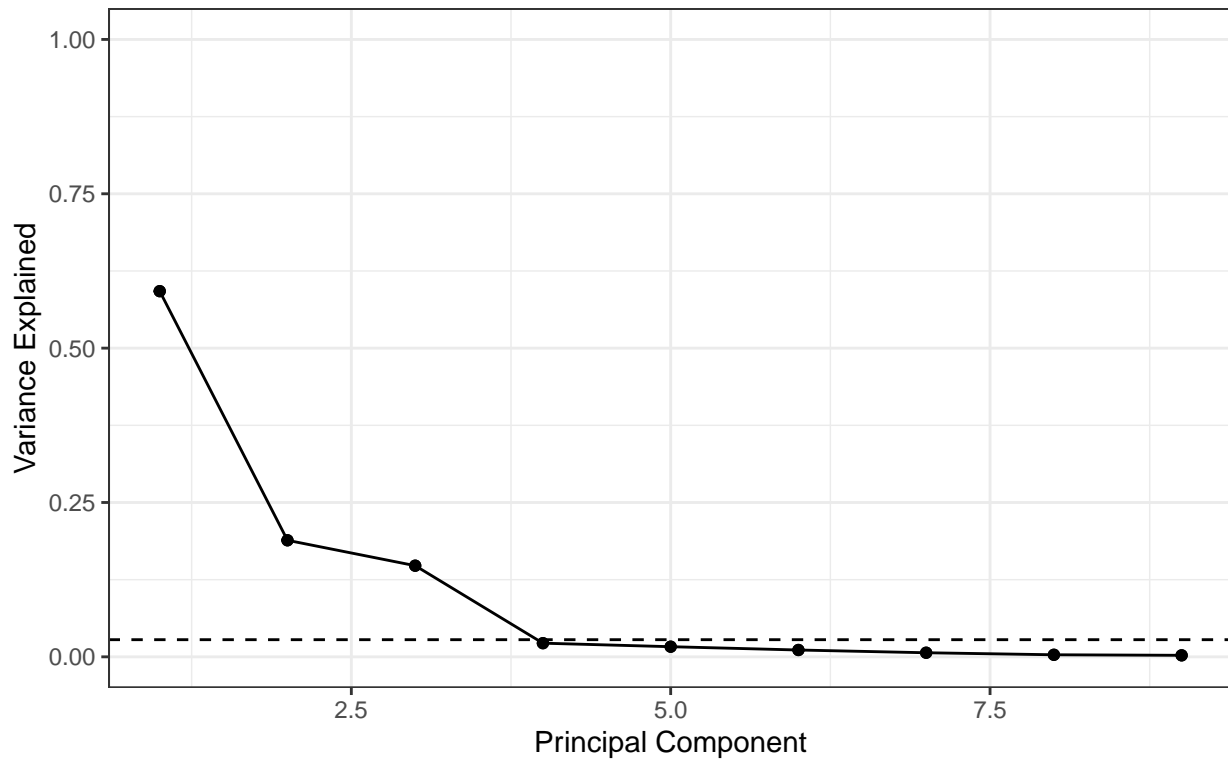
SRP118614 + SRP002628 + SRP212704

Number of SVs: 1

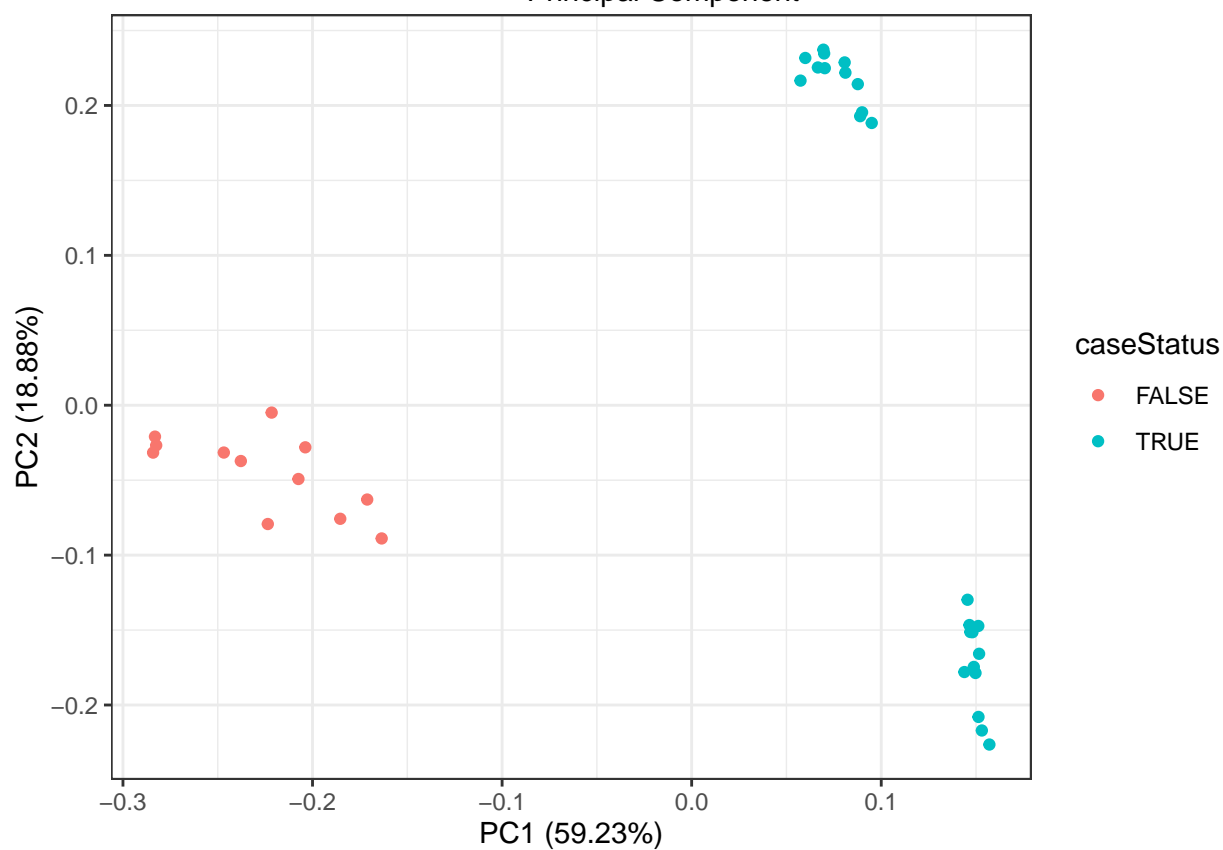
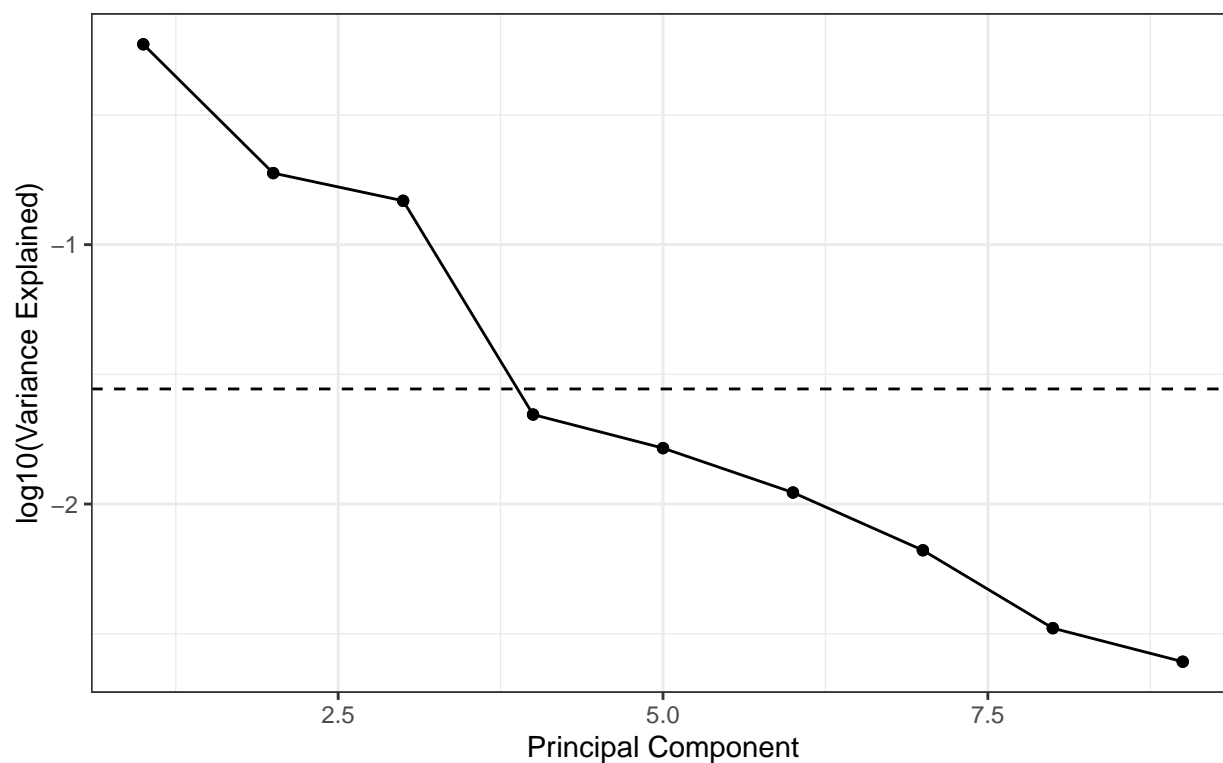


```
## Number of significant surrogate variables is: 1
## Iteration (out of 5 ):1 2 3 4 5 Number of significant surrogate variables is: 10
## Iteration (out of 5 ):1 2 3 4 5 4
## [1] "SizeFactor distribution:"
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.7349  0.9231  1.0293  1.0309  1.1130  1.5355
```

