

CMScaller: an R package for tumor microenvironment-independent colorectal cancer consensus molecular subtyping

15 July 2021

Peter W. Eide^{1,2}, Seyed H. Moosavi^{1,2,3}, Ragnhild A. Lothe^{1,2,3}, and Anita Sveen^{1,2,3}

¹ Department of Molecular Oncology, Institute for Cancer Research; ² K.G.Jebsen Colorectal Cancer Research Centre, Oslo University Hospital, Oslo, NO-0424, Norway; and ³ Institute for Clinical Medicine, University of Oslo, Oslo, N-0318, Norway

- contact: anita.sveen@rr-research.no (mailto:anita.sveen@rr-research.no)
- Date: 2021-07-15
- package: CMScaller 2.0.1

1 Introduction

Colorectal cancers (CRCs) can be divided into four gene expression-based, biologically distinct consensus molecular subtypes (CMS (<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4636487/>)) [1]. This classification provides prognostic stratification of the patients and presents a potential basis for stratified treatment. The original CMS classifier was developed from primary tumor tissue samples and is dependent on gene expression signals from the immune and stromal compartments. It often fails to identify the poor-prognostic CMS4 mesenchymal group in cancer cell cultures (immortalized cell lines and patient-derived organoids) and metastatic tumor samples. CMScaller uses cancer cell-intrinsic, subtype-specific gene expression markers as features for classification of diverse CRC sample types. For pre-clinical models, nearest template prediction algorithm is employed [2, 3]. From version 2, CMScaller also includes a random forest model for CMS classification of CRC liver metastases [4, 5].

2 Input data

For pre-clinical models, CMScaller provides robust *cross platform and sample-type* performance given a balanced, homogeneous dataset of >40 unique samples [3, 6]. For less than ~40 samples, sampling variance (by-chance subtype depletion/enrichment) is a concern. Similarly, selection, e.g. excluding microsatellite instable (MSI) samples or including only aggressive cancers, would break an underlying assumption and bias the resulting predictions [7]. For CRC liver metastases, background signals from the liver tumor microenvironment is an additional concern, and non-malignant liver tissue samples was used for “background adjustment” in the initial CMS classification of metastases. However, the random forest model implemented in CMScaller v2.0.1 is not dependent on background adjustment, and can be applied directly to gene expression profiles from liver metastases.

3 Quick start

3.1 Installation and dependencies

The following packages are required in order to run examples in this vignette.

- Bioconductor [8]: Biobase (<https://bioconductor.org/packages/3.13/Biobase>), limma (<https://bioconductor.org/packages/3.13/limma>)
- CRAN randomForest (<https://CRAN.R-project.org/package=randomForest>), survival (<https://CRAN.R-project.org/package=survival>)

In addition, edgeR (<https://bioconductor.org/packages/3.13/edgeR>) is needed for specific RNA-sequencing normalization methods and snow (<https://CRAN.R-project.org/package=snow>) for ntp parallelization on Windows systems.

```
# dependencies: run if not already installed (R >= 3.5)
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install(c("Biobase", "limma"))

# install dependencies with biocLite (R < 3.5)
# source("https://bioconductor.org/biocLite.R")
# biocLite(c("Biobase", "limma"))

# install package
devtools::install_github("Lothelab/CMScaller")
```

3.2 CMS classification

CMScaller (for pre-clinical models) and lmCMScaller (for liver metastases) functions both require an expression matrix or ExpressionSet as input (emat). Gene names in rownames(emat) must be NCBI Entrez (<https://www.ncbi.nlm.nih.gov/gene>), Ensembl (<http://www.ensembl.org/index.html>) or HGNC symbol (<http://www.genenames.org/>) identifiers. For gene symbols or Ensembl identifiers, parameter rowNames must be set to symbol or ensq, respectively. The code chunk below demonstrates how to perform classification using TCGA primary colorectal cancer example data [9].

- microarray data input should be pre-processed and normalized (often \log_2 transformed)[\wedge log].
- RNA-sequencing counts/RSEM values could be used directly by setting RNAseq=TRUE which activates quantile normalization and \log_2 transformation.

Hoshida [2] does not explicitly state whether input should be \log_2 transformed or not and examples in paper include both. Such transformation reduces the weight of genes with large deviations and will affect results at the margins.

