# SIS NMF Final: Diagnosis to DSD

November 29, 2014

## 1  Preparation

```
################################################################ LIBRARIES
options(java.parameters = "-Xmx4G")

library(survival)

## Loading required package:  splines

library(energy)
library(NMF)

## Loading required package:  methods
## Loading required package:  pkgmaker
## Loading required package:  registry
## Loading required package:  rngtools
## Loading required package:  cluster
## NMF - BioConductor layer [OK] | Shared memory capabilities [NO: bigmemory] | Cores 31/32
##  To enable shared memory capabilities, try:  install.extras('
## NMF
## ')

library(nnls)

library(glmulti)

## Loading required package:  rJava
##
## Attaching package:  'glmulti'
##
## The following object is masked from 'package:NMF':
##
##     consensus

library(glmnet)

## Loading required package:  Matrix
## Loaded glmnet 1.9-8

library(RColorBrewer)
library(gplots)

##
## Attaching package:  'gplots'
##
```

```
## The following object is masked from 'package:stats':
##
##     lowess

library(xtable)
library(stargazer)

##
## Please cite as:
##
##  Hlavac, Marek (2014).  stargazer:  LaTeX code and ASCII text for well-formatted regression
and summary statistics tables.
##  R package version 5.1.  http://CRAN.R-project.org/package=stargazer

load("image.rda")
```

# 2 Probe selection
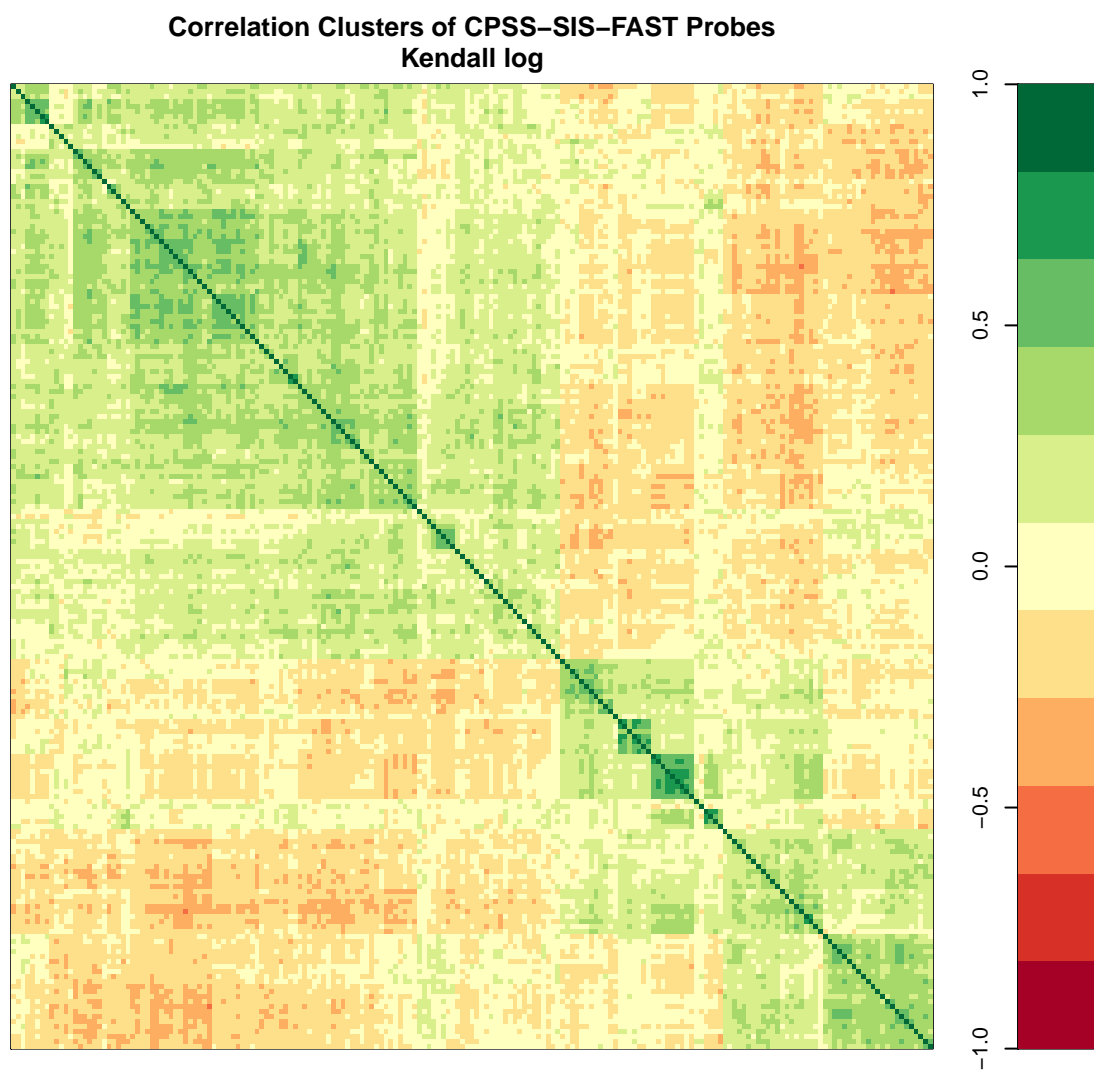
```
table(cpss.sis$sel)

##
## FALSE   TRUE
## 12807    193

mean(cpss.sis$sel)

## [1] 0.01485
```
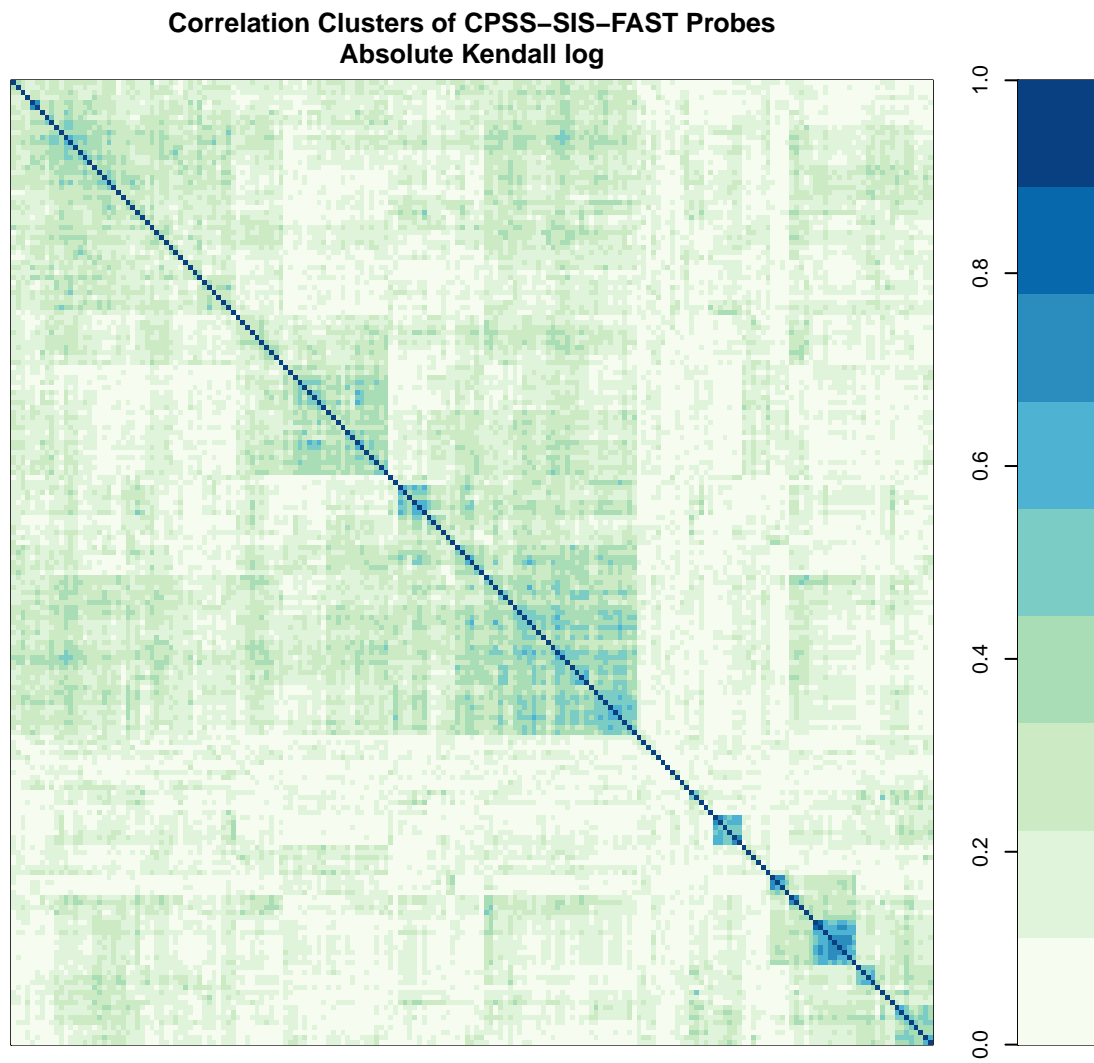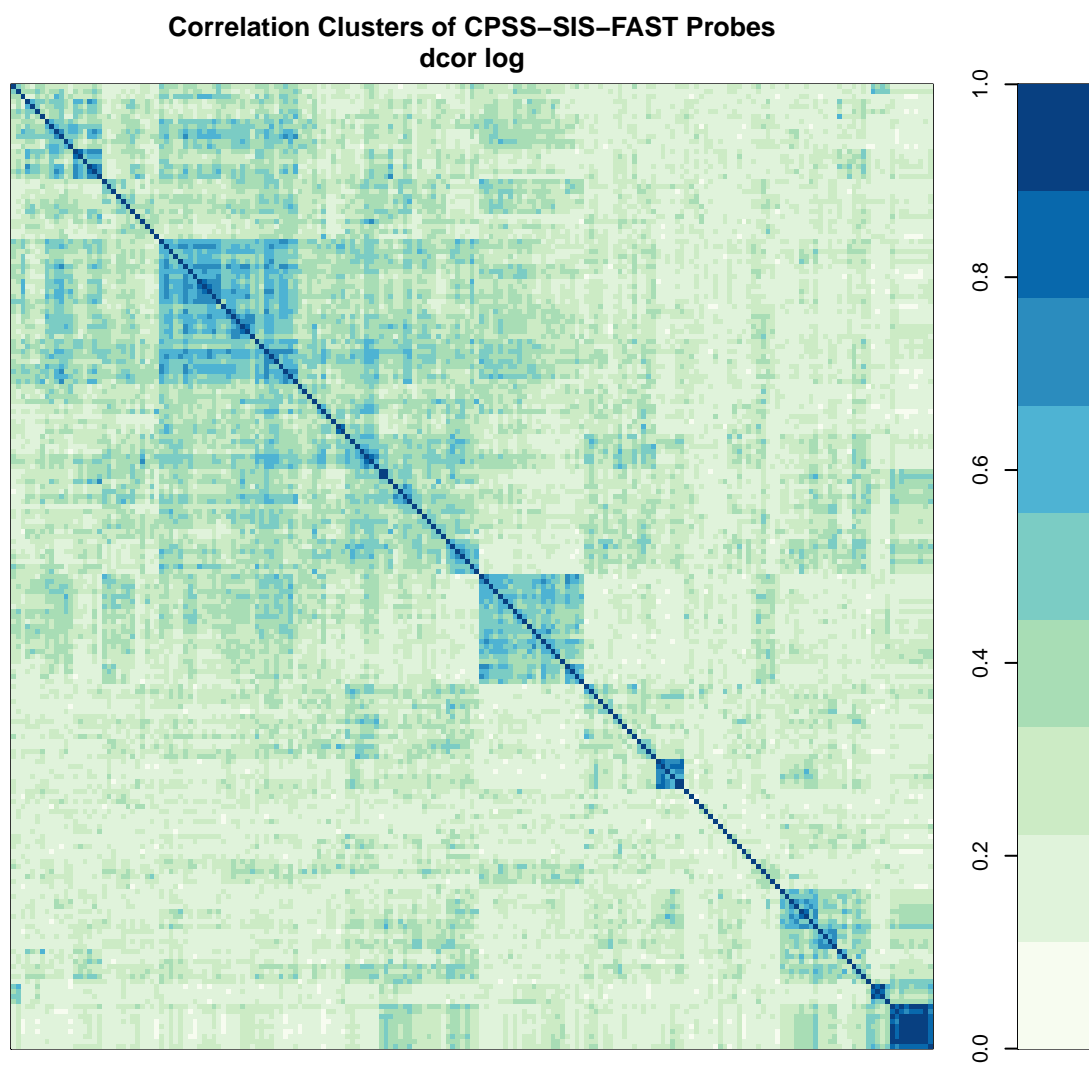
# 3 Expression correlation

```
corPlot(x.sel.kcor, main = "Correlation Clusters of CPSS-SIS-FAST Probes\nKendall log",
    useRaster = FALSE)
```