SRP027258
Number of SVs: 2



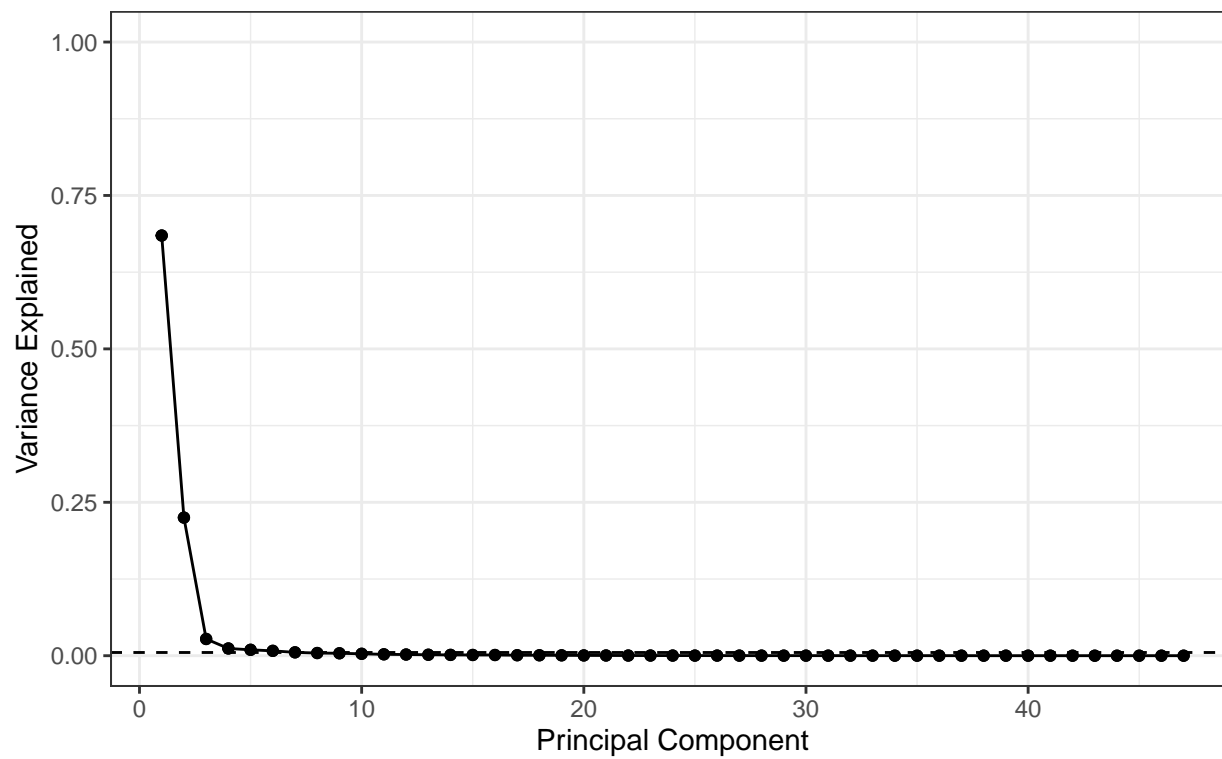
SRP027258
Number of SVs: 2



```
## Number of significant surrogate variables is: 2
## Iteration (out of 5 ):1 2 3 4 5
```

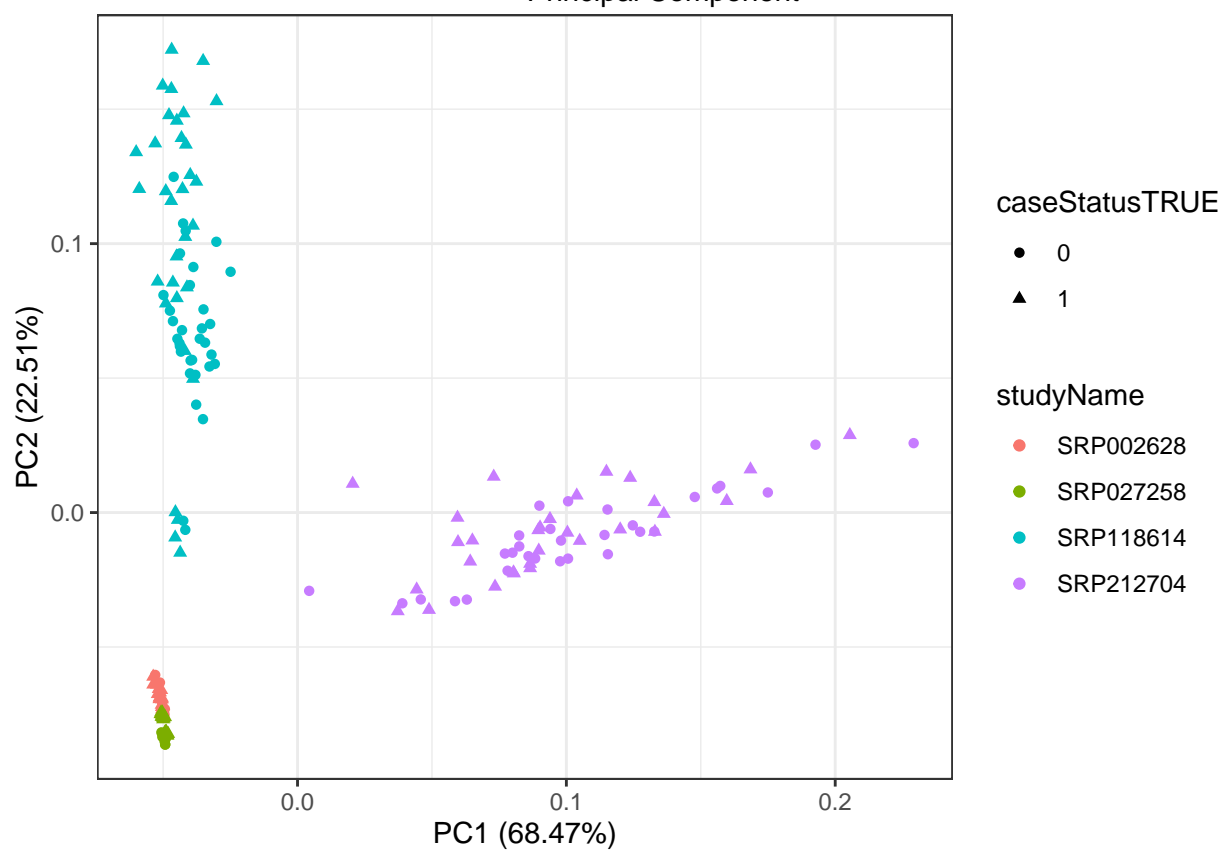
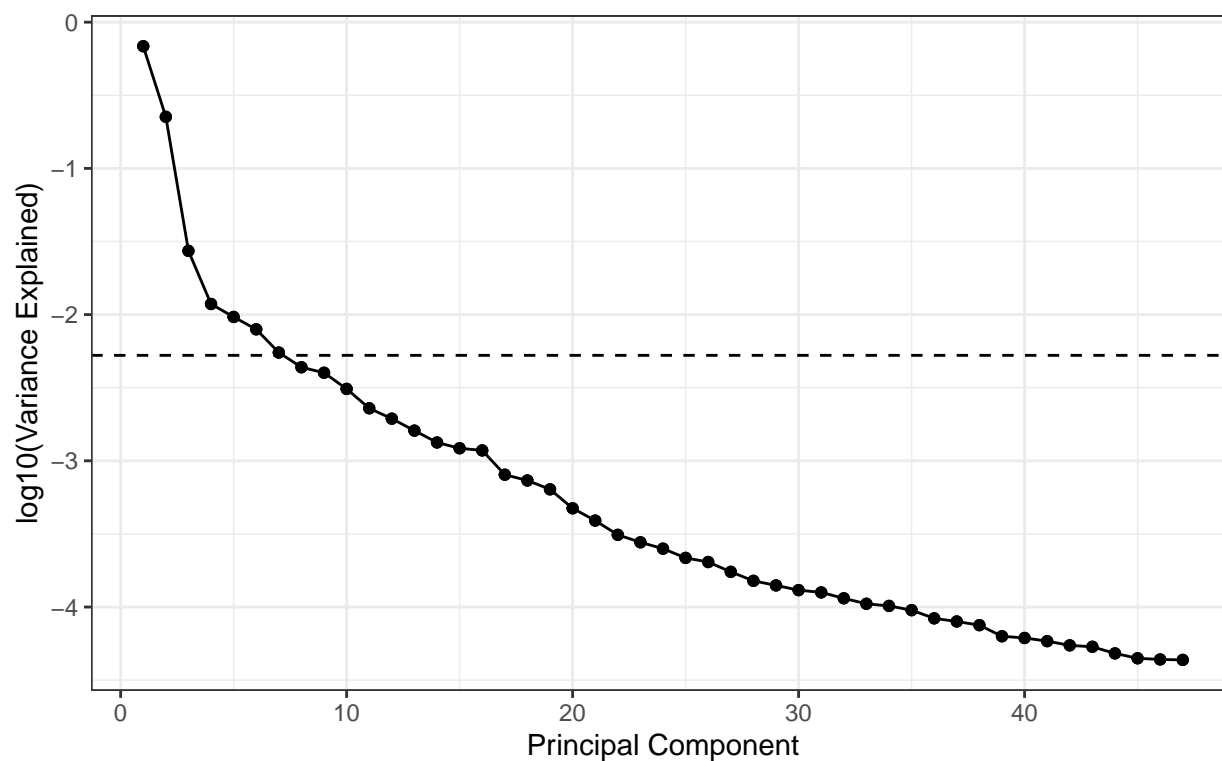
SRP118614 + SRP002628 + SRP212704 + SRP027258