3.2.1 Pre-clinical models

```
library(Biobase) # if input is ExpressionSet
library(CMScaller)
## CMScaller v2.0.1; genome annotation: GENCODE v32/GRCh38.p13
# get RNA-seq counts from TCGA example data
counts <- exprs(crcTCGAsubset)
head(counts[,1:2])
##          TCGA-4N-A93T-01A-11R TCGA-4T-AA8H-01A-11R
## 1                      354                      28
## 29974                   208                      238
## 54715                    21                       1
## 87769                    767                     404
## 144568                    1                      50
## 2                      3088                     1823
# prediction and gene set analysis
par(mfrow=c(1,2))
res <- CMScaller(emat=counts, RNAseq=TRUE, FDR=0.05)
## performing log2-transform and quantile normalization...
## cosine correlation distance
## 92 samples; 4 classes; 82-237 features/class
## serial processing; 1000 permutation(s)...
## predicted samples/class (FDR<0.05)
##
## CMS1 CMS2 CMS3 CMS4 <NA>
## 15  23  15  28  11
## 11/92 samples set to NA
cam <- CMSgsa(emat=counts, class=res$prediction, RNAseq=TRUE)
## 11 samples with class or batch NA's excluded
```

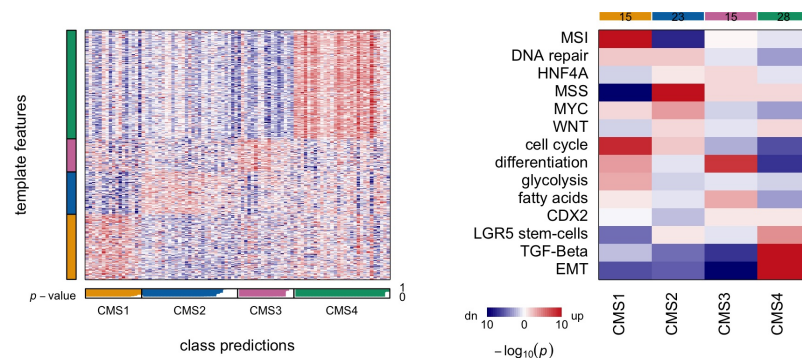


Figure 1: **CMScaller graphic output**

Left heatmap shows relative expression for template genes. Samples (columns) are ordered according to class predictions and confidence. The height of the white bars below gives the unadjusted prediction p -values. Genes (rows) are ordered according to class template. Heatmap color saturation indicates magnitude while red and blue indicate genes up and down relative to the sample mean. Right heatmap shows results for Camera gene set analysis. Heatmap color saturation indicates statistical significance and red and blue indicates direction of change.

```
# comparison with true class
table(pred=res$prediction, true=crcTCGAsubset$CMS)
##      true
## pred  CMS1 CMS2 CMS3 CMS4
## CMS1    8    0    3    0
## CMS2    0   20    0    0
## CMS3    0    1    6    0
## CMS4    0    2    0   19
head(res, n=3)
##      prediction d.CMS1 d.CMS2 d.CMS3 d.CMS4 p.value  FDR
## TCGA-4N-A93T-01A-11R <NA>  0.76  0.70  0.72  0.89  0.927 0.9477
## TCGA-4T-AA8H-01A-11R CMS3   0.76  0.69  0.63  0.90  0.001 0.0013
## TCGA-5M-AAT4-01A-11R <NA>  0.74  0.69  0.73  0.83  0.776 0.8115
```

- rownames(res) equals colnames(emat)
- class predictions with NA for samples with adjusted-*p*-value > threshold
- templates distances
- prediction *p*-values¹
- prediction FDR-adjusted *p*-values

¹ lowest possible estimate of the *p*-value is 1/permutations

3.2.2 Colorectal cancer liver metastases

The package includes a *random forest* [4, 10] model intended for CMS classification of colorectal cancer liver metastases. The model was trained on Affymetrix Human Transcriptome array 2.0 data.