**Correlation Clusters of CPSS–SIS–FAST Probes**
**Kendall log**

```r
corPlot(abs(x.sel.kcor), zlim = c(0, 1), pal = "GnBu", main = "Correlation Clusters of CPSS-SIS-FAST Pro
    useRaster = FALSE)
```

# Correlation Clusters of CPSS–SIS–FAST Probes
## Absolute Kendall log
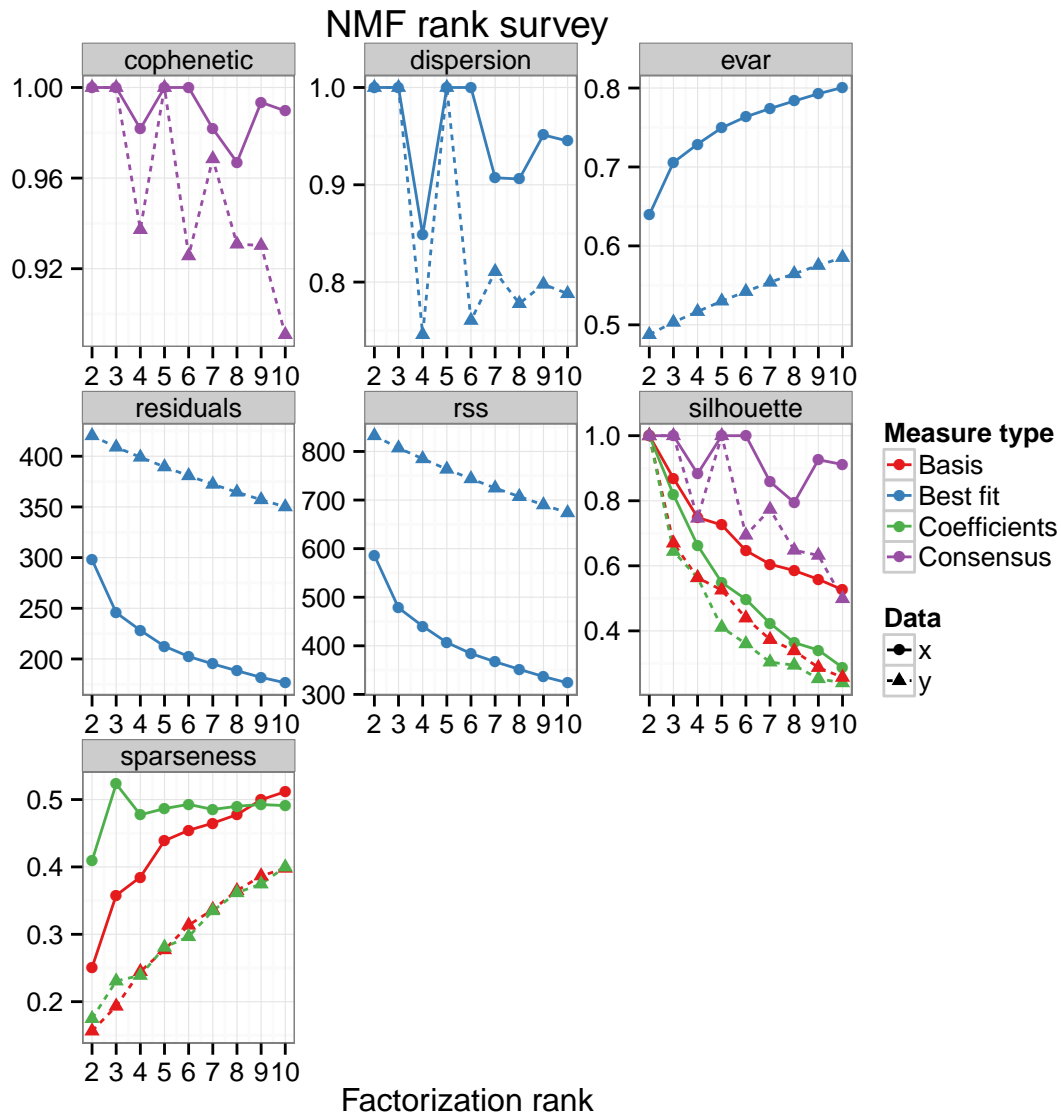


```
corPlot(x.sel.dcor, zlim = c(0, 1), pal = "GnBu", main = "Correlation Clusters of CPSS-SIS-FAST Probes\r
    useRaster = FALSE)
```