Number of SVs: 2

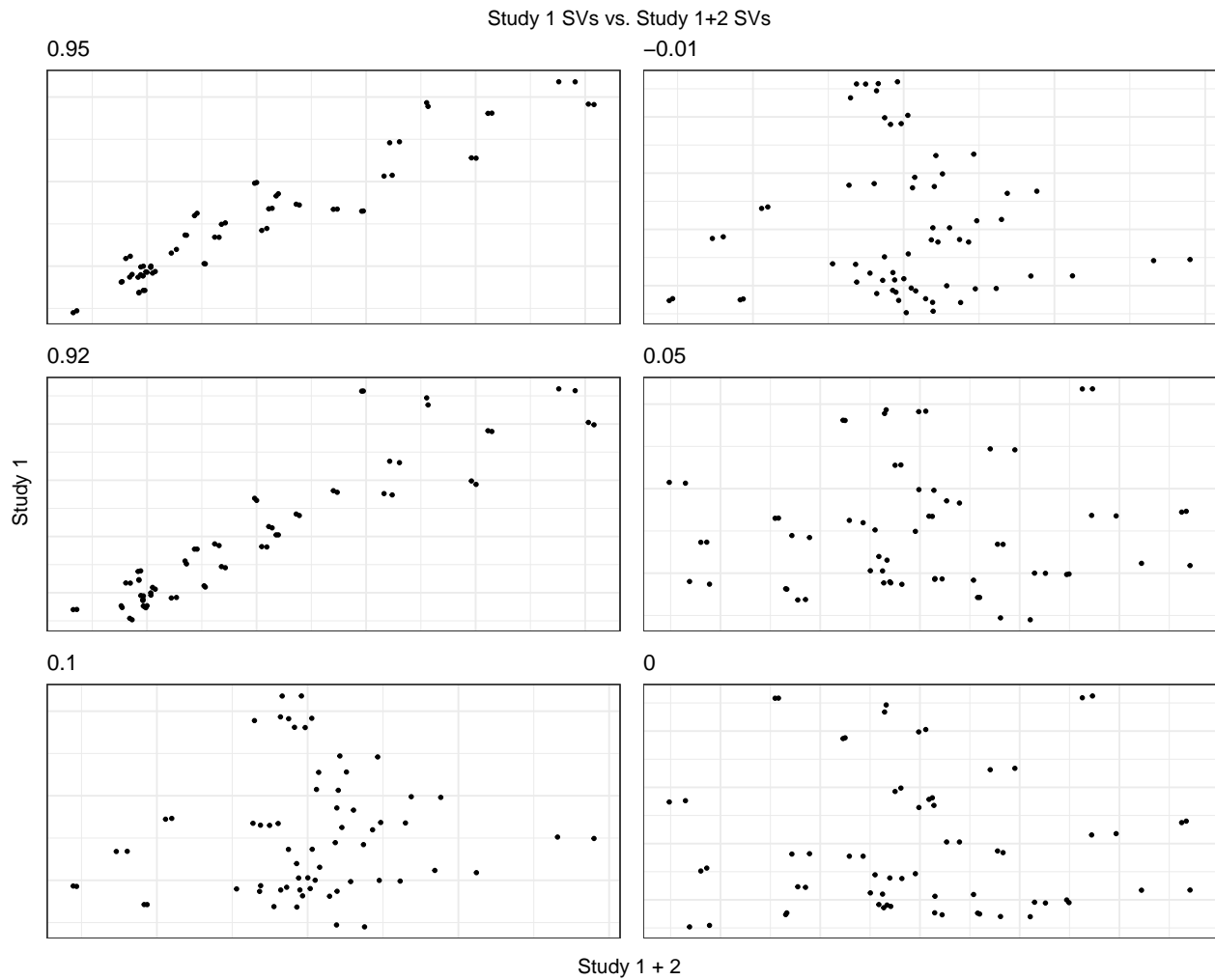


SRP118614 + SRP002628 + SRP212704 + SRP027258

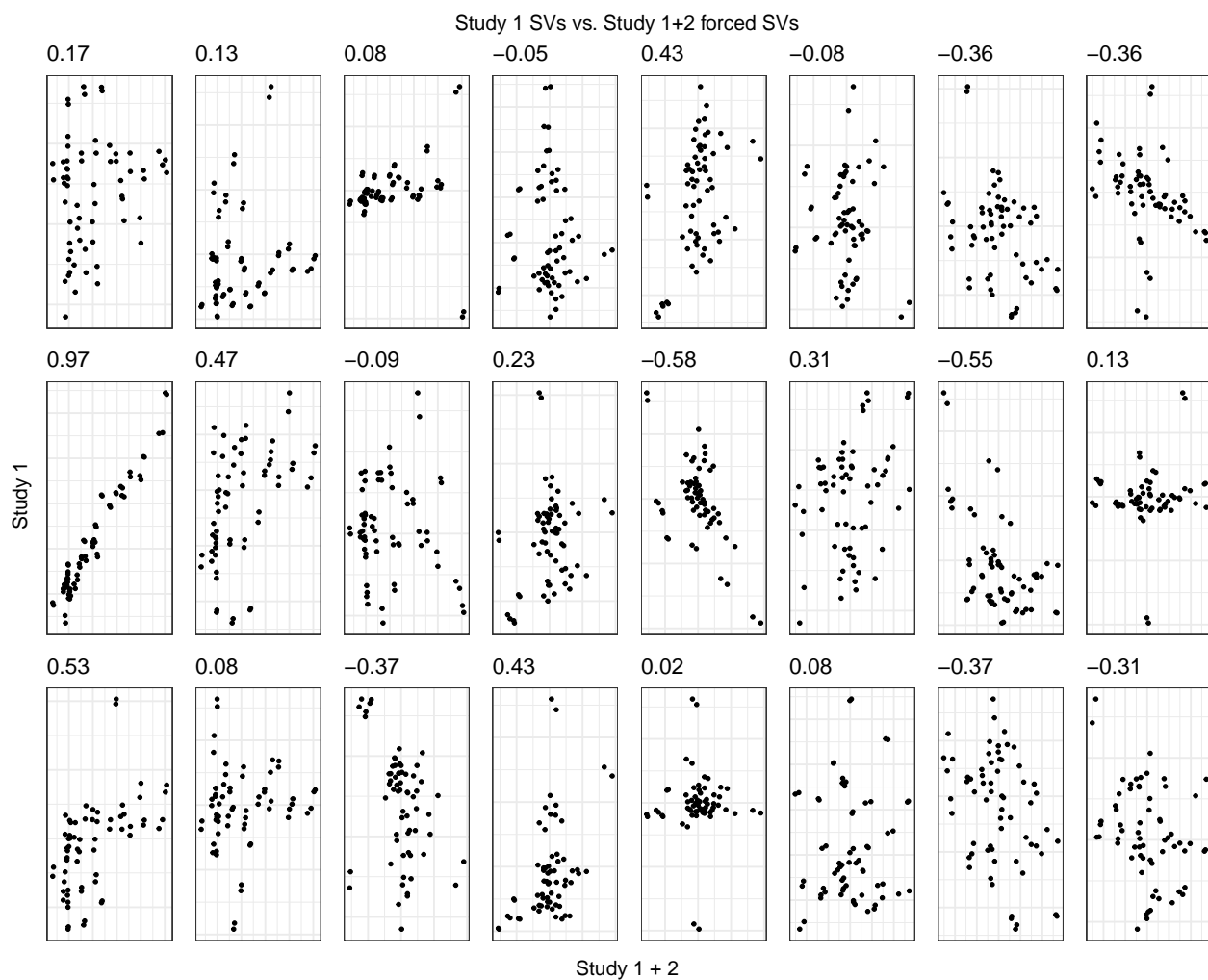
Number of SVs: 2



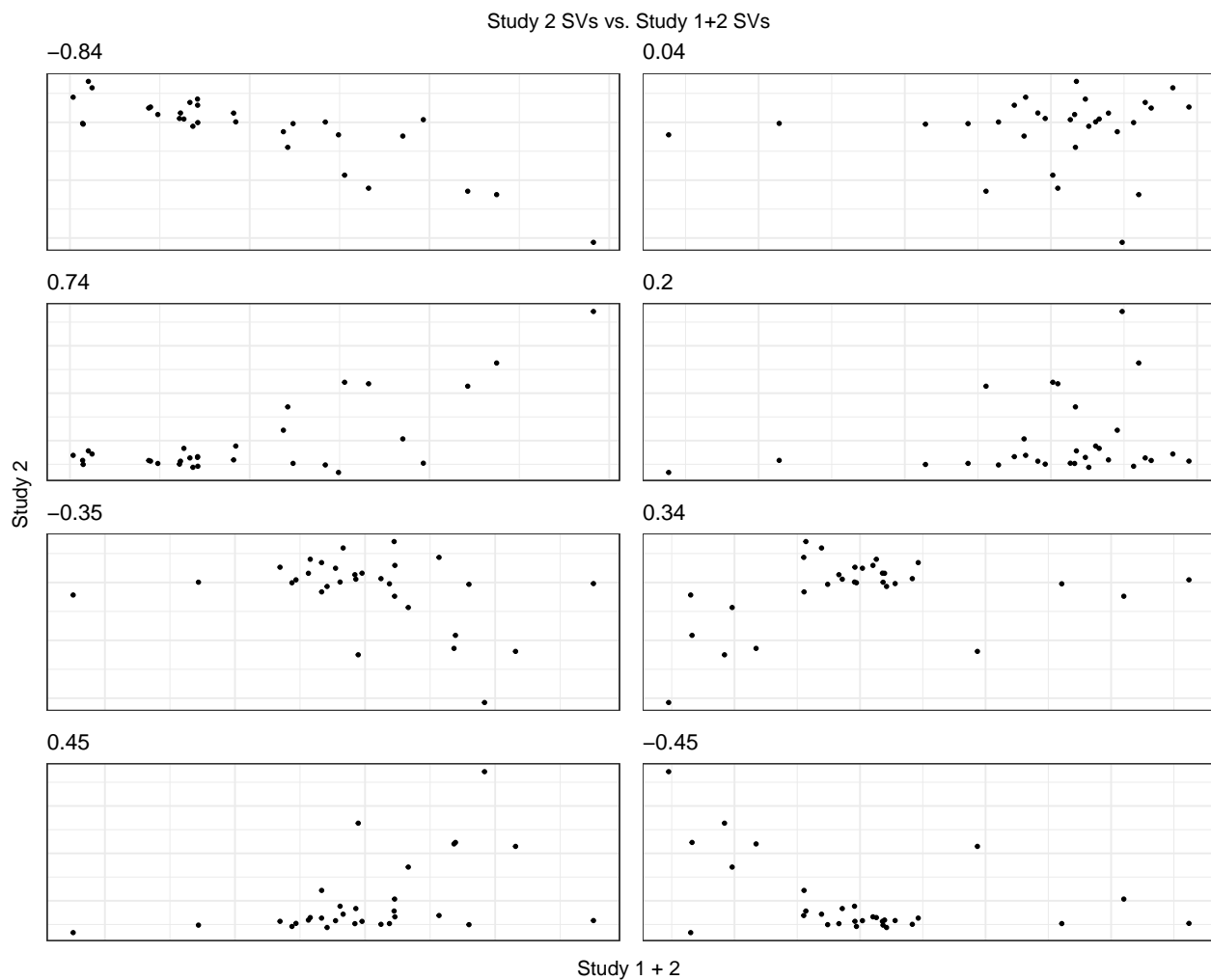
```
## Number of significant surrogate variables is: 2
## Iteration (out of 5):1 2 3 4 5 Number of significant surrogate variables is: 13
## Iteration (out of 5):1 2 3 4 5
```