mcrc0SL0subset serves as example data and is an ExpressionSet-class object holding normalized, log₂-transformed gene expression measurements for 1373 genes × 295 samples (176 patients). Clinical and molecular annotations are available using command Biobase::pData(mcrc0SL0subset).

```
data("mcrc0SL0subset")
dim(mcrc0SL0subset)
## Features Samples
## 1373 295
par(mfrow=c(1,2))
res <- lmCMScaller(mcrc0SL0subset, posterior=.6)
# comparison with PAM model presented in Eide et al. Metastatic heterogeneity..
.
table(RF=res$prediction, PAM=mcrc0SL0subset$cms_pam, useNA="alw")
##      PAM
## RF    CMS1 CMS2 CMS3 CMS4 <NA>
## CMS1  20    0    0    0   10
## CMS2   0   93    0    0   26
## CMS3   0    0    5    0   12
## CMS4   0    0    0   79    1
## <NA>   0   10    0   23   16
cam <- CMSgsa(emat=mcrc0SL0subset, class=res$prediction,
              keepN=!duplicated(mcrc0SL0subset$`Patient ID`), returnMatrix=TRUE)
## 32 samples with class or batch NA's excluded

subPCA(mcrc0SL0subset, res$prediction)
```

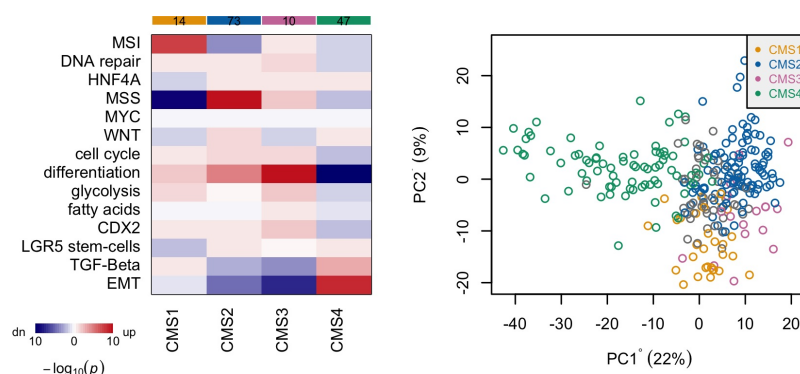


Figure 2: **CMScaller graphic output**

Left heatmap shows results for Camera gene set analysis. Heatmap color saturation indicates statistical significance and red and blue indicates direction of change. Right scatterplot shows the first two principal components from analysis based on CMS-classification genes with samples colored according to CMS.

```
head(res)
##          CMS1 CMS2 CMS3 CMS4 prediction posterior
## 001_T1 0.054 0.472 0.052 0.422      <NA>      0.47
## 002_T1 0.038 0.730 0.168 0.064      CMS2      0.73
## 003_T1 0.094 0.038 0.864 0.004      CMS3      0.86
## 003_T2 0.130 0.044 0.806 0.020      CMS3      0.81
## 004_T3 0.030 0.074 0.886 0.010      CMS3      0.89
## 004_T4 0.078 0.070 0.840 0.012      CMS3      0.84
```

lmCMScaller takes a normalized gene expression matrix or ExpressionSet object and returns a data.frame with samples in rows and columns of class-wise posterior probabilities, classification and the corresponding (max) posterior probabilities. Since there is more than one sample per patient, heterogeneity can be assessed:

```
cms_by_patient <- split(res, mrcrOSLOsubset$`Patient ID`)
cms_by_patient <- cms_by_patient[sapply(cms_by_patient, nrow)>1]
cms_by_patient <- t(sapply(cms_by_patient, function(x)
  table(x$prediction, useNA="alw"))))
head(cms_by_patient) # row names indicate patient number
##      CMS1 CMS2 CMS3 CMS4 <NA>
## 3      0    0    2    0    0
## 4      0    0    2    0    0
## 9      0    1    0    1    0
## 16     0    1    0    0    3
## 25     0    1    0    1    0
## 27     0    2    0    0    0
```