**Correlation Clusters of CPSS–SIS–FAST Probes**
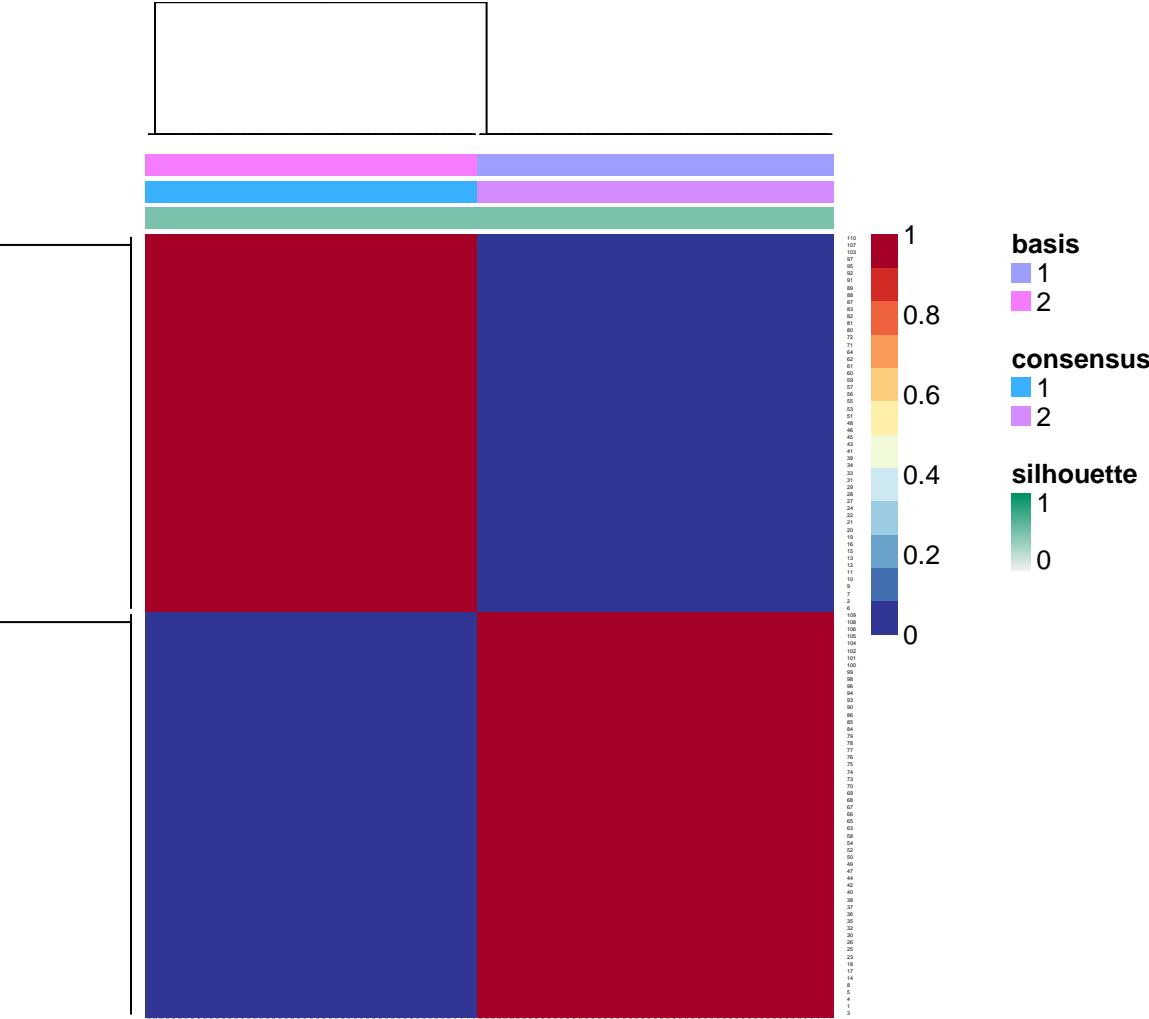**dcor log**

# 4 Factorization

```
plot(nmf.runs.rank, nmf.runs.rank.random[[1]])
```
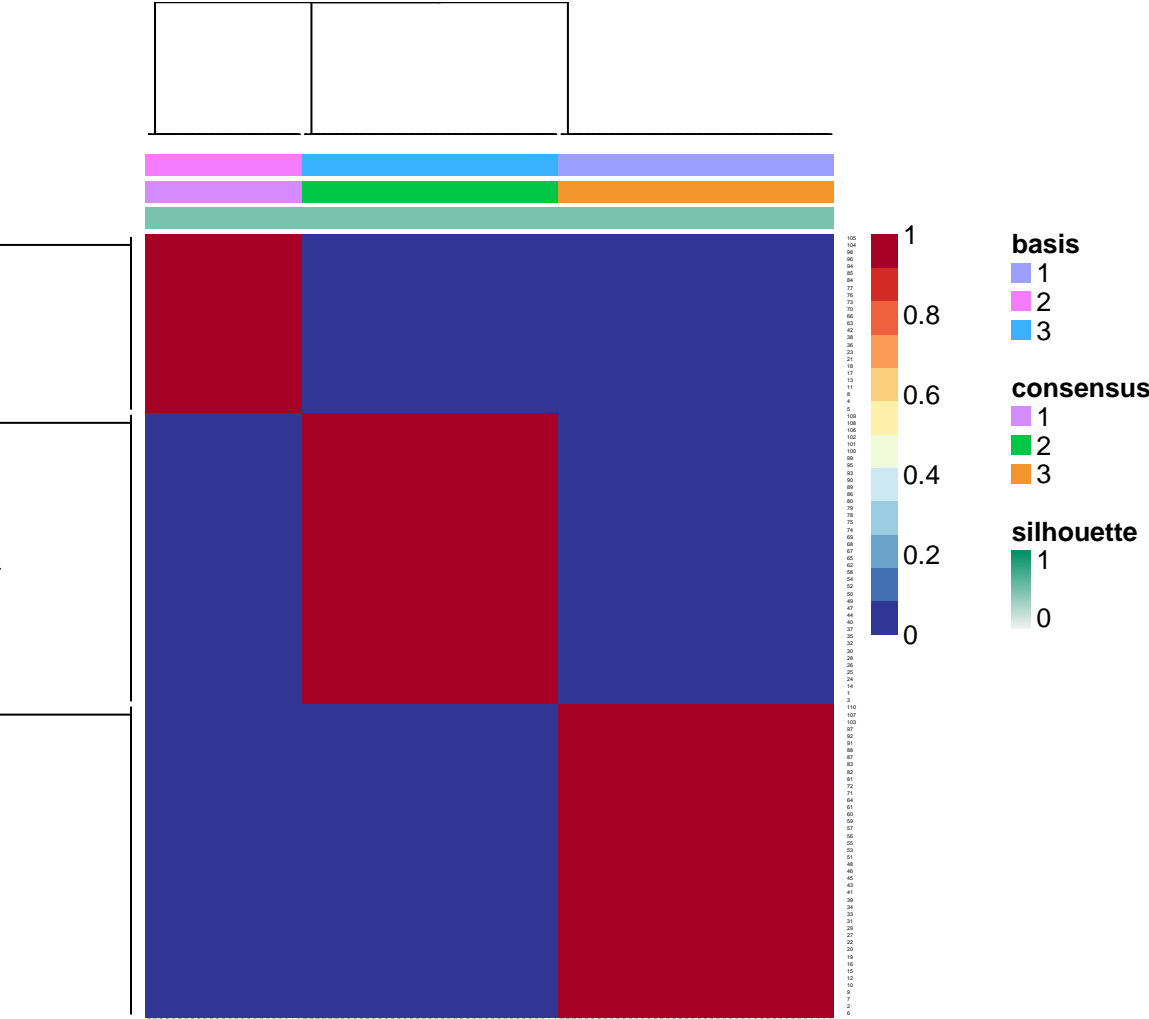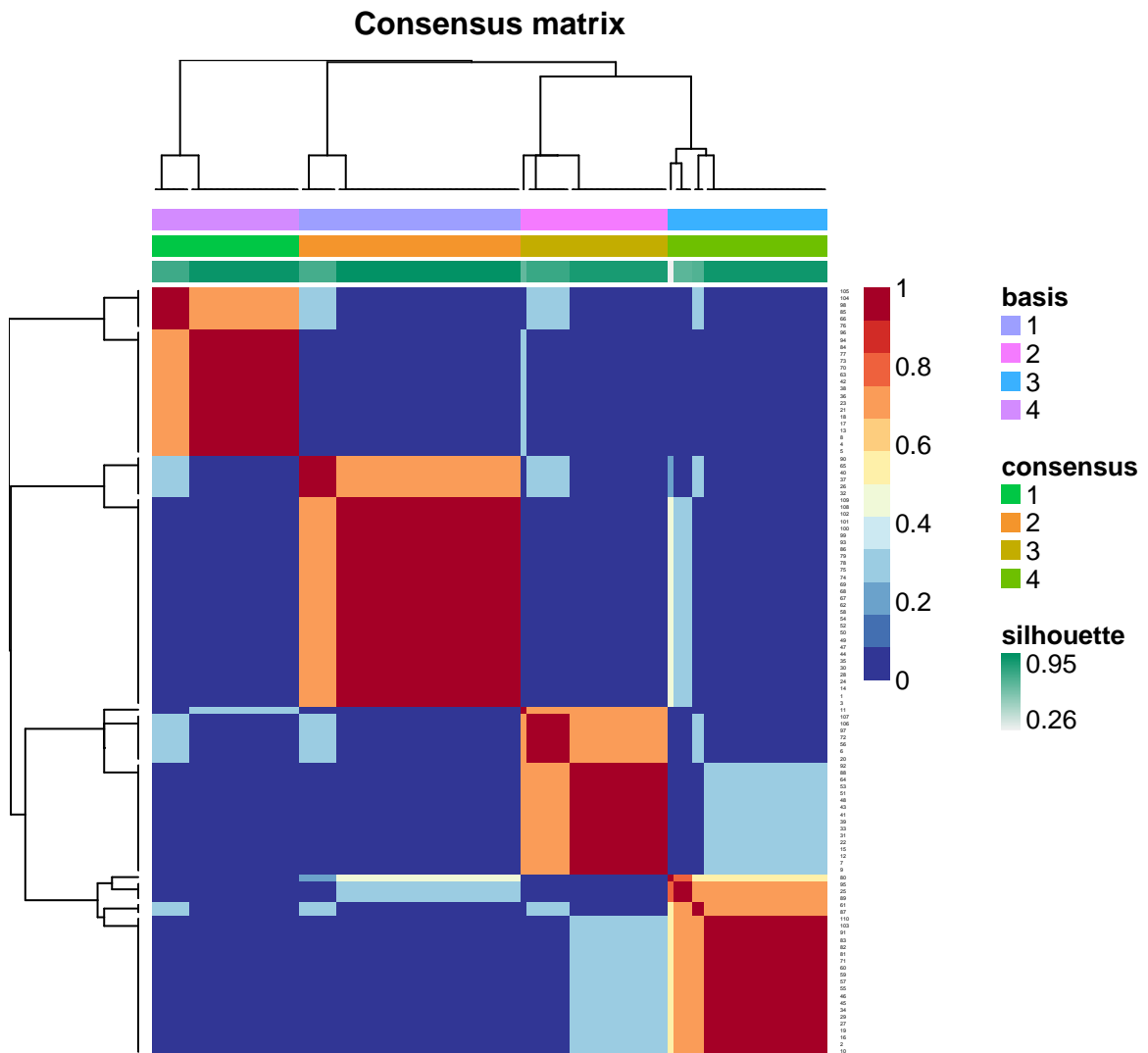
NMF rank survey

```
for (i in nmf.runs.rank$fit) {
    consensusmap(i)
}
```