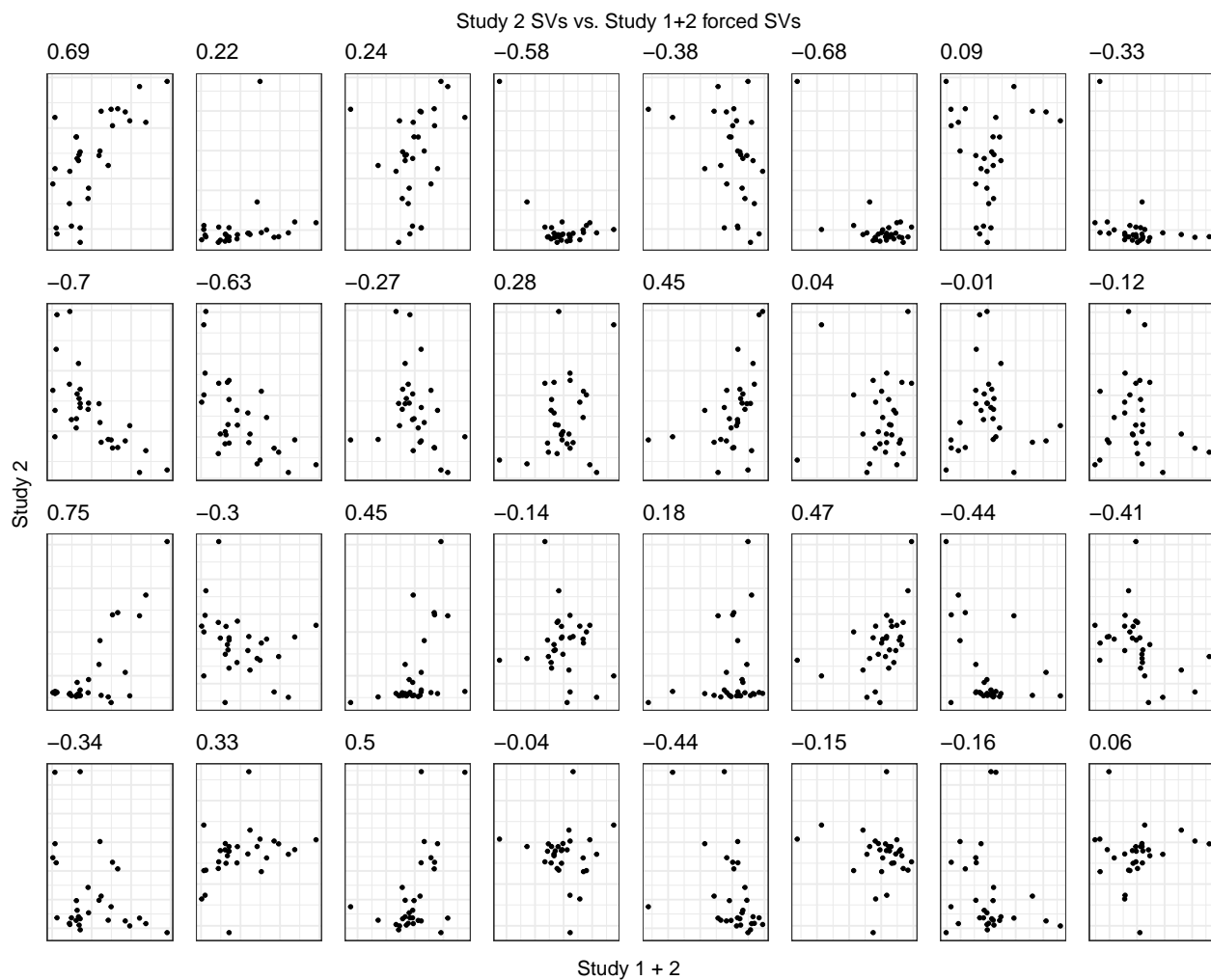
```
## [[1]]
## NULL
```

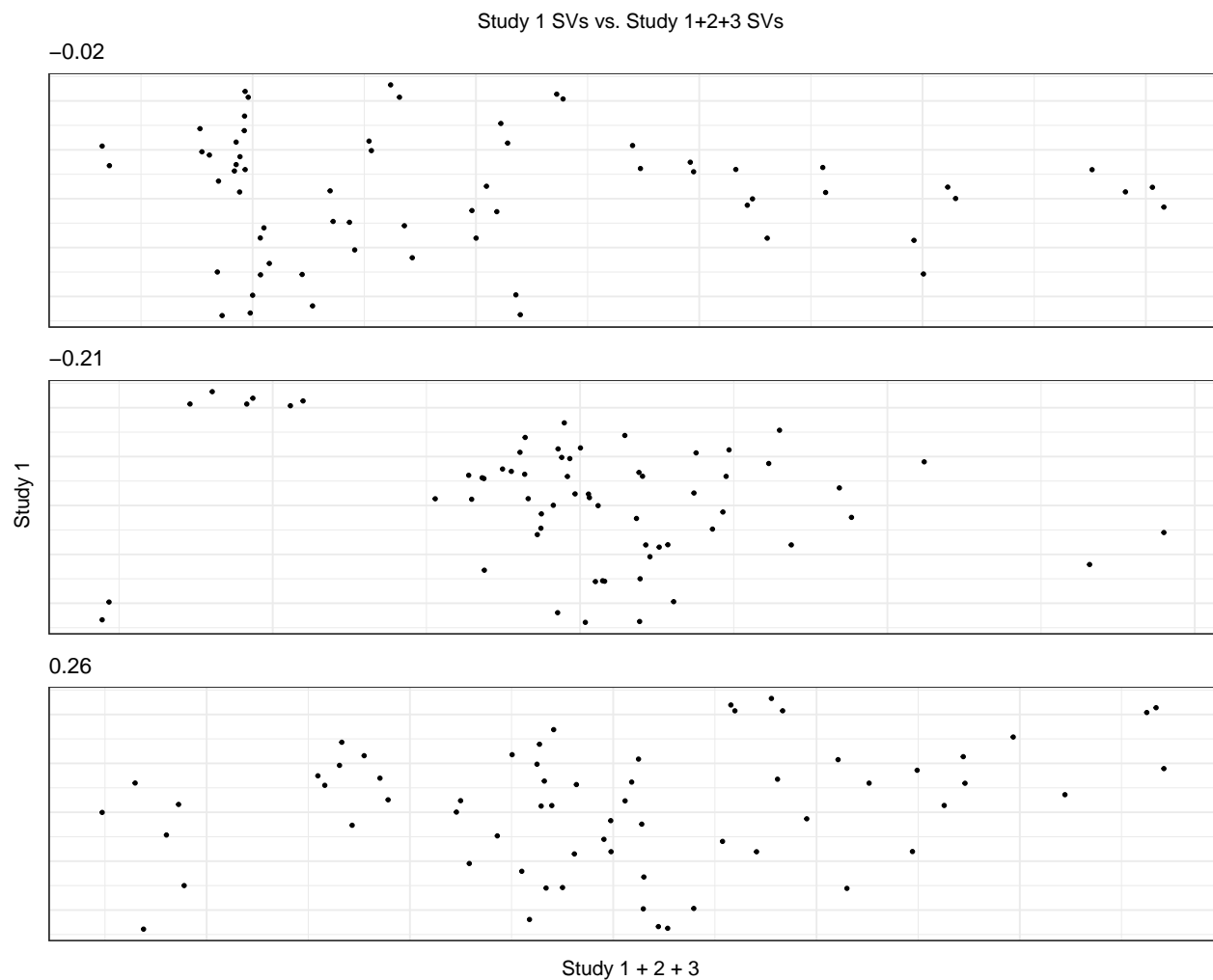
```
## [[1]]
## NULL
```



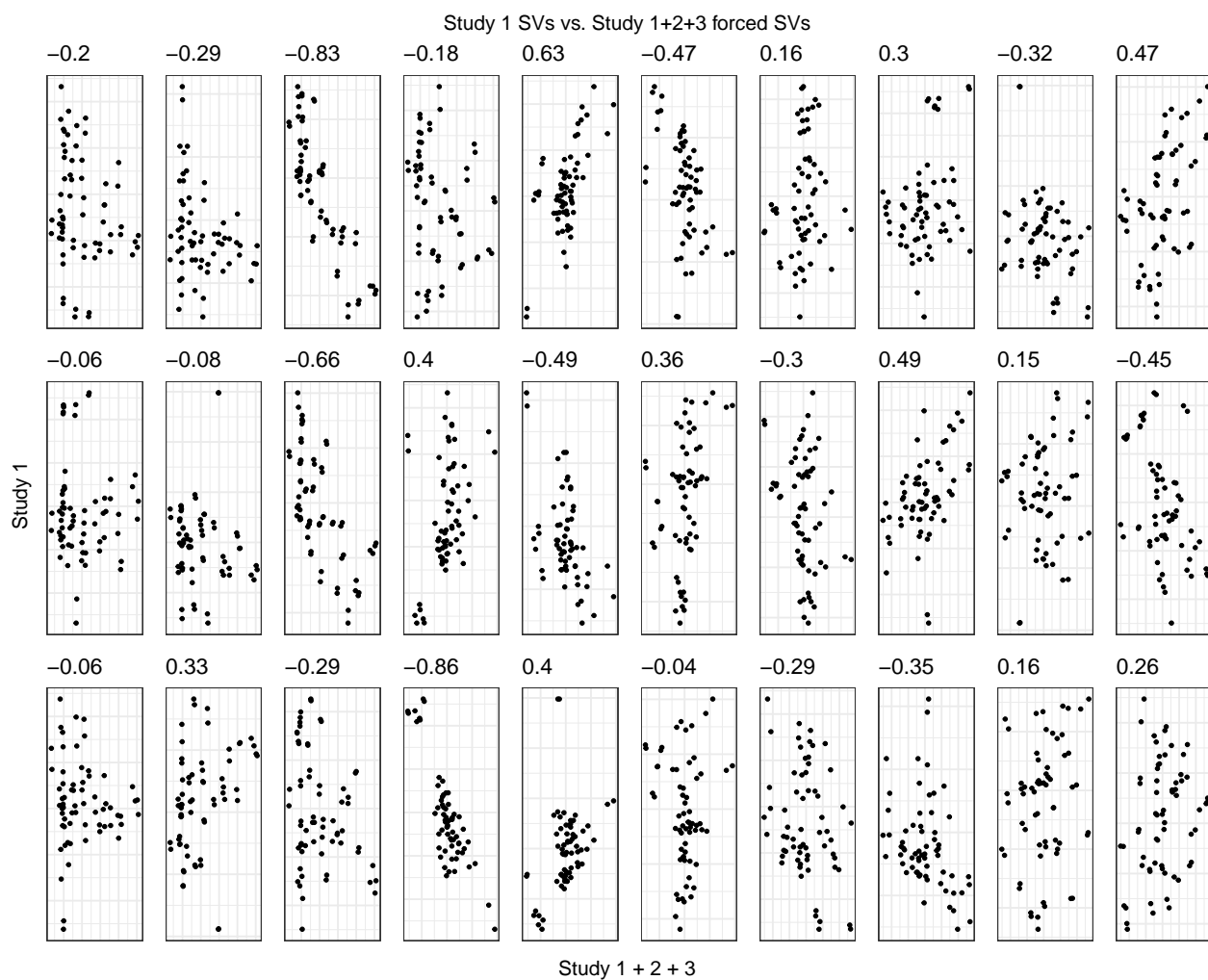
```
## [[1]]
## NULL
```