For instance, for patients 3 and 4, both samples are classified as CMS3 and these cancers are therefore homogeneous in terms of CMS. Patients 9 and 25 have cancers that are heterogeneous. Below, Kaplan-Meier analyses are included for illustration purposes. Patients are stratified according to CMS for a single random-sample and "worst-case" base on the following rule: if any lesion CMS1/CMS3, then patient CMS1/3 (combining the two subtypes with shortest median survival); if no CMS1/3 and any lesion CMS4, then patient CMS4; otherwise CMS2.

```
par(mfrow=c(1,2), bty="n")
library(survival)

# define worst-case subtype
worstCMSfun <- function(cms, patient) {
  if (length(cms) != length(patient)) stop ("cms and patient differ in length")
  # split by patient and determine worst CMS
  cms <- lapply(split(cms, patient), table)
  cms <- sapply(cms, function(x) {
    ifelse(x[1]>0|x[3]>0, "CMS1.3",
           ifelse(x[4]>0, "CMS4",
                  ifelse(x[2]>0, "CMS2", NA)))
  })
  names(cms) <- gsub("\\.CMS.*", "", names(cms))
  return(factor(cms[match(patient, names(cms))],
               levels=c("CMS1.3", "CMS2", "CMS4")))
}

# define dataset, here we use *nearest*, not necessarily confident CMS
res <- lmCMScaller(mrcrOSLOsubset, posterior=0)
df <- data.frame(
  CMS=res$prediction,
  worstCMS=worstCMSfun(res$prediction, mrcrOSLOsubset$`Patient ID`),
  patient=mrcrOSLOsubset$`Patient ID`,
  sens=mrcrOSLOsubset$`60 months overall survival, status`,
  time=mrcrOSLOsubset$`60 months overall survival, time`
)

# remove patient duplicates
df <- df[!is.na(df$time),]
df <- df[!duplicated(df$patient),]
df$Surv <- survival::Surv(time=df$time, event=df$sens)

# CMS and worst-case CMS
plot(survival::survfit(Surv~CMS, data=df), col=CMScaller:::subData$classCol,
     xlab="time (months)", ylab="overall survival (proportion)")
plot(survival::survfit(Surv~worstCMS, data=df),
     col=CMScaller:::subData$classCol[-3],
     xlab="time (months)")
```

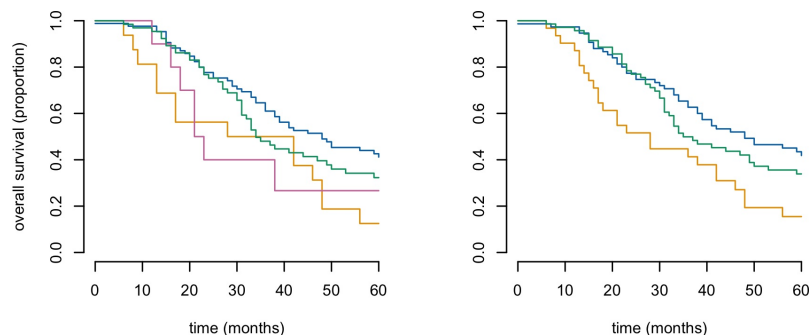


Figure 3: **Kaplan-Meier survival analyses for CMS (left) and worst-CMS stratifications (right)**

4 Package details

CMScaller provides a wrapper function for `ntp` in R. `lmCMScaller` calls the random forest function. CMSgsa provides some presets for `subCamera`.

4.1 Preparing custom templates

Templates consist of sets of subtype-specific marker genes. CMScaller and `lmCMScaller` provide build-in templates for CMS classification. However, below, is an example on how to prepare custom templates based on a training set with known class labels. `subDEG` performs *limma* differential expression analysis for identification of such markers. `doVoom=TRUE` enables *voom* transformation - required for proper *limma* modeling of RNA-sequencing counts [11].

```
emat <- crcTCGAsubset
cms <- emat$CMS.Syn
train <- sample(seq_along(cms), size=length(cms)/(2))
deg <- subDEG(emat[,train], class=cms[train], doVoom=TRUE)
## 14 samples with class or batch NA's excluded
templates <- ntpMakeTemplates(deg, resDEG=TRUE, topN=50)
templates$symbol <- fromTo(templates$probe)
tail(templates,n=3)
##      probe class symbol
## 198  9444  CMS4    QKI
## 199 1009   CMS4  CDH11
## 200 23768  CMS4  FLRT2
```