# Consensus matrix

**Consensus matrix**

# Consensus matrix

# Consensus matrix

# Consensus matrix

**Consensus matrix**
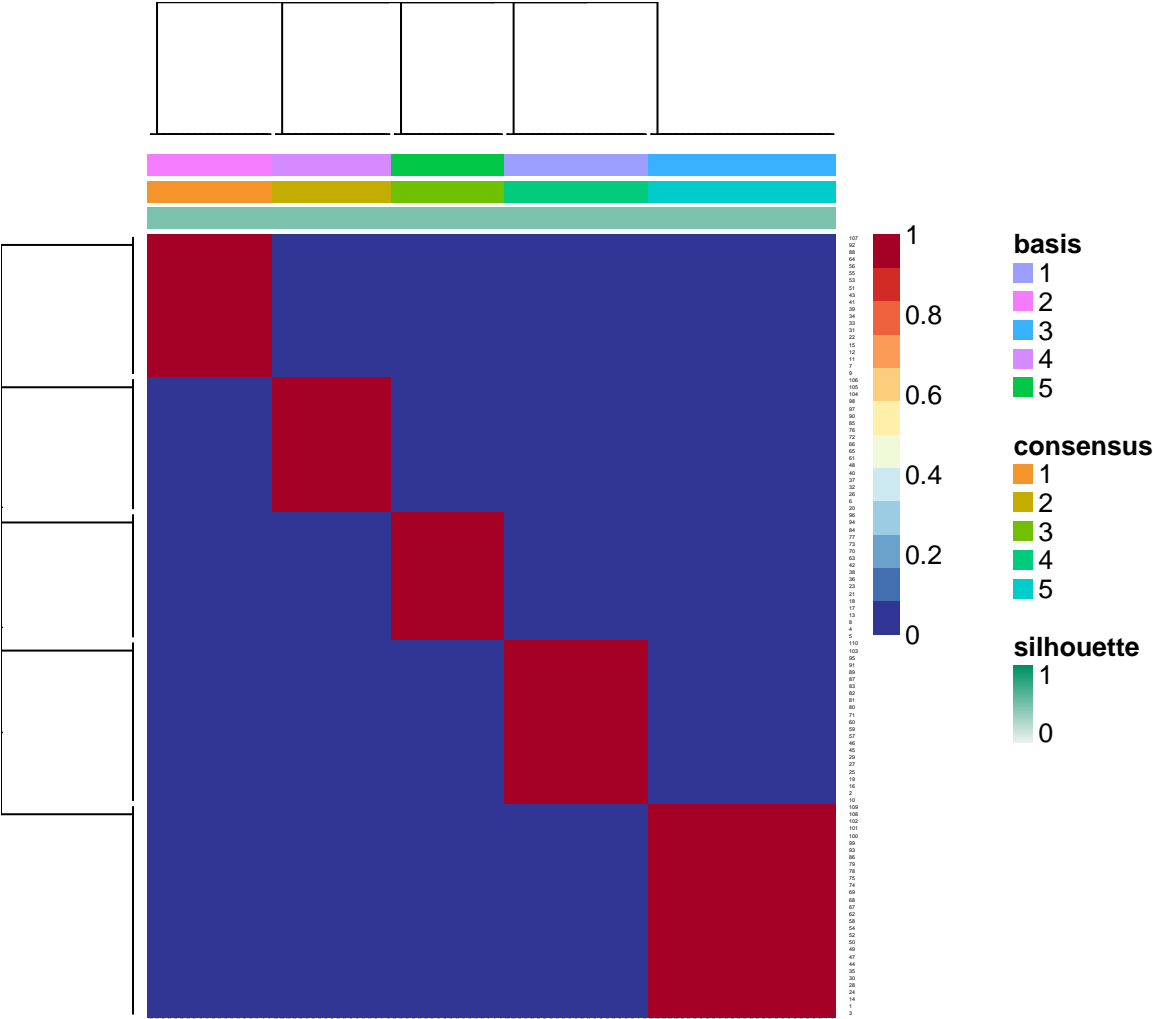
basis
1
2
3
4
5
6
7

consensus
1
2
3
4
5
6
7

silhouette
0.99
0.26

# Consensus matrix

**Consensus matrix**

basis
1
2
3
4
5
6
7
8
9

consensus
1
2
3
4
5
6
7
8
9

silhouette
0.99
0.74

**Consensus matrix**



```
plot(nmf.rankrange[-1], -temp.orig_resids.delta, type = "o", col = "black",
    pch = 21, ylim = range(-c(temp.orig_resids.delta, temp.perm_resids.delta.mean)),
    xlab = "Factorization Rank Added", ylab = "Improvement in Total Residual Error")
lines(nmf.rankrange[-1], -temp.perm_resids.delta.mean, col = "red", type = "o",
    pch = 21, lwd = 1)
for (i in 1:ncol(temp.perm_resids)) {
    lines(nmf.rankrange[-1], -temp.perm_resids.delta[, i], type = "o", col = rgb(1,
        0, 0, 0.25))
}
lines(nmf.rankrange[-1], -temp.perm_resids.delta.threshold, col = "red", lty = "dotted")
if (nmf.rank.wasauto == TRUE) {
    temp.col = "green"
} else {
    temp.col = "blue"
}
abline(v = nmf.rank, col = temp.col, lwd = 2)
legend("topright", legend = c("Original data", "Permuted data", sprintf("Selected rank (%s)",
```

```
ifelse(nmf.rank.wasauto == TRUE, "auto", "fixed"))), col = c("black", "red",
temp.col), lty = "solid", pch = 21, inset = 0.05)
```



## 4.1 Fit

```
consensusmap(nmf.final)
```

# Consensus matrix



```
basismap(nmf.final)
```

**Basis components**



```
coefmap(nmf.final)
```

**Mixture coefficients**

```
coefs.diag_dsd = apply(xlin.diag_dsd.sel, 2, function(xcol) nnls(basis(nmf.final),
    xcol)$x)
coefs.diag_rec = apply(xlin.diag_rec.sel, 2, function(xcol) nnls(basis(nmf.final),
    xcol)$x)
coefs.recr_dsd = apply(xlin.recr_dsd.sel, 2, function(xcol) nnls(basis(nmf.final),
    xcol)$x)
coefs.pdac_au = apply(xlin.pdac_au.sel, 2, function(xcol) nnls(basis(nmf.final),
    xcol)$x)
```

```
temp.pred.pairs = t(rbind(coefs.pdac_au, metapcna.scores[colnames(coefs.pdac_au)]))
colnames(temp.pred.pairs) = paste("mg", 1:ncol(temp.pred.pairs), sep = ".")
colnames(temp.pred.pairs)[ncol(temp.pred.pairs)] = "PCNA"
temp.pred.pairs = cbind(temp.pred.pairs, qpure = samps.pdac_au$purity_qpure)
pairs(temp.pred.pairs, pch = 16, cex = 1, col = ifelse(rownames(temp.pred.pairs) %in%
    colnames(xlin.diag_dsd.sel), rgb(0, 0, 0, 0.5), rgb(1, 0, 1, 0.5)))
```

## 4.2 Prediction on training set

```
nmf.final.cpv.pvals = data.frame(surv.diag_rec.p = apply(coefs.diag_rec, 1,
    function(xc) pchisq(2 * diff(coxph(y.diag_rec ~ xc)$loglik), df = 1, lower.tail = FALSE)),
    surv.diag_rec.c = apply(coefs.diag_rec, 1, function(xc) coef(coxph(y.diag_rec ~
        xc))), surv.diag_dsd.p = apply(coefs.diag_dsd, 1, function(xc) pchisq(2 *
        diff(coxph(y.diag_dsd ~ xc)$loglik), df = 1, lower.tail = FALSE)), surv.diag_dsd.c = apply(coefs
        1, function(xc) coef(coxph(y.diag_dsd ~ xc))), surv.recr_dsd.p = apply(coefs.recr_dsd,
        1, function(xc) pchisq(2 * diff(coxph(y.recr_dsd ~ xc)$loglik), df = 1,
            lower.tail = FALSE)), surv.recr_dsd.c = apply(coefs.recr_dsd, 1,
        function(xc) coef(coxph(y.recr_dsd ~ xc))), pure.p = apply(coefs.pdac_au,
        1, function(xc) cor.test(samps.pdac_au$purity_qpure, xc, method = "kendall")$p.value),
    pure.s = apply(coefs.pdac_au, 1, function(xc) cor.test(samps.pdac_au$purity_qpure,
        xc, method = "kendall")$statistic))
temp.pvals = as.matrix(nmf.final.cpv.pvals[, grepl("\\.p$", colnames(nmf.final.cpv.pvals))])
```

```r
temp.pvals.FWER = matrix(p.adjust(as.vector(temp.pvals), "holm"), nrow = nrow(temp.pvals))
colnames(temp.pvals.FWER) = paste(colnames(temp.pvals), "Holm", sep = ".")
temp.pvals.BY = matrix(p.adjust(as.vector(temp.pvals), "BY"), nrow = nrow(temp.pvals))
colnames(temp.pvals.BY) = paste(colnames(temp.pvals), "BY", sep = ".")
nmf.final.cpv.pvals = cbind(nmf.final.cpv.pvals, temp.pvals.FWER, temp.pvals.BY)
nmf.final.cpv.pvals = nmf.final.cpv.pvals[, order(colnames(nmf.final.cpv.pvals))]
```

| | pure.p | pure.p.BY | pure.p.Holm | pure.s | surv.diag_dsd.c | surv.diag_dsd.p | surv.diag_dsd.p.BY | surv.diag_dsd.p.Holm | surv.diag_rec.c | surv.diag_rec.p | surv.diag_rec.p.BY | surv.diag_rec.p.Holm | surv.recr_dsd.c | surv.recr_dsd.p | surv.recr_dsd.p.BY | surv.recr_dsd.p.Holm |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0.00 | 0.01 | 0.01 | -3.42 | -4.56 | 0.00 | 0.04 | 0.06 | -3.78 | 0.01 | 0.07 | 0.11 | -2.36 | 0.16 | 0.75 | 0.96 |
| 2 | 0.05 | 0.30 | 0.50 | -1.96 | 3.10 | 0.02 | 0.15 | 0.25 | 1.99 | 0.11 | 0.55 | 0.77 | 2.62 | 0.06 | 0.35 | 0.56 |
| 3 | 0.00 | 0.00 | 0.00 | 3.91 | -5.46 | 0.00 | 0.00 | 0.00 | -2.64 | 0.03 | 0.19 | 0.33 | -4.66 | 0.00 | 0.04 | 0.04 |
| 4 | 0.03 | 0.19 | 0.33 | -2.19 | -0.95 | 0.44 | 1.00 | 1.00 | -1.68 | 0.16 | 0.75 | 0.96 | 0.56 | 0.68 | 1.00 | 1.00 |
| 5 | 0.07 | 0.38 | 0.57 | -1.81 | 2.84 | 0.02 | 0.12 | 0.21 | 2.73 | 0.01 | 0.12 | 0.20 | 1.37 | 0.29 | 1.00 | 1.00 |
| 6 | 0.78 | 1.00 | 1.00 | 0.28 | 6.61 | 0.00 | 0.00 | 0.00 | 5.49 | 0.00 | 0.00 | 0.00 | 3.97 | 0.00 | 0.04 | 0.06 |