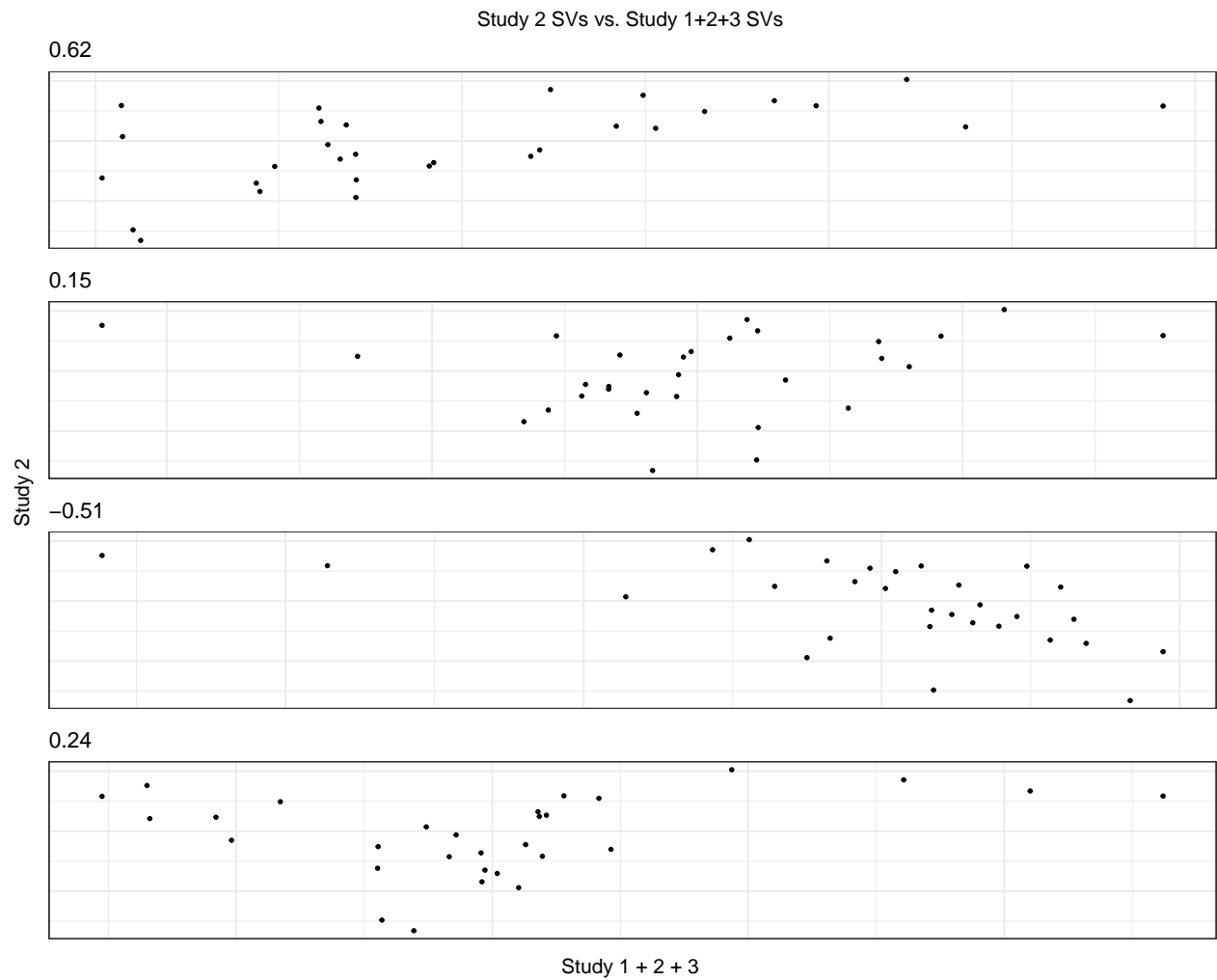
```
## [[1]]
## NULL
```



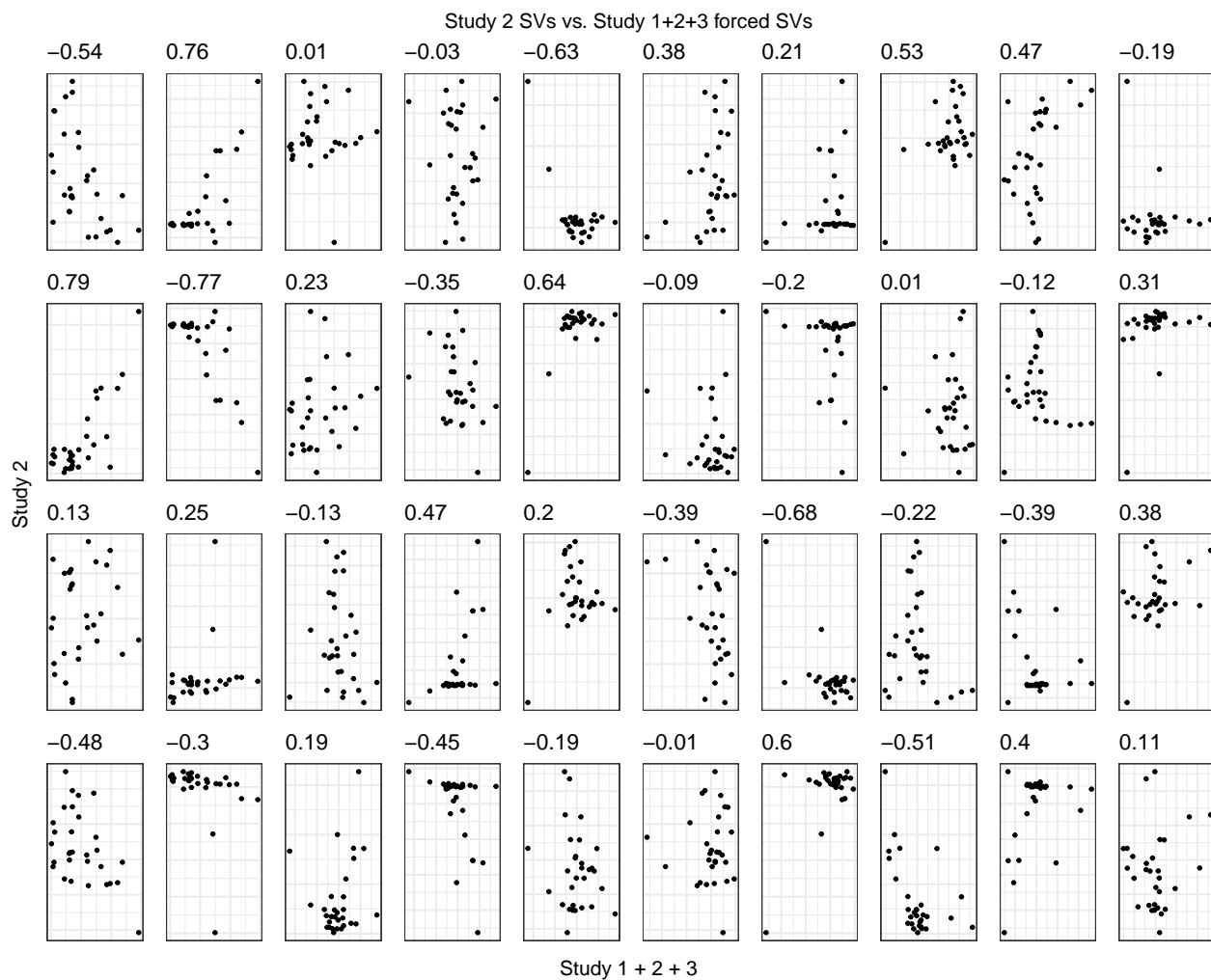
```
## [[1]]
## NULL
```



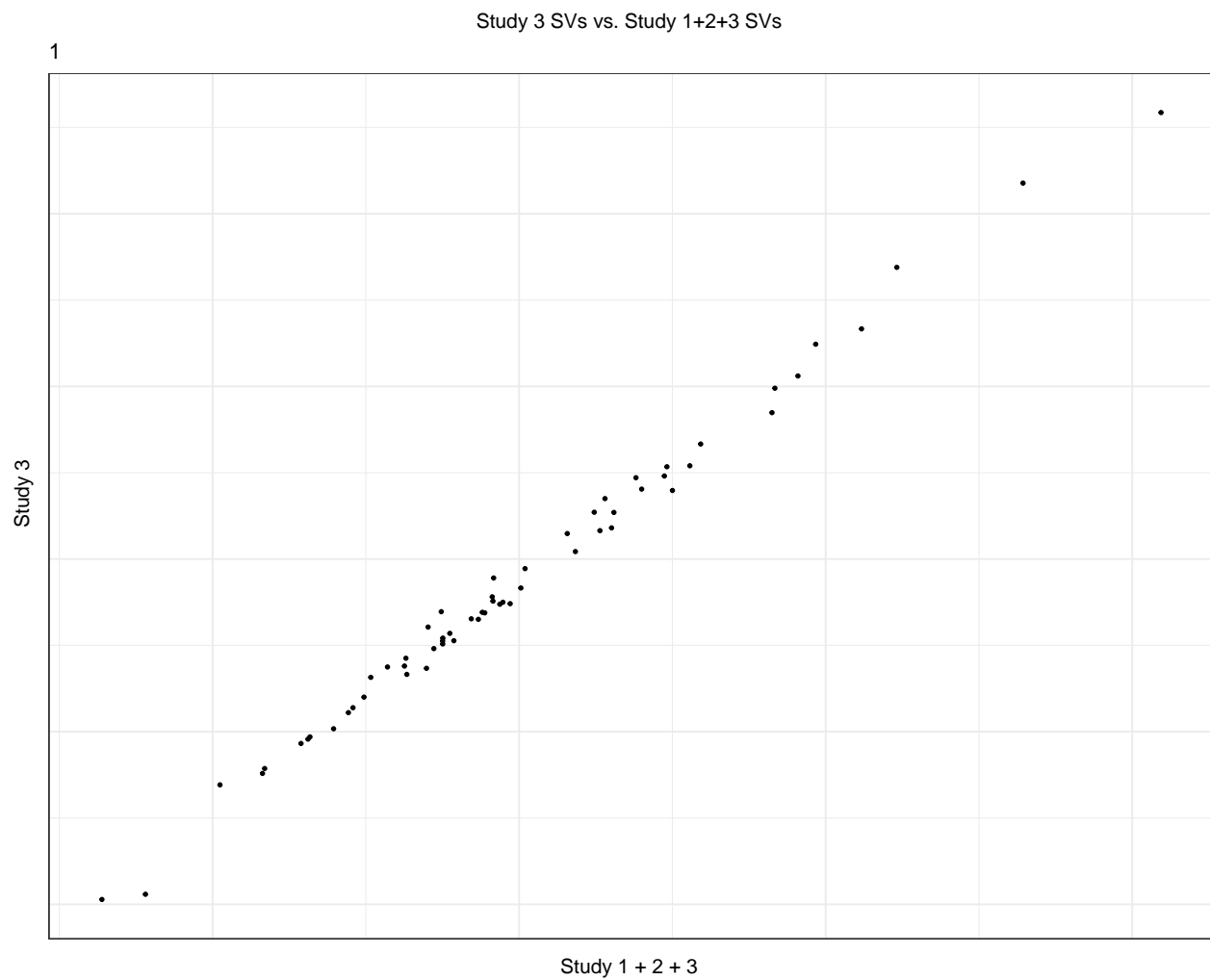
```
## [[1]]
## NULL
```



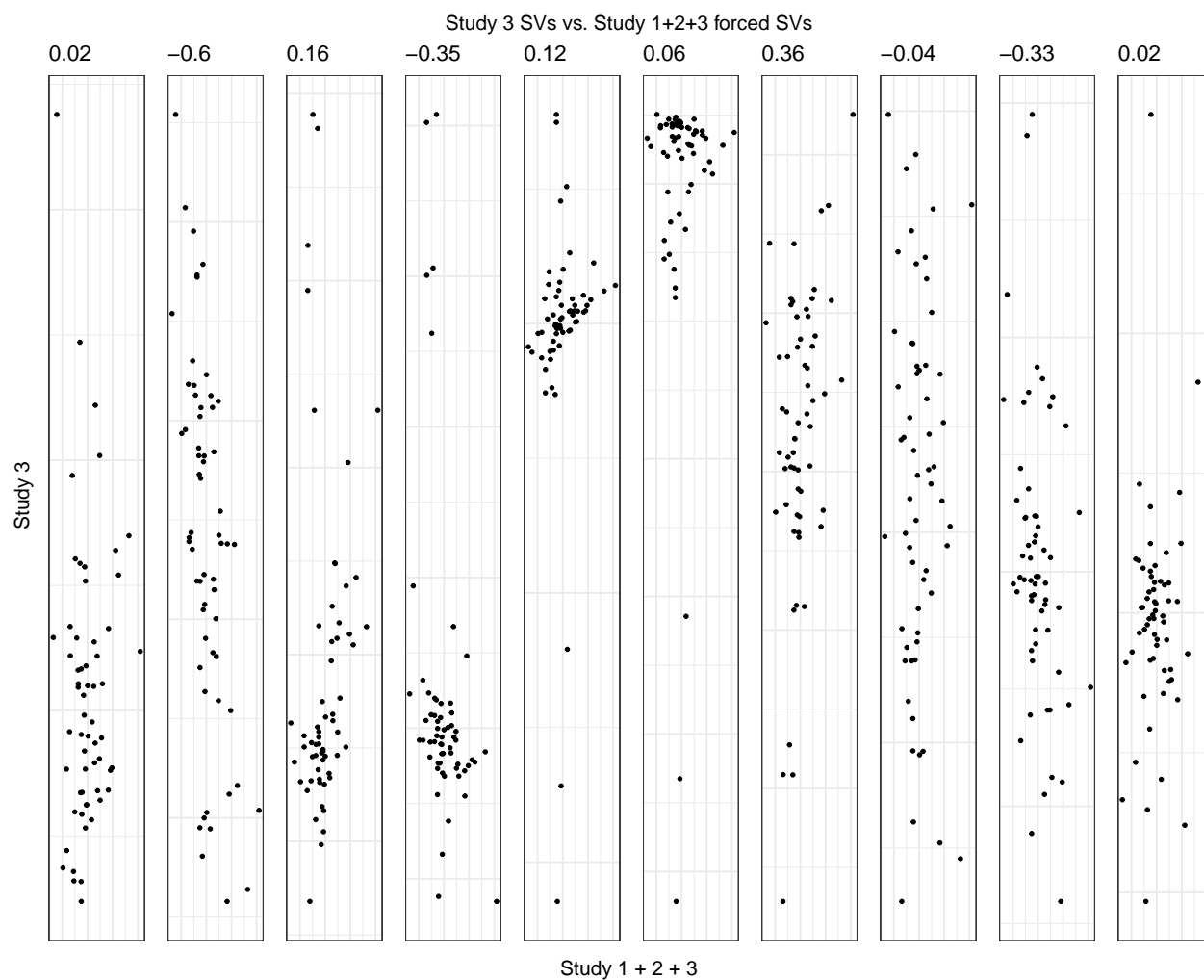