4.2 Gene Set Analysis

`subCamera` provides gene set analysis and visualization and is a wrapper functions for `camera` in the *limma* (<https://bioconductor.org/packages/3.13/limma>) package. `camera` controls for intra-set gene-wise correlations in order to reduce false-positive rate while retaining statistical power [12, 13]. CMSgsa provides preset gene sets to `subCamera`.

```
# increase left margins to accommodate gene set names
par.old <- par()
par(mfrow=c(1,1), mar=par.old$mar+c(0,4,0,0))
subCamera(counts, cms, geneList=geneSets.CRC, doVoom=TRUE)
## 29 samples with class or batch NA's excluded
```

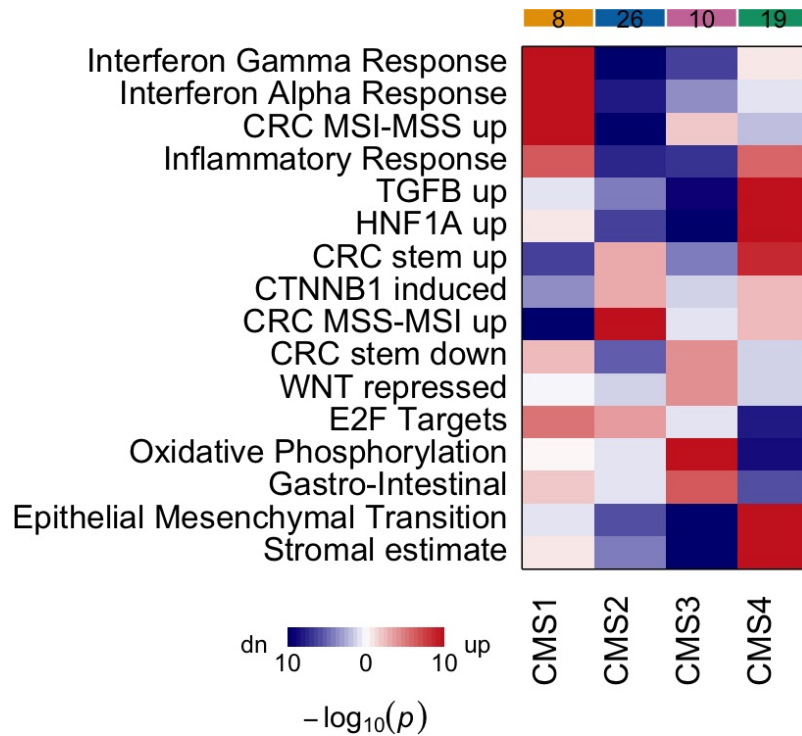


Figure 4: **Gene Set Analysis (GSA)** shows that CMS are biologically distinct

```
# restore margins
par(mar=par.old$mar)
```

4.3 PCA Analysis

subPCA and plotPC provide convenient wrapper functions for prcomp. subPCA perform PCA on the input data and plot the resulting low-dimensional projection with samples colored according to either a continuous covariate (e.g. expression of gene of interest) or group (such as CMS). plotPC visualizes the most important variables.

```
# increase left margins to accommodate gene set names
par(mfrow=c(1,2))
p <- subPCA(emat = crcTCGAsubset, class = crcTCGAsubset$CMS.Syn,
            normMethod = "quantile", pch=16, frame=FALSE)
plotPC(p, n=6, entrez=TRUE)
```

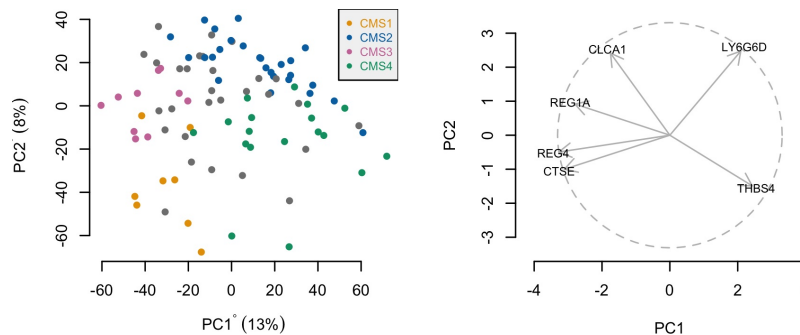


Figure 5: **Principal component analysis (PCA) and CMS**
First two principal components separates CMS (left) with CMS4 characterized by high levels of THBS4 and low levels of CLCA1 (right).