Table 1: Resubstitution prediction, all tests

| | pure.p.Holm | pure.s | surv.diag_dsd.c | surv.diag_dsd.p.Holm | surv.diag_rec.c | surv.diag_rec.p.Holm | surv.recr_dsd.c | surv.recr_dsd.p.Holm |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.01 | -3.42 | -4.56 | 0.06 | -3.78 | 0.11 | -2.36 | 0.96 |
| 2 | 0.50 | -1.96 | 3.10 | 0.25 | 1.99 | 0.77 | 2.62 | 0.56 |
| 3 | 0.00 | 3.91 | -5.46 | 0.00 | -2.64 | 0.33 | -4.66 | 0.04 |
| 4 | 0.33 | -2.19 | -0.95 | 1.00 | -1.68 | 0.96 | 0.56 | 1.00 |
| 5 | 0.57 | -1.81 | 2.84 | 0.21 | 2.73 | 0.20 | 1.37 | 1.00 |
| 6 | 1.00 | 0.28 | 6.61 | 0.00 | 5.49 | 0.00 | 3.97 | 0.06 |

Table 2: Resubstitution prediction, Holm MTC only

## MTC P-values

```r
print(asreg.result)
```

```
## glmulti.analysis
## Method: h / Fitting: coxph / IC used: bic
## Level: 2 / Marginality: TRUE
## From 100 models:
## Best IC: 551.558245978867
## Best model:
## [1] "Surv(time, event) ~ 1 + mg.5 + mg.6"
## Evidence weight: 0.121525761025609
## Worst IC: 561.093163081812
## 5 models within 2 IC units.
## 71 models to reach 95% of evidence weight.
```

```r
coef(asreg.result)
```

```
##           Estimate Uncond. variance Nb models Importance +/- (alpha=0.05)
## mg.1:mg.4 -0.00342         0.001098         1    0.001595           0.06569
## mg.2:mg.4  0.06847         0.024924         2    0.003030           0.31297
## mg.4:mg.5 -0.05064         0.019737         1    0.003679           0.27850
## mg.4:mg.6  0.02227         0.015093         2    0.004699           0.24355
## mg.2:mg.3  0.08587         0.057059         3    0.004885           0.47354
## mg.1:mg.2  0.10761         0.069012         2    0.004997           0.52078
## mg.1:mg.5 -0.01322         0.015809         2    0.005116           0.24925
## mg.3:mg.4 -0.25668         0.347670         4    0.008845           1.16890
## mg.3:mg.5  0.01013         0.047858         3    0.009870           0.43368
## mg.2:mg.5  0.09367         0.064793         4    0.010069           0.50461
```

```
## mg.5:mg.6  0.10404          0.413552        8   0.036678          1.27485
## mg.1:mg.6  0.75549          3.893974        8   0.039107          3.91192
## mg.1:mg.3 -1.05886          7.166918        9   0.043683          5.30712
## mg.2:mg.6 -1.53508          9.333667       13   0.074749          6.05646
## mg.3:mg.6  2.60281         29.616218       14   0.097762         10.78842
## mg.4        0.09036         0.113089       33   0.136430          0.66666
## mg.2        0.33327         0.778397       42   0.219557          1.74902
## mg.1       -1.85625         6.374836       54   0.439670          5.00527
## mg.5        1.08749         2.127350       49   0.444550          2.89143
## mg.3       -2.66620         7.511723       65   0.592820          5.43329
## mg.6        5.31298         5.463222       90   0.932557          4.63359
```

```
summary(asreg.result@objects[[1]])
```

```
## Call:
## fitfunc(formula = as.formula(x), data = data)
##
##   n= 110, number of events= 70
##
##        coef exp(coef) se(coef)    z Pr(>|z|)
## mg.5   2.81     16.65     1.08 2.60   0.0093
## mg.6   6.99   1089.19     1.19 5.87  4.4e-09
##
##      exp(coef) exp(-coef) lower .95 upper .95
## mg.5      16.7   0.060050         2       139
## mg.6    1089.2   0.000918       105     11264
##
## Concordance= 0.702  (se = 0.038 )
## Rsquare= 0.265   (max possible= 0.995 )
## Likelihood ratio test= 33.9  on 2 df,   p=4.29e-08
## Wald test            = 38.8  on 2 df,   p=3.77e-09
## Score (logrank) test = 42.9  on 2 df,   p=4.87e-10
```

```
plot(asreg.result, type = "p")
```

**IC profile**



```
plot(asreg.result, type = "s")
```

**Model−averaged importance of terms**



```
plot(asreg.result, type = "w")
```

**Profile of model weights**



```
glmnet.coef.1se

## 6 x 1 sparse Matrix of class "dgCMatrix"
##             1
## mg.1  .
## mg.2  .
## mg.3 -0.1635
## mg.4  .
## mg.5  .
## mg.6  3.0808

glmnet.coef.min
```

```
## 6 x 1 sparse Matrix of class "dgCMatrix"
##           1
## mg.1 -2.372
## mg.2   .
## mg.3 -3.002
## mg.4   .
## mg.5  1.102
## mg.6  4.391
```

```r
plot(glmnet.fit.cv)
```

**LASSO**



```r
plot(glmnet.fit.cv$glmnet.fit, label = TRUE)
abline(v = sum(abs(glmnet.coef.1se)))
abline(v = sum(abs(glmnet.coef.min)))
```