```
## [[1]]
## NULL
```



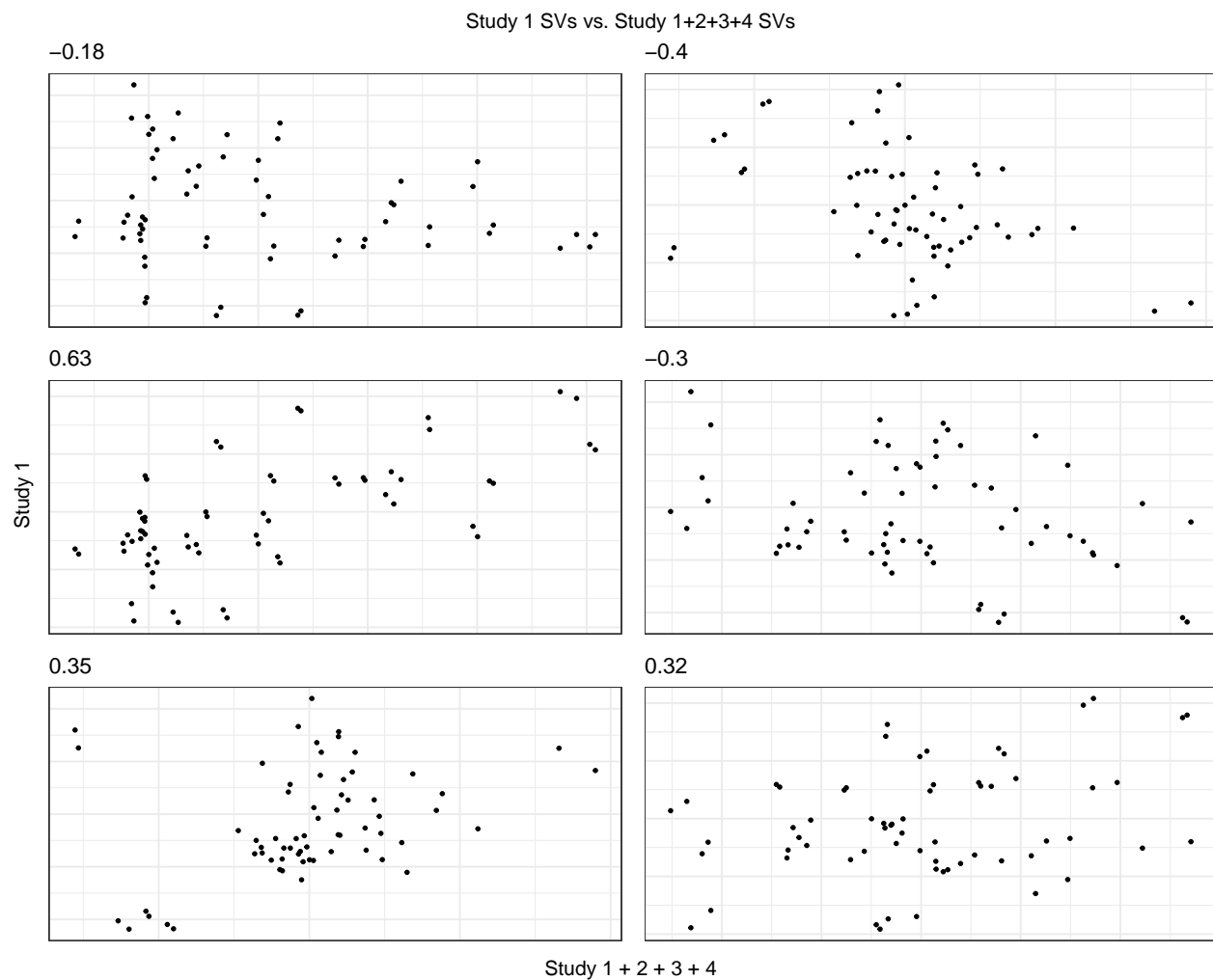
```
## [[1]]
## NULL
```



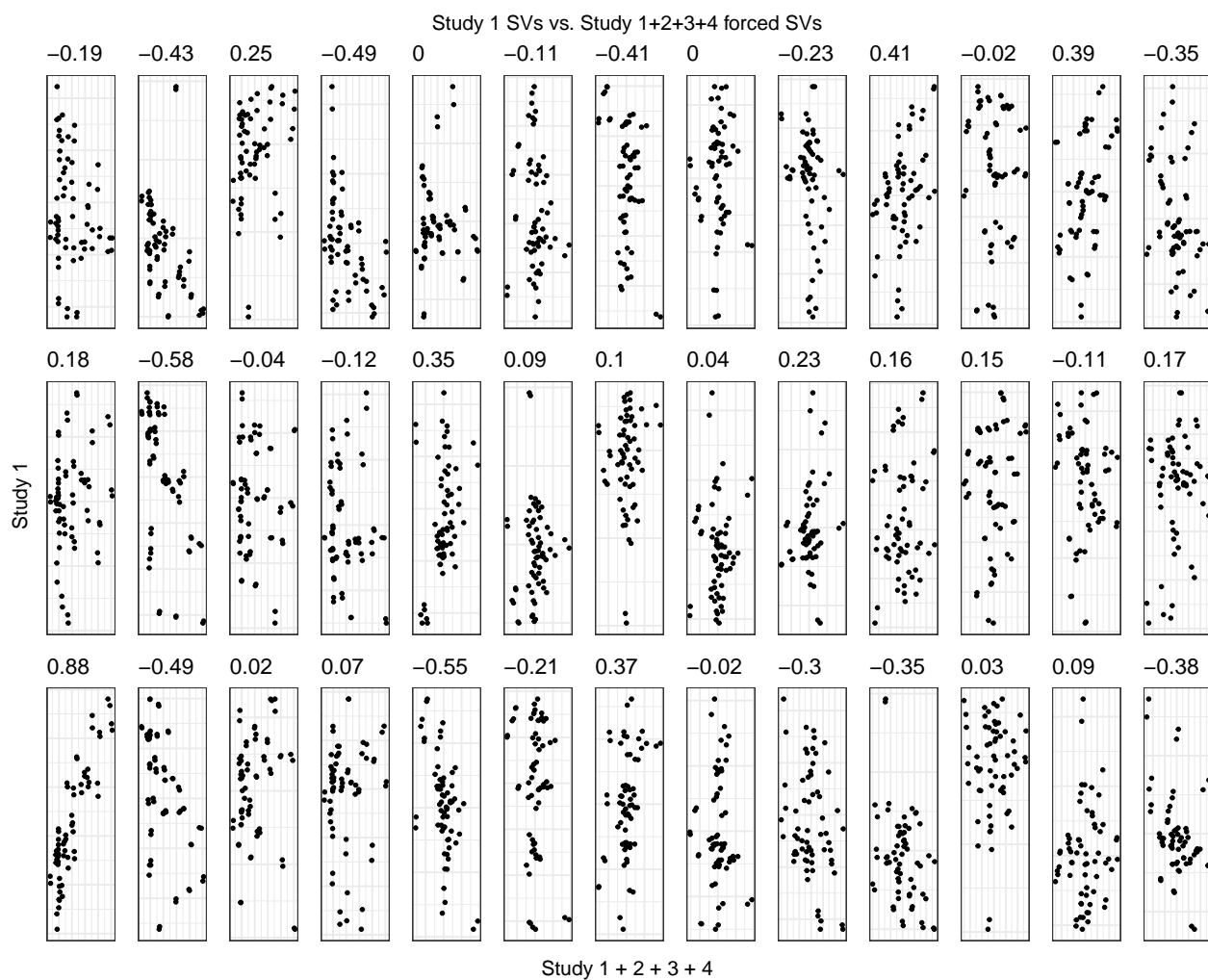
```
## [[1]]  
## NULL
```

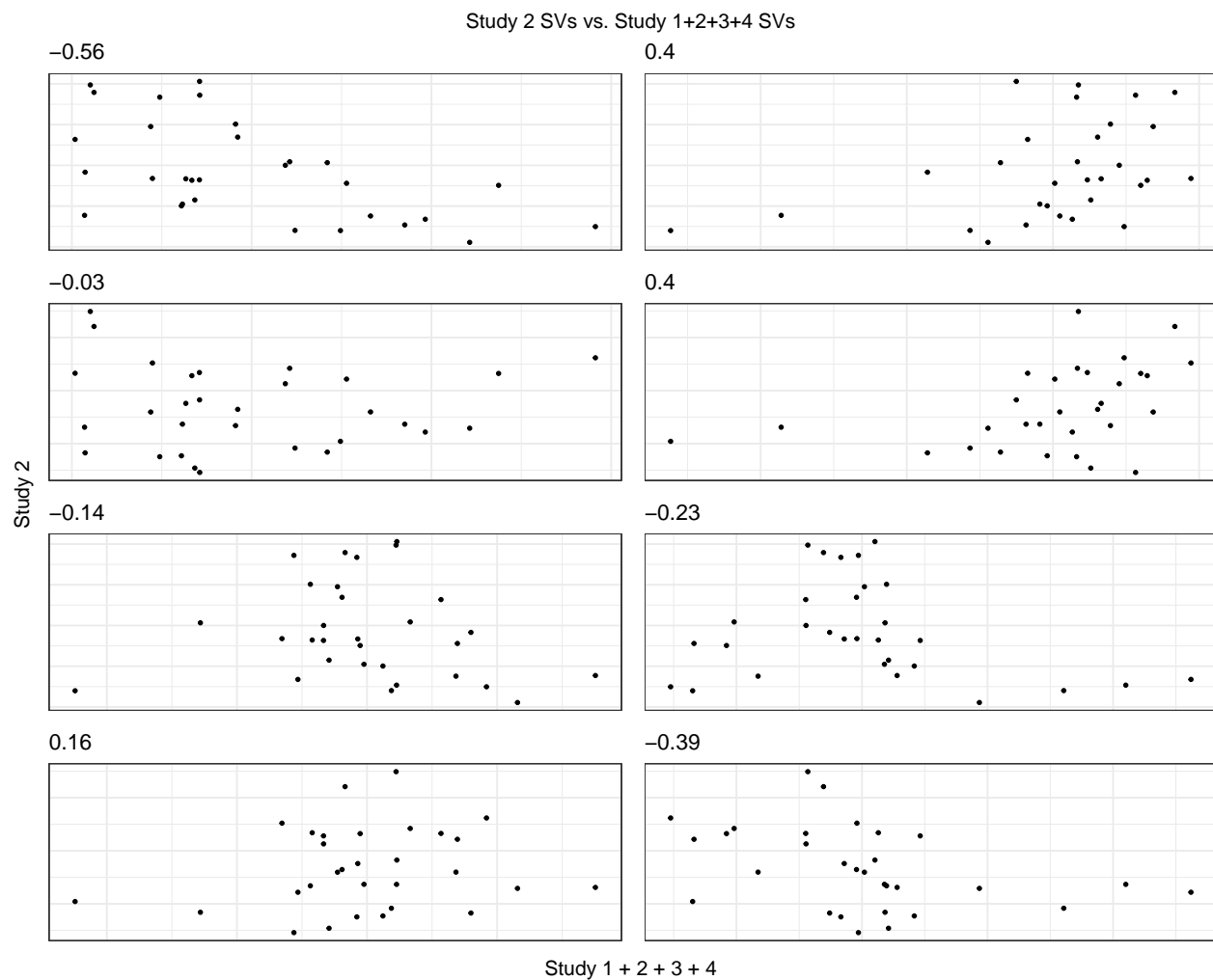
```