4.4 Nearest Template Prediction

ntp matches templates\$probe against rownames(emat). Missing features and features with NA/NaN's are ignored in the prediction. emat should be row-wise centered and scaled.

```
# loads included emat, scales and centers
emat <- crcTCGAsubset
emat_sc <- ematAdjust(emat, normMethod="quantile")
head(emat_sc[,1:2])
##          TCGA-4N-A93T-01A-11R  TCGA-4T-AA8H-01A-11R
## 1                4.2710                1.0117
## 29974             -0.0033                0.5152
## 54715              1.2946                0.0041
## 87769              1.9697                1.5102
## 144568             0.2288                2.4779
## 2               -1.5602               -1.6659
```

ntp function requires an expression matrix and templates. Since prediction confidence is estimated from permutations, strict p -value reproducibility requires `set.seed`.

```
# test set prediction
res <- ntp(emat_sc[, -train], templates, nPerm=1000)
res <- subSetNA(res, pValue=.1)
## 3/46 samples set to NA
table(pred=res$prediction, true=cms[-train])
##           true
## pred   CMS1 CMS2 CMS3 CMS4
## CMS1    4    0    1    0
## CMS2    0   10    0    0
## CMS3    1    1    3    0
## CMS4    0    0    0   10
head(res)
##           prediction d.CMS1 d.CMS2 d.CMS3 d.CMS4 p.value  FDR
## TCGA-4N-A93T-01A-11R   CMS3  0.75  0.67  0.62  0.88  0.003 0.0038
## TCGA-5M-AAT4-01A-11R   CMS2  0.80  0.63  0.70  0.85  0.035 0.0402
## TCGA-A6-5657-01A-01R   CMS4  0.69  0.80  0.74  0.50  0.001 0.0014
## TCGA-A6-5667-01A-21R   CMS2  0.78  0.52  0.84  0.63  0.001 0.0014
## TCGA-A6-6780-01A-11R   CMS1  0.45  0.78  0.72  0.81  0.001 0.0014
## TCGA-A6-6782-01A-11R   CMS4  0.73  0.65  0.82  0.51  0.001 0.0014
```

ntp output is a `data.frame` with $3 + K$ columns where K is the number of classes. Rows represent columns in input `emat`.

- `rownames(res)` equals `colnames(emat)`
- class predictions with `levels(res$prediction)` equaling `levels(templates$class)`
- templates distances (defaults to cosine correlation distance)
- prediction p -values
- prediction FDR-adjusted p -values

`subSetNA` function resets predictions with p -value or FDR above some arbitrary threshold to `NA`.

5 Nearest Template Prediction

Nearest template prediction (NTP) was proposed as a classification algorithm by Dr. Yujin Hoshida and published in *PLoS ONE* in 2010 [2]. It aims to provide robust single-sample class prediction for high-dimensional, noisy gene expression data. In brief, first, for each subclass, a *template*, a list of genes coherently upregulated is determined. Then, for each sample, the *distance* to each template is calculated and class is assigned based on the smallest distance. Finally, prediction confidence is assessed based on the distance of the null-distribution, estimated from *permutation tests* (feature permutation). The default distance metric selected by Hoshida was a cosine similarity-derived distance (see below). When applied to a reasonably well-balanced homogeneous dataset, row-wise centering and scaling is performed (gene means and standard deviations $\mu = 0$, $\sigma = 1$). In case of single-sample prediction, feature-wise means and standard deviations from a previous sample set are used to perform sample-wise scaling and centering. The key advantages of the NTP algorithm are conceptual simplicity, biological plausibility, ease of implementation and robustness.

Formally, N samples with expression values for P genes divided into K different classes.

- $\mathbf{X}_{[P,N]}$ centered and scaled expression matrix where column vector $x_{[P]}$ is the expression for sample n .
- M is a list of K vectors where each element m is a set of marker features with higher expression in samples belonging to class k as compared to remaining samples. m 's may be of uneven length, but are typically $\ll P$
- $\mathbf{Y}_{[P,K]}$ template matrix where $y_{[P]} = [p \in m]$ for class k (0 if not marker, 1 otherwise).