## 4.3 Prediction on validation sets

```r
load("../../data/15_validation.rda")
```

```r
val.basis = basis(nmf.final)
rownames(GSE21501.lingex) = GSE21501.feat$Gene.symbol
rownames(GSE28735.lingex) = GSE28735.feat$Gene.symbol
GSE21501.lingex.for_basis = GSE21501.lingex[match(rownames(val.basis), rownames(GSE21501.lingex)),
    ]
GSE28735.lingex.for_basis = GSE28735.lingex[match(rownames(val.basis), rownames(GSE28735.lingex)),
    ]
GSE21501.lingex.for_basis[is.na(GSE21501.lingex.for_basis)] = 0
GSE28735.lingex.for_basis[is.na(GSE28735.lingex.for_basis)] = 0

GSE21501.coefs = apply(GSE21501.lingex.for_basis, 2, function(xcol) nnls(val.basis,
```

```
    xcol)$x)
GSE28735.coefs = apply(GSE28735.lingex.for_basis, 2, function(xcol) nnls(val.basis,
    xcol)$x)
```

```
apply(GSE21501.coefs, 1, function(xc) coxph(Surv(time, event) ~ xc, data = GSE21501.samp))

## [[1]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE21501.samp)
##
##
##     coef exp(coef) se(coef)    z    p
## xc -5.31   0.00495     3.44 -1.54 0.12
##
## Likelihood ratio test=2.57  on 1 df, p=0.109  n= 102, number of events= 66
##
## [[2]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE21501.samp)
##
##
##    coef exp(coef) se(coef)     z   p
## xc 2.17      8.79      2.6 0.836 0.4
##
## Likelihood ratio test=0.68  on 1 df, p=0.41  n= 102, number of events= 66
##
## [[3]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE21501.samp)
##
##
##     coef exp(coef) se(coef)      z    p
## xc -1.09     0.337     2.73 -0.399 0.69
##
## Likelihood ratio test=0.16  on 1 df, p=0.688  n= 102, number of events= 66
##
## [[4]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE21501.samp)
##
##
##      coef exp(coef) se(coef)       z    p
## xc -0.226     0.798     2.47 -0.0914 0.93
##
## Likelihood ratio test=0.01  on 1 df, p=0.927  n= 102, number of events= 66
##
## [[5]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE21501.samp)
##
##
##    coef exp(coef) se(coef)     z    p
## xc 2.17      8.75     3.04 0.713 0.48
```

```
##
## Likelihood ratio test=0.49  on 1 df, p=0.486  n= 102, number of events= 66
##
## [[6]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE21501.samp)
##
##
##    coef exp(coef) se(coef)    z    p
## xc 2.57      13.1     2.08 1.24 0.22
##
## Likelihood ratio test=1.45  on 1 df, p=0.229  n= 102, number of events= 66
```

```
apply(GSE21501.coefs, 1, function(xc) coxph(Surv(time, event) ~ tstage + nstage +
    xc, data = GSE21501.samp))
```

```
## [[1]]
## Call:
## coxph(formula = Surv(time, event) ~ tstage + nstage + xc, data = GSE21501.samp)
##
##
##          coef exp(coef) se(coef)      z     p
## tstage -0.166   0.84664    0.281 -0.593 0.550
## nstage  0.653   1.92135    0.315  2.074 0.038
## xc     -5.237   0.00532    3.545 -1.477 0.140
##
## Likelihood ratio test=7.13  on 3 df, p=0.0678  n= 97, number of events= 63
##    (5 observations deleted due to missingness)
##
## [[2]]
## Call:
## coxph(formula = Surv(time, event) ~ tstage + nstage + xc, data = GSE21501.samp)
##
##
##          coef exp(coef) se(coef)      z     p
## tstage -0.161     0.851    0.287 -0.560 0.580
## nstage  0.643     1.903    0.316  2.039 0.041
## xc      1.025     2.788    2.658  0.386 0.700
##
## Likelihood ratio test=4.94  on 3 df, p=0.176  n= 97, number of events= 63
##    (5 observations deleted due to missingness)
##
## [[3]]
## Call:
## coxph(formula = Surv(time, event) ~ tstage + nstage + xc, data = GSE21501.samp)
##
##
##          coef exp(coef) se(coef)      z     p
## tstage -0.153     0.858    0.286 -0.535 0.590
## nstage  0.659     1.933    0.314  2.097 0.036
## xc     -1.379     0.252    2.796 -0.493 0.620
##
## Likelihood ratio test=5.04  on 3 df, p=0.169  n= 97, number of events= 63
##    (5 observations deleted due to missingness)
```

```
##
## [[4]]
## Call:
## coxph(formula = Surv(time, event) ~ tstage + nstage + xc, data = GSE21501.samp)
##
##
##           coef exp(coef) se(coef)      z     p
## tstage -0.158     0.854    0.285 -0.554 0.580
## nstage  0.678     1.970    0.319  2.126 0.034
## xc     -1.041     0.353    2.463 -0.423 0.670
##
## Likelihood ratio test=4.98  on 3 df, p=0.173  n= 97, number of events= 63
##     (5 observations deleted due to missingness)
##
## [[5]]
## Call:
## coxph(formula = Surv(time, event) ~ tstage + nstage + xc, data = GSE21501.samp)
##
##
##           coef exp(coef) se(coef)      z     p
## tstage -0.179     0.836    0.289 -0.621 0.530
## nstage  0.641     1.898    0.315  2.033 0.042
## xc      1.878     6.543    3.313  0.567 0.570
##
## Likelihood ratio test=5.11  on 3 df, p=0.164  n= 97, number of events= 63
##     (5 observations deleted due to missingness)
##
## [[6]]
## Call:
## coxph(formula = Surv(time, event) ~ tstage + nstage + xc, data = GSE21501.samp)
##
##
##           coef exp(coef) se(coef)      z     p
## tstage -0.110     0.896    0.284 -0.388 0.700
## nstage  0.657     1.928    0.316  2.077 0.038
## xc      2.510    12.310    2.164  1.160 0.250
##
## Likelihood ratio test=6.07  on 3 df, p=0.108  n= 97, number of events= 63
##     (5 observations deleted due to missingness)

apply(GSE21501.coefs, 1, function(xc) anova(coxph(Surv(time, event) ~ tstage +
    nstage + xc, data = GSE21501.samp)))

## [[1]]
## Analysis of Deviance Table
##  Cox model: response is Surv(time, event)
## Terms added sequentially (first to last)
##
##        loglik Chisq Df Pr(>|Chi|)
## NULL     -242
## tstage   -242  0.01  1      0.928
## nstage   -239  4.79  1      0.029
## xc       -238  2.34  1      0.126
##
```

```
## [[2]]
## Analysis of Deviance Table
##  Cox model: response is Surv(time, event)
## Terms added sequentially (first to last)
##
##        loglik Chisq Df Pr(>|Chi|)
## NULL     -242
## tstage   -242  0.01  1      0.928
## nstage   -239  4.79  1      0.029
## xc       -239  0.15  1      0.702
##
## [[3]]
## Analysis of Deviance Table
##  Cox model: response is Surv(time, event)
## Terms added sequentially (first to last)
##
##        loglik Chisq Df Pr(>|Chi|)
## NULL     -242
## tstage   -242  0.01  1      0.928
## nstage   -239  4.79  1      0.029
## xc       -239  0.25  1      0.619
##
## [[4]]
## Analysis of Deviance Table
##  Cox model: response is Surv(time, event)
## Terms added sequentially (first to last)
##
##        loglik Chisq Df Pr(>|Chi|)
## NULL     -242
## tstage   -242  0.01  1      0.928
## nstage   -239  4.79  1      0.029
## xc       -239  0.18  1      0.668
##
## [[5]]
## Analysis of Deviance Table
##  Cox model: response is Surv(time, event)
## Terms added sequentially (first to last)
##
##        loglik Chisq Df Pr(>|Chi|)
## NULL     -242
## tstage   -242  0.01  1      0.928
## nstage   -239  4.79  1      0.029
## xc       -239  0.31  1      0.577
##
## [[6]]
## Analysis of Deviance Table
##  Cox model: response is Surv(time, event)
## Terms added sequentially (first to last)
##
##        loglik Chisq Df Pr(>|Chi|)
## NULL     -242
## tstage   -242  0.01  1      0.928
## nstage   -239  4.79  1      0.029
```

```
## xc          -239  1.27  1        0.259
```

```
apply(GSE28735.coefs, 1, function(xc) coxph(Surv(time, event) ~ xc, data = GSE28735.samp))
```

```
## [[1]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE28735.samp)
##
##
##     coef exp(coef) se(coef)    z    p
## xc -4.5    0.0111     3.74 -1.2 0.23
##
## Likelihood ratio test=1.55   on 1 df, p=0.213   n= 42, number of events= 29
##
## [[2]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE28735.samp)
##
##
##     coef exp(coef) se(coef)    z    p
## xc 3.46      31.7     2.53 1.36 0.17
##
## Likelihood ratio test=1.63   on 1 df, p=0.201   n= 42, number of events= 29
##
## [[3]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE28735.samp)
##
##
##      coef exp(coef) se(coef)     z    p
## xc -7.79  0.000415     3.22 -2.42 0.015
##
## Likelihood ratio test=6.32   on 1 df, p=0.0119   n= 42, number of events= 29
##
## [[4]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE28735.samp)
##
##
##     coef exp(coef) se(coef)     z   p
## xc 0.72      2.05     2.81 0.256 0.8
##
## Likelihood ratio test=0.06   on 1 df, p=0.801   n= 42, number of events= 29
##
## [[5]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE28735.samp)
##
##
##       coef exp(coef) se(coef)     z    p
## xc -0.469     0.625     2.61 -0.18 0.86
##
## Likelihood ratio test=0.03   on 1 df, p=0.857   n= 42, number of events= 29
##
```

```
## [[6]]
## Call:
## coxph(formula = Surv(time, event) ~ xc, data = GSE28735.samp)
##
##
##    coef exp(coef) se(coef)    z    p
## xc 5.73       308     2.32 2.47 0.014
##
## Likelihood ratio test=5.68  on 1 df, p=0.0171  n= 42, number of events= 29
```