## [[1]]
## NULL
```



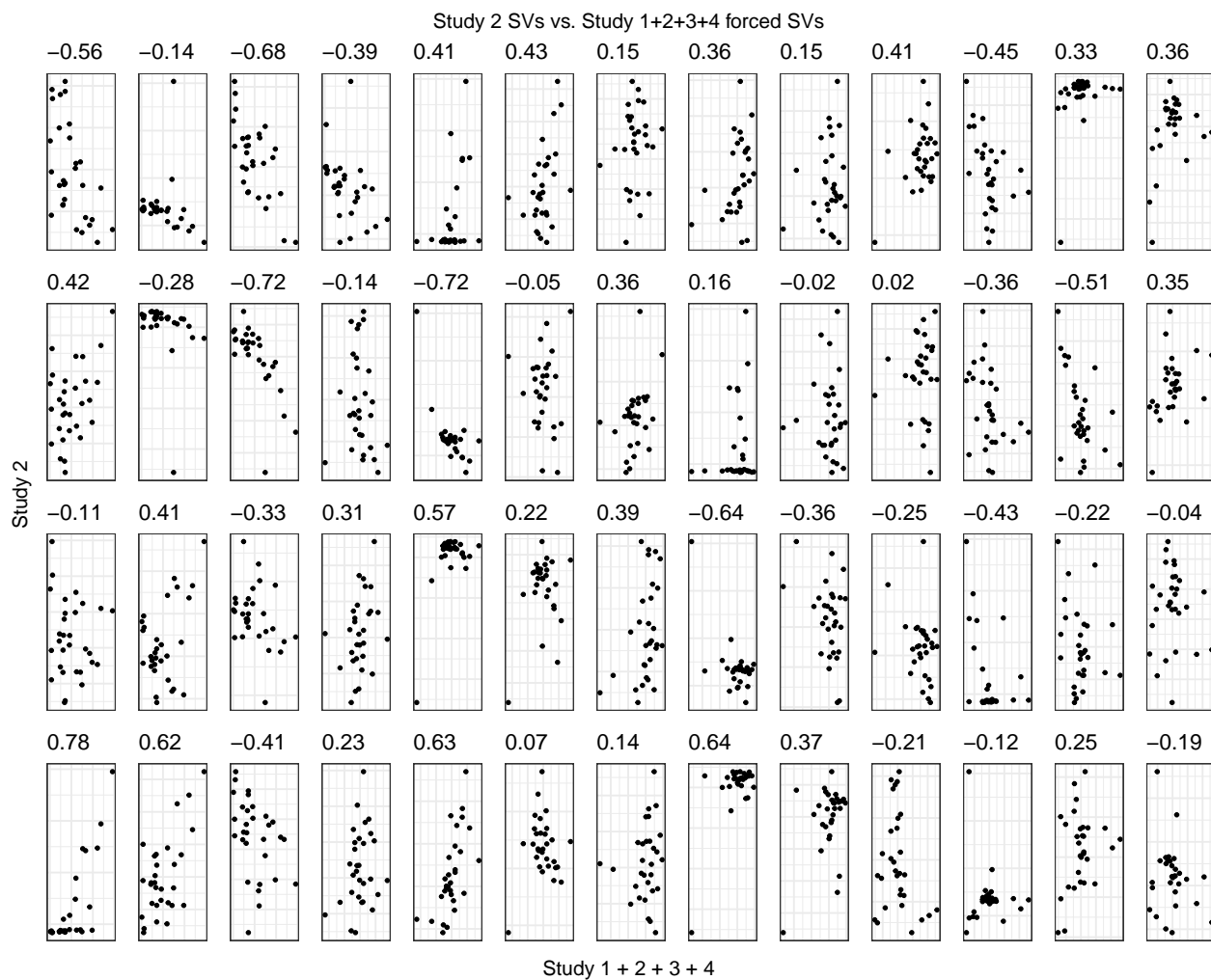
```
## [[1]]
## NULL
```



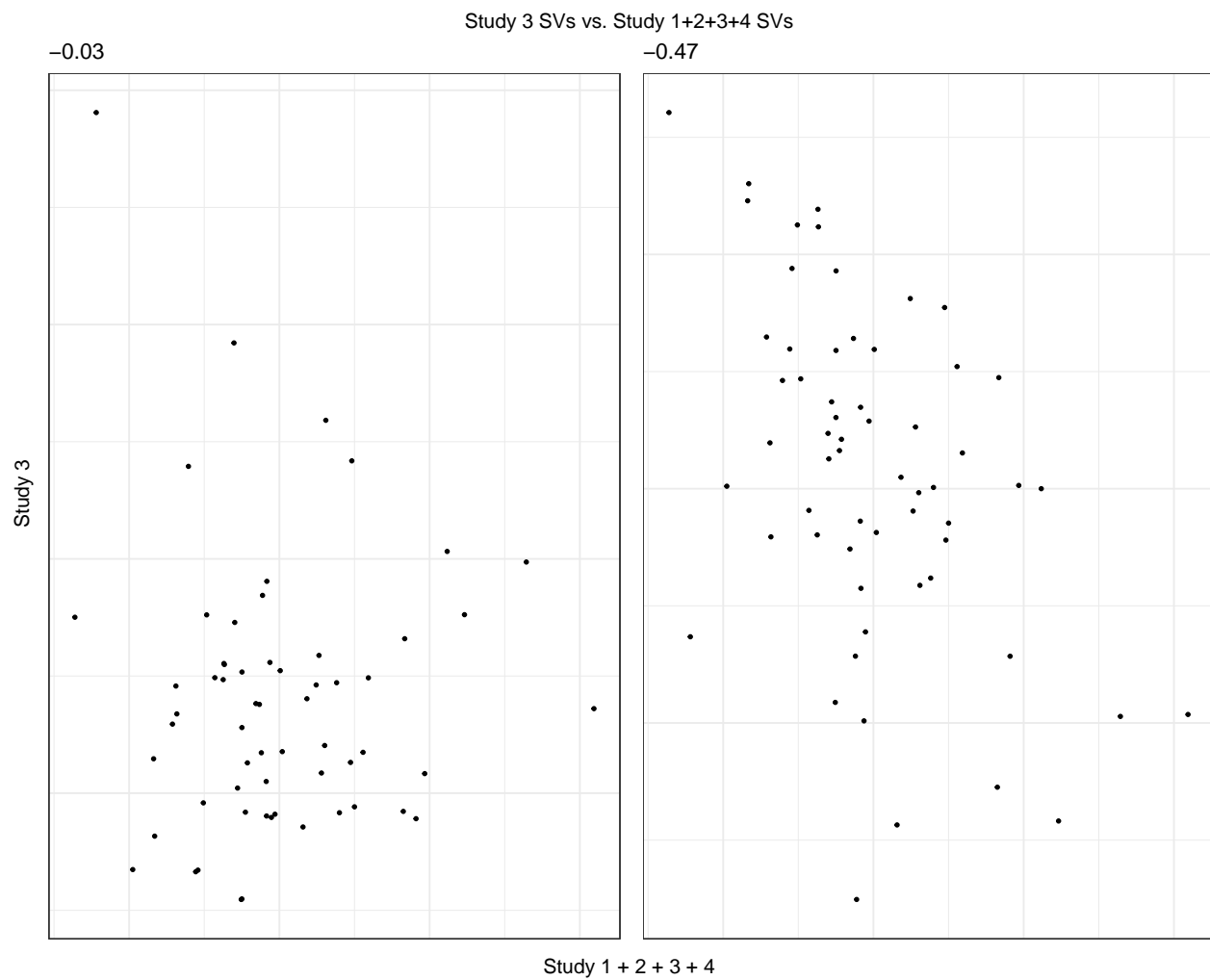
```
## [[1]]
## NULL
```



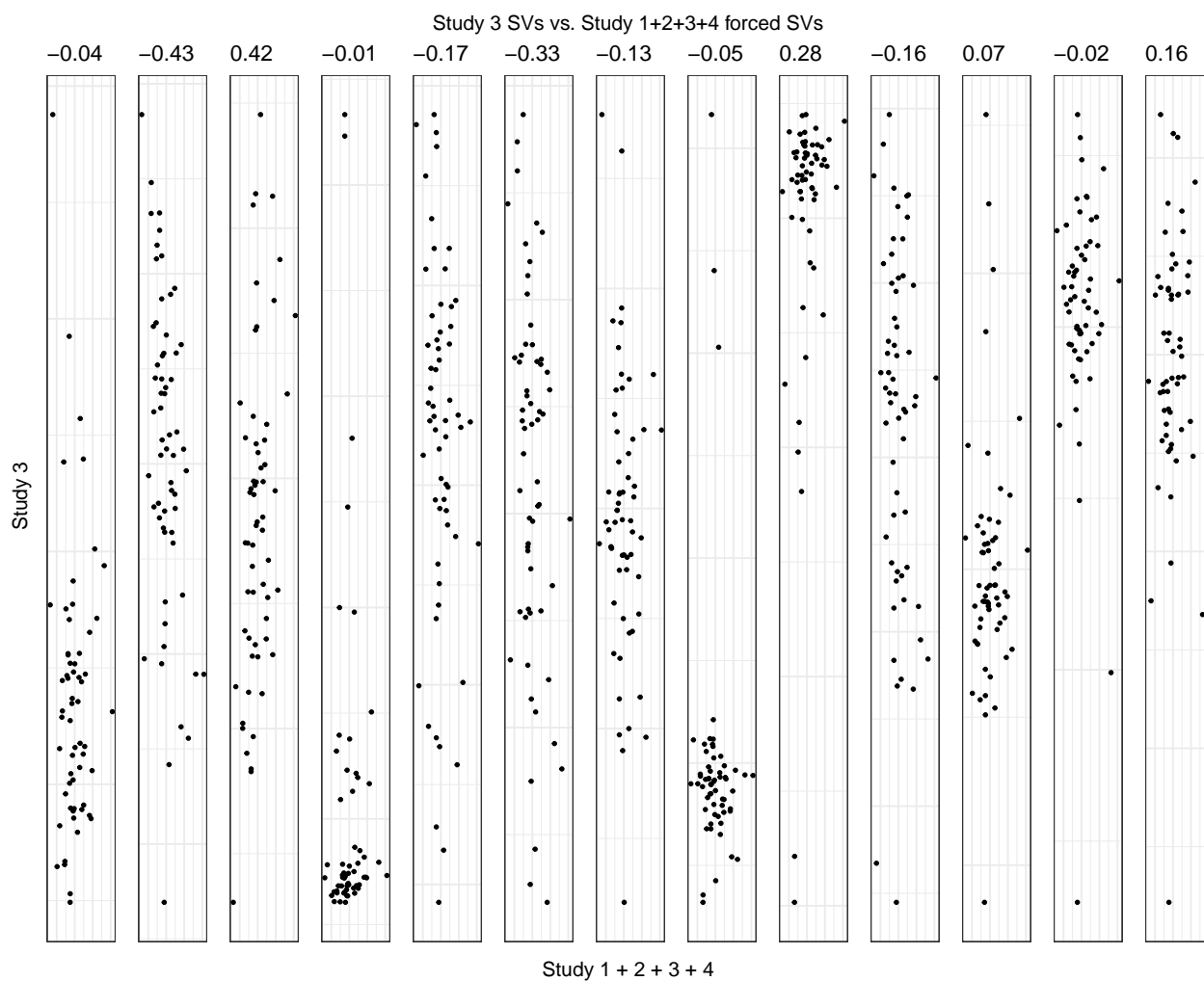
```
## [[1]]
## NULL
```