For the sample and template vectors x and y , a proper distance metric, $d_{x,y}$ for the similarity function $f(x,y)$ is given by $d = \sqrt{\frac{1}{2}(1 - f(x,y))}$ [14]. Here f is either cosine, Kendall,

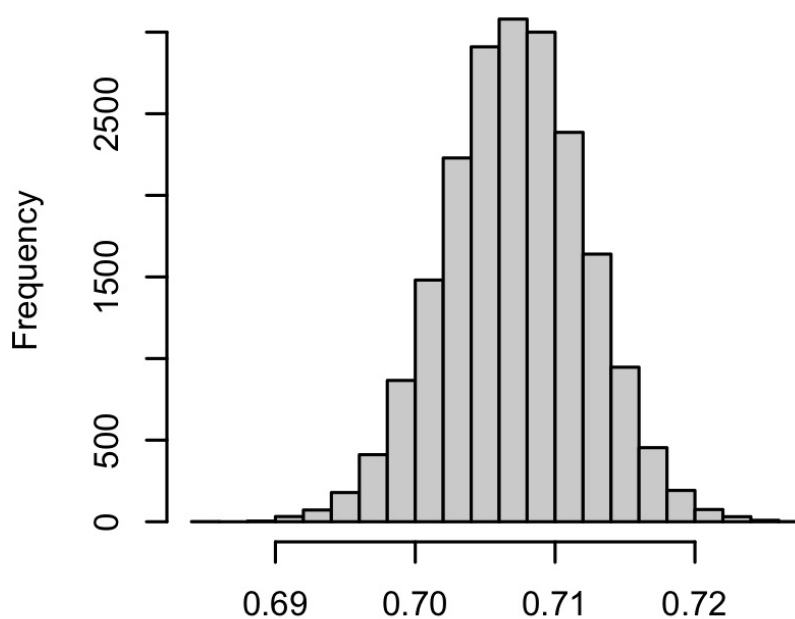
Pearson or Spearman correlation. Cosine similarity, the angle between two Euclidean vectors is given by

$$f(x,y) = \cos(\theta) = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

The following code chunks demonstrate NTP in code.

```
# random centered/scaled expression matrix and templates
set.seed(42)
N <- 5000; P <- 5000; K <- 4; nPerm <- 1000; n <- 1
X <- matrix(rnorm(P*N, mean=0, sd=1), ncol=N)
Y <- matrix(rbinom(P*K, size=1, prob=.01), ncol=K)
# sample-template correlations (implemented in corCosine)
cos.sim <- crossprod(X,Y) / outer(
  sqrt(apply(X, 2, crossprod)),
  sqrt(apply(Y, 2, crossprod)))
# sample-template distances (vectorized)
simToDist <- function(cos.sim) sqrt(1/2 * (1-cos.sim))
cos.dist <- simToDist(cos.sim)
hist(cos.dist, xlab="cosine correlation distance")
```

Histogram of cos.dist



cosine correlation distance

For centered, scaled and uncorrelated data, the cosine correlation distance is

$$\sqrt{0.5 \times (1 - 0)} \approx 0.707$$

Resulting distances are ranked among distances of permuted samples and used to estimate prediction confidence. The lowest possible p -value estimate is therefore $1/\text{permutations}$

```
# estimate prediction confidence
pred.class <- apply(cos.dist, 1, which.min)
pred.dists <- apply(cos.dist, 1, min)
null.dist <- replicate(nPerm, min(simToDist(corCosine(sample(X[,n]), Y))))
p <- rank(c(pred.dists[n], null.dist))[1]/(length(null.dist))
```

Code and plot below illustrate the uniform p -value distribution for centered and scaled uncorrelated input².

```
# rearrange matrix and templates for ntp input
rownames(X) <- make.names(seq_len(P))
templates <- lapply(seq_len(K), function(k) rownames(X)[Y[,k]==1])
names(templates) <- paste("k", seq_len(K))
templates <- ntpMakeTemplates(templates, resDEG=FALSE)
# permutations set to 100 to reduce processing for vignette
par(mfrow=c(1,2))
res <- ntp(X, templates, nCores=1L, nPerm=100, doPlot=TRUE)
# expect uniform distribution
hist(res$p.value, main="", xlab="prediction p-values")
```

² present NTP implementation provides more conservative p -value estimates than Hoshida[2].