## 4.4   MSigDB score correlation thresholding

```
temp.sel_cols = apply(abs(nmf.final.msigdb.corr) >= sig.corr.threshold, 2, any)
heatmap.2(nmf.final.msigdb.corr[, temp.sel_cols], trace = "none", scale = "none",
    useRaster = TRUE, col = brewer.pal(11, "PiYG"), symbreaks = TRUE)
```

```
heatmap.2(nmf.final.msigdb.corr[, temp.sel_cols], trace = "none", scale = "none",
    useRaster = TRUE, col = brewer.pal(3, "PiYG"), breaks = c(-1, -sig.corr.threshold,
        sig.corr.threshold, 1))
```



```
temp.sig_id = colnames(nmf.final.msigdb.corr)
temp.sig_class = gsub("\\..*", "", temp.sig_id)
temp.nsigs = length(temp.sig_id)
temp.nmeta = nrow(nmf.final.msigdb.corr)
tables = lapply(1:temp.nmeta, function(metagene_i) {
    tapply(1:temp.nsigs, temp.sig_class, function(sig_class_is) {
        all_cors = nmf.final.msigdb.corr[, sig_class_is]
        this_cors = all_cors[metagene_i, ]
        this_ids = temp.sig_id[sig_class_is]

        all_sig_cors = abs(all_cors) >= sig.corr.threshold
        this_sig_cors = all_sig_cors[metagene_i, ]
```

```
        sigs_to_report = which(this_sig_cors)

        if (length(sigs_to_report) == 0) {
           table = data.frame(GeneSet = c(), Correlation = c(), Metagenes = c())
        } else {
           table = data.frame(GeneSet = this_ids[sigs_to_report], Correlation = this_cors[sigs_to_repor
               Metagenes = apply(all_cors[, sigs_to_report, drop = FALSE],
                 2, function(cors) {
                   sel = abs(cors) >= sig.corr.threshold
                   # A positive number implies that positive GSVA signal is associated with
                   # worse prognosis
                   paste(which(sel) * sign(cors[which(sel)]) * sign(nmf.final.cpv.pvals$surv.diag_dsd.c
                     collapse = ",")
                 }))
           table = table[order(-(table$Correlation)), ]
           rownames(table) <- NULL
        }
        table
    }, simplify = FALSE)
})
tables

## [[1]]
## [[1]]$c1
## data frame with 0 columns and 0 rows
##
## [[1]]$c2
##                                            GeneSet Correlation Metagenes
## 1                 c2.KATSANOU_ELAVL1_TARGETS_SIGNED      0.5096        -1
## 2               c2.ZHAN_MULTIPLE_MYELOMA_CD2_SIGNED      0.5086        -1
## 3          c2.SMID_BREAST_CANCER_NORMAL_LIKE_SIGNED      0.5080        -1
## 4                  c2.GREENBAUM_E2A_TARGETS_SIGNED     -0.5009         1
## 5   c2.MARTORIATI_MDM4_TARGETS_NEUROEPITHELIUM_SIGNED   -0.5012         1
## 6                        c2.YU_MYC_TARGETS_SIGNED     -0.5029         1
## 7             c2.SABATES_COLORECTAL_ADENOMA_SIGNED     -0.5036         1
## 8                        c2.WINTER_HYPOXIA_SIGNED     -0.5241         1
## 9              c2.WEST_ADRENOCORTICAL_TUMOR_SIGNED     -0.5416         1
## 10     c2.SHIPP_DLBCL_VS_FOLLICULAR_LYMPHOMA_SIGNED     -0.5500         1
## 11             c2.LEE_EARLY_T_LYMPHOCYTE_SIGNED     -0.5517         1
## 12             c2.HAHTOLA_SEZARY_SYNDROM_SIGNED     -0.5641         1
##
## [[1]]$c3
##          GeneSet Correlation Metagenes
## 1 c3.V$STAT5A_01      0.5234        -1
##
## [[1]]$c4
##        GeneSet Correlation Metagenes
## 1 c4.MODULE_51      0.5399        -1
##
## [[1]]$c5
##                                  GeneSet Correlation Metagenes
## 1 c5.PHOSPHORIC_DIESTER_HYDROLASE_ACTIVITY      0.5113        -1
##
## [[1]]$c6
```

```
## data frame with 0 columns and 0 rows
##
## [[1]]$c7
##                                                           GeneSet Correlation
## 1          c7.GSE20715_0H_VS_48H_OZONE_TLR4_KO_LUNG_SIGNED        0.5160
## 2 c7.GSE22886_IGM_MEMORY_BCELL_VS_BLOOD_PLASMA_CELL_SIGNED        0.5019
## 3          c7.GSE34205_HEALTHY_VS_RSV_INF_INFANT_PBMC_SIGNED        0.5002
##   Metagenes
## 1        -1
## 2        -1
## 3        -1
##
##
## [[2]]
## [[2]]$c1
## data frame with 0 columns and 0 rows
##
## [[2]]$c2
##                                                                                                    Gen
## 1                          c2.REACTOME_EXTRACELLULAR_MATRIX_ORGANIZATION/c2.REACTOME_COLLAGEN_FORMA
## 2                                                                           c2.PID_SYNDECAN_1_PAT
## 3                                                          c2.VERRECCHIA_DELAYED_RESPONSE_TO_T
## 4                                                                             c2.PID_INTEGRIN1_PAT
## 5                                                                        c2.PID_AVB3_INTEGRIN_PAT
## 6   c2.FARMER_BREAST_CANCER_CLUSTER_5/c2.ANASTASSIOU_CANCER_MESENCHYMAL_TRANSITION_SIGNATURE/c4.GNF2_G
## 7                                             c2.YAO_TEMPORAL_RESPONSE_TO_PROGESTERONE_CLUSTE
## 8                                                             c2.KEGG_ECM_RECEPTOR_INTERAC
## 9                                                            c2.VERRECCHIA_RESPONSE_TO_TGFE
## 10                                                  c2.VERRECCHIA_EARLY_RESPONSE_TO_T
## 11                                                                     c2.KEGG_FOCAL_ADHE
## 12                                                      c2.MAHADEVAN_GIST_MORPHOLOGICAL_SW
## 13                                                          c2.CAIRO_LIVER_DEVELOPMENT_SI
## 14                                                                      c2.PID_INTEGRIN3_PAT
## 15                                                         c2.KEGG_BASAL_CELL_CARCI
## 16                                                              c2.BURTON_ADIPOGENES
## 17                                                        c2.VERRECCHIA_RESPONSE_TO_TGFE
## 18                                                          c2.CROMER_TUMORIGENESIS_SI
## 19                                                     c2.WIEDERSCHAIN_TARGETS_OF_BMI1_AND_F
## 20                                                            c2.ROZANOV_MMP14_TARGETS_SU
## 21                                                                c2.PID_WNT_SIGNALING_PAT
## 22                            c2.KEGG_ARRHYTHMOGENIC_RIGHT_VENTRICULAR_CARDIOMYOPATHY_
## 23                                                      c2.LABBE_TARGETS_OF_TGFB1_AND_WNT3A_SI
## 24                                                         c2.LIEN_BREAST_CARCINOMA_METAPLA
## 25                                                                  c2.PID_INTEGRIN5_PAT
## 26                                                       c2.LINDGREN_BLADDER_CANCER_HIGH_RECURF
## 27                                                               c2.POTTI_TOPOTECAN_SENSITI
## 28                                                    c2.MIYAGAWA_TARGETS_OF_EWSR1_ETS_FUSIONS_SI
## 29                                                   c2.BERTUCCI_MEDULLARY_VS_DUCTAL_BREAST_CANCER_SI
## 30                                                              c2.PASINI_SUZ12_TARGETS_SI
##     Correlation Metagenes
## 1        0.6490         2
## 2        0.6355         2
## 3        0.6178         2
## 4        0.6067       2,-3
```

```
## 5         0.6020          2
## 6         0.5990          2
## 7         0.5963          2
## 8         0.5953          2
## 9         0.5849          2
## 10        0.5829          2
## 11        0.5758          2
## 12        0.5587          2
## 13        0.5429          2
## 14        0.5409       2,-3
## 15        0.5396          2
## 16        0.5346          2
## 17        0.5312          2
## 18        0.5258          2
## 19        0.5242          2
## 20        0.5228          2
## 21        0.5171          2
## 22        0.5114          2
## 23        0.5081          2
## 24        0.5077          2
## 25        0.5074          2
## 26        0.5047          2
## 27        0.5017          2
## 28       -0.5087         -2
## 29       -0.5436         -2
## 30       -0.5916         -2
##
## [[2]]$c3
## data frame with 0 columns and 0 rows
##
## [[2]]$c4
##                        GeneSet Correlation Metagenes
## 1              c4.GNF2_PTX3        0.5533      2,-3
## 2              c4.MODULE_122       0.5369         2
## 3               c4.GNF2_MMP1       0.5366         2
## 4              c4.MODULE_562       0.5178         2
## 5 c4.MODULE_419/c4.MODULE_524   0.5128         2
## 6               c4.MODULE_47       0.5003         2
##
## [[2]]$c5
##                                                                                    GeneSet
## 1                                                                              c5.COLLAGEN
## 2 c5.PROTEINACEOUS_EXTRACELLULAR_MATRIX/c5.EXTRACELLULAR_MATRIX_PART/c5.EXTRACELLULAR_MATRIX
## 3                                                                      c5.BASEMENT_MEMBRANE
## 4                                                                    c5.SKELETAL_DEVELOPMENT
##   Correlation Metagenes
## 1      0.6496          2
## 2      0.5336          2
## 3      0.5148          2
## 4      0.5101          2
##
## [[2]]$c6
## data frame with 0 columns and 0 rows
##
```

```
## [[2]]$c7
## data frame with 0 columns and 0 rows
##
##
## [[3]]
## [[3]]$c1
## data frame with 0 columns and 0 rows
##
## [[3]]$c2
##                                                  GeneSet Correlation
## 1   c2.CHARAFE_BREAST_CANCER_LUMINAL_VS_MESENCHYMAL_SIGNED      0.5882
## 2              c2.REACTOME_GLYCEROPHOSPHOLIPID_BIOSYNTHESIS      0.5269
## 3              c2.SMID_BREAST_CANCER_RELAPSE_IN_BONE_SIGNED      0.5215
## 4                            c2.LIU_PROSTATE_CANCER_SIGNED      0.5202
## 5          c2.WAMUNYOKOLI_OVARIAN_CANCER_GRADES_1_2_SIGNED      0.5178
## 6   c2.WANG_BARRETTS_ESOPHAGUS_AND_ESOPHAGUS_CANCER_SIGNED      0.5175
## 7              c2.WAMUNYOKOLI_OVARIAN_CANCER_LMP_SIGNED      0.5165
## 8                       c2.WALLACE_PROSTATE_CANCER_SIGNED      0.5155
## 9                 c2.DOANE_BREAST_CANCER_CLASSES_SIGNED      0.5111
## 10              c2.WOO_LIVER_CANCER_RECURRENCE_SIGNED     -0.5000
## 11                                c2.PID_UPA_UPAR_PATHWAY     -0.5011
## 12          c2.SUZUKI_RESPONSE_TO_TSA_AND_DECITABINE_1A     -0.5141
## 13               c2.HUANG_DASATINIB_RESISTANCE_SIGNED     -0.5145
## 14              c2.LIM_MAMMARY_STEM_CELL_SIGNED     -0.5175
## 15                          c2.PID_INTEGRIN3_PATHWAY     -0.5175
## 16              c2.ROY_WOUND_BLOOD_VESSEL_SIGNED     -0.5235
## 17                          c2.PID_INTEGRIN1_PATHWAY     -0.5248
## 18   c2.LIEN_BREAST_CARCINOMA_METAPLASTIC_VS_DUCTAL_SIGNED     -0.6110
## 19       c2.VECCHI_GASTRIC_CANCER_ADVANCED_VS_EARLY_SIGNED     -0.6217
##      Metagenes
## 1            -3
## 2            -3
## 3            -3
## 4            -3
## 5            -3
## 6            -3
## 7            -3
## 8            -3
## 9            -3
## 10            3
## 11            3
## 12            3
## 13            3
## 14            3
## 15         -2,3
## 16            3
## 17         -2,3
## 18            3
## 19            3
##
## [[3]]$c3
## data frame with 0 columns and 0 rows
##
```

```
## [[3]]$c4
##                          GeneSet Correlation Metagenes
## 1 c4.MODULE_139/c4.MODULE_180        0.5195         -3
## 2                 c4.GNF2_PTX3      -0.5155       -2,3
##
## [[3]]$c5
## data frame with 0 columns and 0 rows
##
## [[3]]$c6
##                 GeneSet Correlation Metagenes
## 1 c6.LEF1_UP.V1_SIGNED      -0.5597          3
##
## [[3]]$c7
## data frame with 0 columns and 0 rows
##
##
## [[4]]
## [[4]]$c1
## data frame with 0 columns and 0 rows
##
## [[4]]$c2
##                              GeneSet Correlation Metagenes
## 1          c2.BERGER_MBD2_TARGETS        0.5646         -4
## 2 c2.TERAMOTO_OPN_TARGETS_CLUSTER_8     0.5274         -4
## 3    c2.LEE_LIVER_CANCER_MYC_SIGNED    -0.5203          4
##
## [[4]]$c3
##          GeneSet Correlation Metagenes
## 1 c3.V$HNF1_Q6      0.5124         -4
##
## [[4]]$c4
##                              GeneSet Correlation Metagenes
## 1 c4.GNF2_SERPINI2/c4.GNF2_SPINK1       0.6959         -4
##
## [[4]]$c5
##                          GeneSet Correlation Metagenes
## 1 c5.CARBOXYPEPTIDASE_ACTIVITY      0.5342         -4
##
## [[4]]$c6
## data frame with 0 columns and 0 rows
##
## [[4]]$c7
## data frame with 0 columns and 0 rows
##
##
## [[5]]
## [[5]]$c1
## data frame with 0 columns and 0 rows
##
## [[5]]$c2
##                                                     GeneSet
## 1                  c2.IVANOVA_HEMATOPOIESIS_LATE_PROGENITOR
## 2                      c2.MARSON_BOUND_BY_FOXP3_STIMULATED
```

```
## 3                                                      c2.SESTO_RESPONSE_TO_UV_C1
## 4                                      c2.NAKAYAMA_SOFT_TISSUE_TUMORS_PCA1_SIGNED
## 5                                          c2.IVANOVA_HEMATOPOIESIS_MATURE_CELL
## 6   c2.MARTINEZ_RB1_AND_TP53_TARGETS_SIGNED/c2.MARTINEZ_TP53_TARGETS_SIGNED
## 7                                          c2.KAMIKUBO_MYELOID_CEBPA_NETWORK
## 8                                          c2.MARSON_BOUND_BY_FOXP3_UNSTIMULATED
## 9                                                      c2.VALK_AML_CLUSTER_5
## 10                                     c2.RAGHAVACHARI_PLATELET_SPECIFIC_GENES
## 11                              c2.LIAN_LIPA_TARGETS_6M/c2.LIAN_LIPA_TARGETS_3M
## 12                                          c2.BROCKE_APOPTOSIS_REVERSED_BY_IL6
## 13                                     c2.SHEDDEN_LUNG_CANCER_GOOD_SURVIVAL_A5
##     Correlation Metagenes
## 1        0.6114         5
## 2        0.5798         5
## 3        0.5491         5
## 4        0.5413         5
## 5        0.5410         5
## 6        0.5304         5
## 7        0.5280         5
## 8        0.5154         5
## 9        0.5134         5
## 10       0.5124         5
## 11       0.5005         5
## 12       0.5001         5
## 13      -0.5080        -5
##
## [[5]]$c3
## data frame with 0 columns and 0 rows
##
## [[5]]$c4
##                      GeneSet Correlation Metagenes
## 1              c4.MODULE_86      0.5240         5
## 2 c4.MODULE_491/c4.MODULE_568      0.5063         5
##
## [[5]]$c5
## data frame with 0 columns and 0 rows
##
## [[5]]$c6
## data frame with 0 columns and 0 rows
##
## [[5]]$c7
##                                                                 GeneSet
## 1                               c7.GSE29618_MONOCYTE_VS_PDC_SIGNED
## 2              c7.GSE29618_MONOCYTE_VS_PDC_DAY7_FLU_VACCINE_SIGNED
## 3                                   c7.GSE3982_MAC_VS_NKCELL_SIGNED
## 4                              c7.GSE15767_MED_VS_SCS_MAC_LN_SIGNED
## 5               c7.GSE3982_BASOPHIL_VS_CENT_MEMORY_CD4_TCELL_SIGNED
## 6                       c7.GSE11057_PBMC_VS_MEM_CD4_TCELL_SIGNED
## 7   c7.GSE9006_HEALTHY_VS_TYPE_1_DIABETES_PBMC_4MONTH_POST_DX_SIGNED
## 8                  c7.GSE11057_NAIVE_CD4_VS_PBMC_CD4_TCELL_SIGNED
## 9                           c7.GSE3982_DC_VS_MAC_LPS_STIM_SIGNED
## 10               c7.GSE17721_POLYIC_VS_GARDIQUIMOD_1H_BMDM_SIGNED
## 11                                  c7.GSE29618_PDC_VS_MDC_SIGNED
## 12               c7.GSE6269_HEALTHY_VS_STREP_AUREUS_INF_PBMC_SIGNED
```

```
## 13                      c7.GSE20366_CD103_KLRG1_DP_VS_DN_TREG_SIGNED
## 14                        c7.GSE22886_NAIVE_CD8_TCELL_VS_DC_SIGNED
## 15                   c7.GSE22886_NAIVE_CD8_TCELL_VS_MONOCYTE_SIGNED
## 16                       c7.GSE11057_CD4_CENT_MEM_VS_PBMC_SIGNED
## 17                        c7.GSE22886_NAIVE_CD4_TCELL_VS_DC_SIGNED
## 18                   c7.GSE11057_NAIVE_VS_MEMORY_CD4_TCELL_SIGNED
## 19                          c7.GSE11057_CD4_EFF_MEM_VS_PBMC_SIGNED
## 20                             c7.GSE10325_BCELL_VS_MYELOID_SIGNED
## 21                      c7.GSE22886_NAIVE_TCELL_VS_MONOCYTE_SIGNED
## 22            c7.GSE10325_LUPUS_CD4_TCELL_VS_LUPUS_MYELOID_SIGNED
## 23                   c7.GSE22886_NAIVE_CD4_TCELL_VS_MONOCYTE_SIGNED
##    Correlation Metagenes
## 1       0.5760         5
## 2       0.5712         5
## 3       0.5573         5
## 4       0.5502         5
## 5       0.5352         5
## 6       0.5314         5
## 7       0.5209         5
## 8      -0.5042        -5
## 9      -0.5042        -5
## 10     -0.5076        -5
## 11     -0.5086        -5
## 12     -0.5086        -5
## 13     -0.5233        -5
## 14     -0.5267        -5
## 15     -0.5274        -5
## 16     -0.5352        -5
## 17     -0.5355        -5
## 18     -0.5379        -5
## 19     -0.5420        -5
## 20     -0.5519        -5
## 21     -0.5610        -5
## 22     -0.5699        -5
## 23     -0.5825        -5
##
##
## [[6]]
## [[6]]$c1
## data frame with 0 columns and 0 rows
##
## [[6]]$c2
##            GeneSet Correlation Metagenes
## 1 c2.LEI_MYB_TARGETS       0.509         6
##
## [[6]]$c3
## data frame with 0 columns and 0 rows
##
## [[6]]$c4
##                    GeneSet Correlation Metagenes
## 1 c4.GNF2_CDH3/c4.GNF2_SERPINB5      0.5532         6
##
## [[6]]$c5
## data frame with 0 columns and 0 rows
```

```
##
## [[6]]$c6
## data frame with 0 columns and 0 rows
##
## [[6]]$c7
## data frame with 0 columns and 0 rows
```