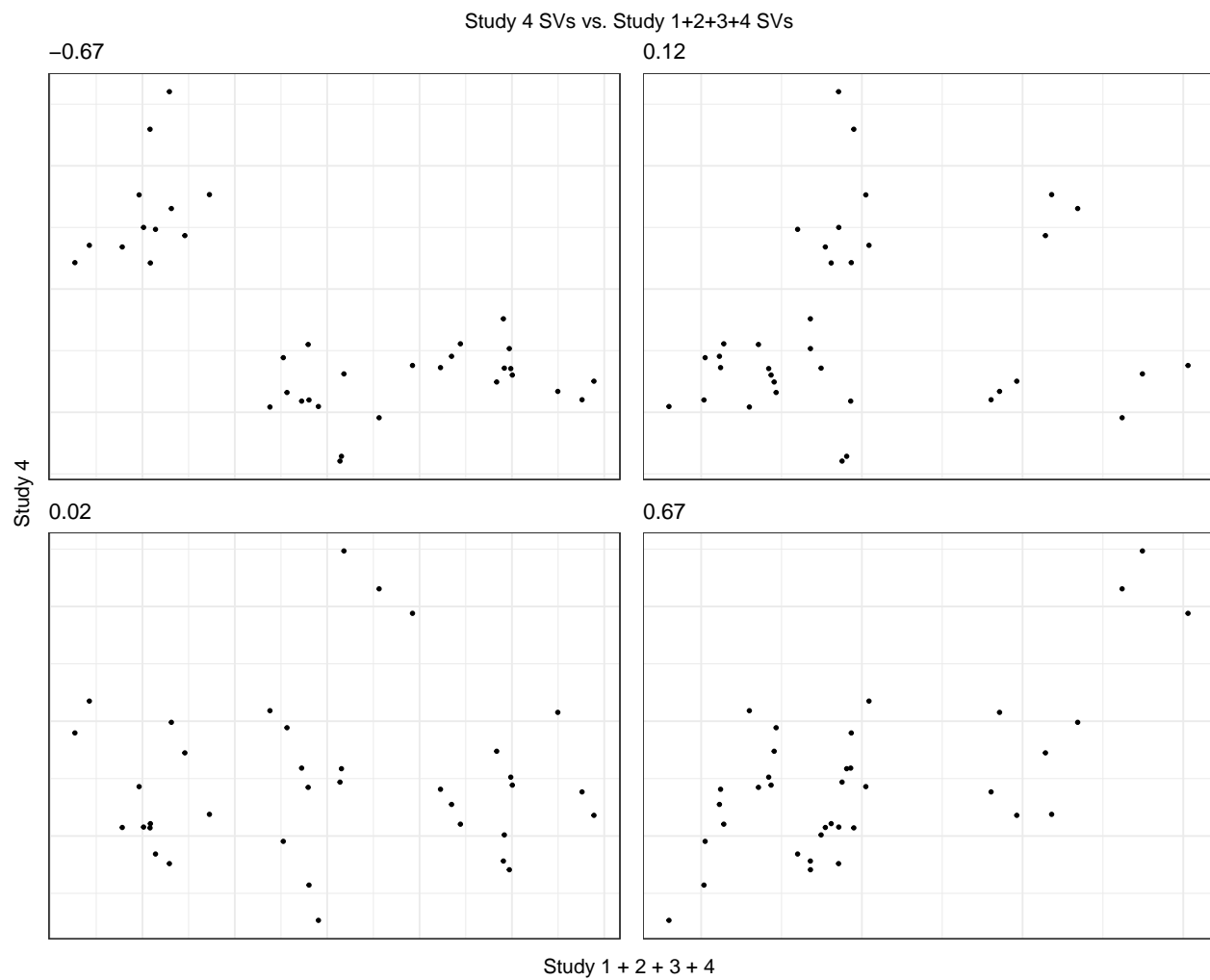
```
## [[1]]
## NULL
```



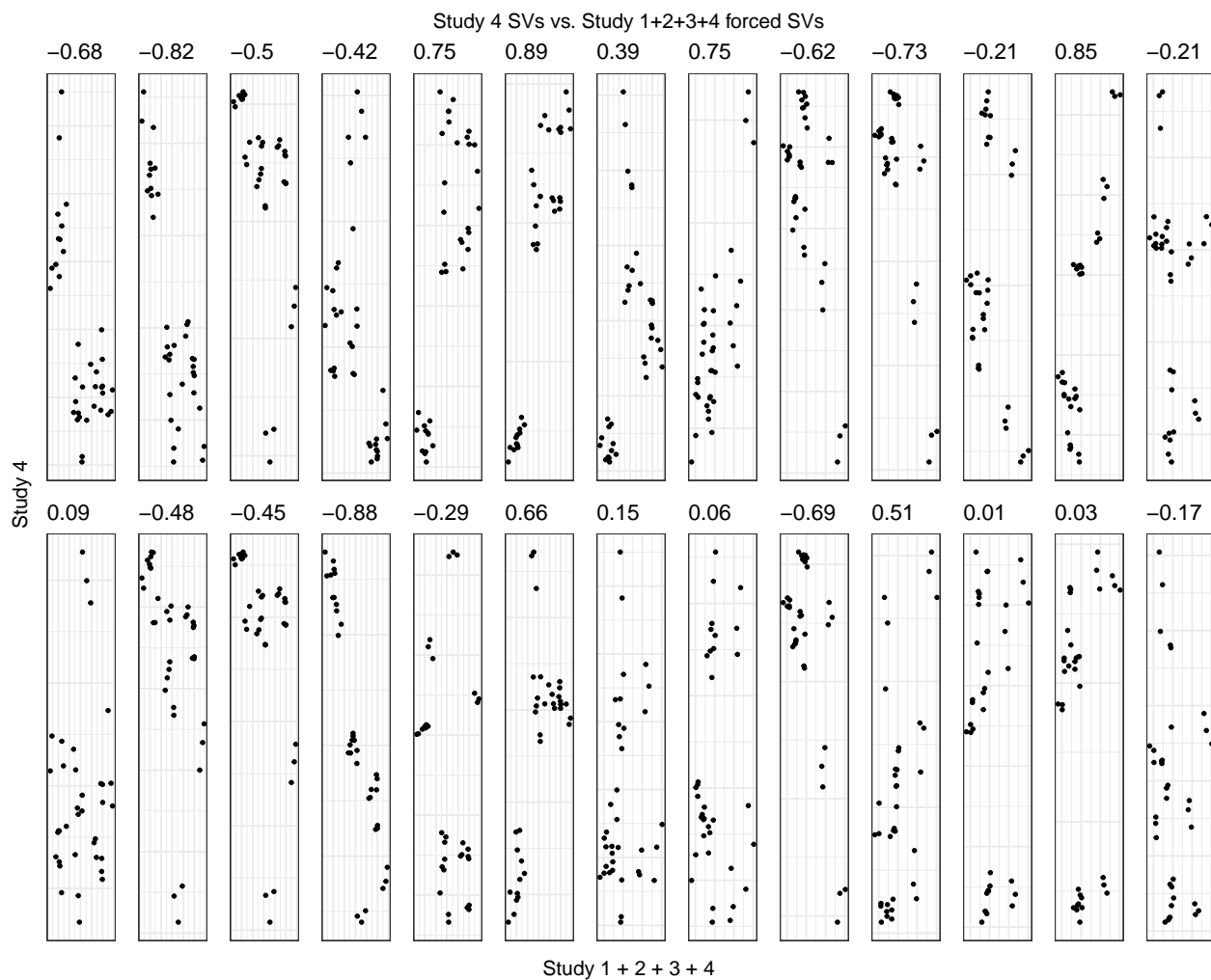
```
## [[1]]
## NULL
```



```
## [[1]]
## NULL
```



```
## [[1]]
## NULL
```

```
## [[1]]
## NULL
```