# 5 Session information

```
session_info

## R version 3.1.1 (2014-07-10)
## Platform: x86_64-unknown-linux-gnu (64-bit)
##
## locale:
##  [1] LC_CTYPE=en_AU.UTF-8         LC_NUMERIC=C
##  [3] LC_TIME=en_AU.UTF-8          LC_COLLATE=en_AU.UTF-8
##  [5] LC_MONETARY=en_AU.UTF-8      LC_MESSAGES=en_AU.UTF-8
##  [7] LC_PAPER=en_AU.UTF-8         LC_NAME=en_AU.UTF-8
##  [9] LC_ADDRESS=en_AU.UTF-8       LC_TELEPHONE=en_AU.UTF-8
## [11] LC_MEASUREMENT=en_AU.UTF-8   LC_IDENTIFICATION=en_AU.UTF-8
##
## attached base packages:
## [1] splines   parallel  methods   stats     graphics  grDevices utils
## [8] datasets  base
##
## other attached packages:
##  [1] doParallel_1.0.8   iterators_1.0.7    foreach_1.4.2
##  [4] ahaz_1.14          survival_2.37-7    NMF_0.20.5
##  [7] Biobase_2.26.0     BiocGenerics_0.12.1 cluster_1.15.3
## [10] rngtools_1.2.4     pkgmaker_0.22      registry_0.2
## [13] energy_1.6.2       glmnet_1.9-8       Matrix_1.1-4
## [16] glmulti_1.0.7      rJava_0.9-6
##
## loaded via a namespace (and not attached):
##  [1] boot_1.3-13        codetools_0.2-9    colorspace_1.2-4
##  [4] compiler_3.1.1     digest_0.6.4       ggplot2_1.0.0
##  [7] grid_3.1.1         gridBase_0.4-7     gtable_0.1.2
## [10] lattice_0.20-29    MASS_7.3-35        munsell_0.4.2
## [13] plyr_1.8.1         proto_0.3-10       RColorBrewer_1.0-5
## [16] Rcpp_0.11.3        reshape2_1.4       scales_0.2.4
## [19] stringr_0.6.2      tools_3.1.1        xtable_1.7-4

sessionInfo()

## R version 3.1.1 (2014-07-10)
## Platform: x86_64-unknown-linux-gnu (64-bit)
##
## locale:
##  [1] LC_CTYPE=en_AU.UTF-8         LC_NUMERIC=C
##  [3] LC_TIME=en_AU.UTF-8          LC_COLLATE=en_AU.UTF-8
##  [5] LC_MONETARY=en_AU.UTF-8      LC_MESSAGES=en_AU.UTF-8
```

```
##  [7] LC_PAPER=en_AU.UTF-8          LC_NAME=en_AU.UTF-8
##  [9] LC_ADDRESS=en_AU.UTF-8         LC_TELEPHONE=en_AU.UTF-8
## [11] LC_MEASUREMENT=en_AU.UTF-8     LC_IDENTIFICATION=en_AU.UTF-8
##
## attached base packages:
## [1] parallel  methods   splines   stats     graphics  grDevices utils
## [8] datasets  base
##
## other attached packages:
##  [1] stargazer_5.1      xtable_1.7-4       gplots_2.14.2
##  [4] RColorBrewer_1.0-5 glmnet_1.9-8       Matrix_1.1-4
##  [7] glmulti_1.0.7      rJava_0.9-6        nnls_1.4
## [10] NMF_0.20.5         Biobase_2.26.0     BiocGenerics_0.12.1
## [13] cluster_1.15.3     rngtools_1.2.4     pkgmaker_0.22
## [16] registry_0.2       energy_1.6.2       survival_2.37-7
## [19] knitr_1.8
##
## loaded via a namespace (and not attached):
##  [1] bitops_1.0-6      boot_1.3-13       caTools_1.17.1
##  [4] codetools_0.2-9   colorspace_1.2-4  digest_0.6.4
##  [7] doParallel_1.0.8  evaluate_0.5.5    foreach_1.4.2
## [10] formatR_1.0       gdata_2.13.3      ggplot2_1.0.0
## [13] grid_3.1.1        gridBase_0.4-7    gtable_0.1.2
## [16] gtools_3.4.1      highr_0.4         iterators_1.0.7
## [19] KernSmooth_2.23-13 labeling_0.3     lattice_0.20-29
## [22] MASS_7.3-35       munsell_0.4.2     plyr_1.8.1
## [25] proto_0.3-10      Rcpp_0.11.3       reshape2_1.4
## [28] scales_0.2.4      stringr_0.6.2     tools_3.1.1
```