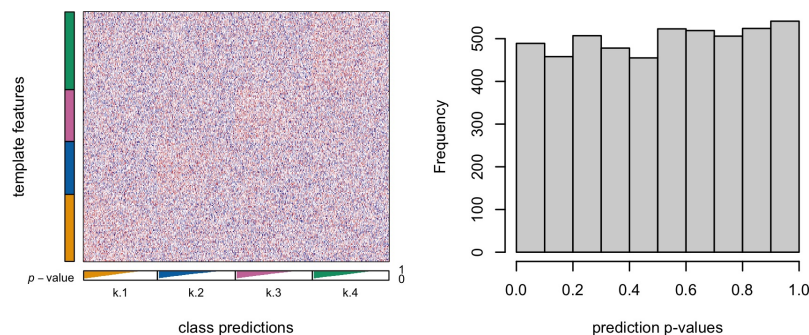


Figure 6: **NTP results for random data**
 Left heatmap shows expression for template genes for random data with rows and columns sorted according to class. Right histogram shows the expected uniform p -value distribution for the random data.

6 Notes

- Default qualitative color palette is now from [15].
- See `news(package="CMScaller")` for additional details.

7 Session

```
toLatex(sessionInfo())
```

References

1. Guinney J, Dienstmann R, Wang X, Reyniès A de, Schlicker A, Soneson C, Marisa L, Roepman P, Nyamundanda G, Angelino P, Bot BM, Morris JS, Simon IM, Gerster S, Fessler E, Melo FDSE, Missiaglia E, Ramay H, Barras D, Homicsko K, Maru D, Manyam GC, Broom B, Boige V, Perez-Villamil B, Laderas T, Salazar R, Gray JW, Hanahan D, Tabernero J, et al.: **The consensus molecular subtypes of colorectal cancer.** *Nat Med* 2015, **21**:1350-1356.
2. Hoshida Y: **Nearest Template Prediction: A Single-Sample-Based Flexible Class Prediction with Confidence Assessment.** *PLoS One* 2010, **5**:e15543.
3. Eide PW, Bruun J, Lothe RA, Svein A: **CMScaller: An r package for consensus molecular subtyping of colorectal cancer pre-clinical models.** *Sci Rep* 2017, **7**:16618.
4. Breiman L: **Random forests.** *Mach Learn* 2001, **45**:5-32.
5. Eide PW, Moosavi SH, Eilertsen IA, Brunzell TH, Langerud J, Berg KCG, Røskok BI, Bjørnbeth BA, Nesbakken A, Lothe RA, Svein A: **Metastatic heterogeneity of the consensus molecular subtypes of colorectal cancer.** *npj Genomic Medicine* 2021, **6**.
6. Svein A, Bruun J, Eide PW, Eilertsen IA, Ramirez L, Murumägi A, Arjama M, Danielsen SA, Kryeziu K, Elez E, Tabernero J, Guinney J, Palmer HG, Nesbakken A, Kallioniemi O, Dienstmann R, Lothe RA: **Colorectal cancer consensus molecular subtypes translated to preclinical models uncover potentially targetable cancer cell dependencies.** *Clin Cancer Res* 2018, **24**:794-806.
7. Zhao X, Rødland EA, Tibshirani R, Plevritis S: **Molecular subtyping for clinically defined breast cancer subgroups.** *Breast Cancer Res* 2015, **17**.
8. Huber W, Carey VJ, Gentleman R, Anders S, Carlson M, Carvalho BS, Bravo HC, Davis S, Gatto L, Girke T, Gottardo R, Hahne F, Hansen KD, Irizarry RA, Lawrence M, Love MI, MacDonald J, Obenchain V, Oleś AK, Pagès H, Reyes A, Shannon P, Smyth GK, Tenenbaum D, Waldron L, Morgan M: **Orchestrating high-throughput genomic analysis with Bioconductor.** *Nat Meth* 2015, **12**:115-121.
9. TCGA: **Comprehensive molecular characterization of human colon and rectal cancer.** *Nature* 2012, **487**:330-337.
10. Liaw A, Wiener M: **Classification and regression by randomForest.** *R News* 2002, **2**:18-22.
11. Law CW, Chen Y, Shi W, Smyth GK: **Voom: Precision weights unlock linear model analysis tools for RNA-seq read counts.** *Genome Biology* 2014, **15**:R29.
12. Wu D, Smyth GK: **Camera: A competitive gene set test accounting for inter-gene correlation.** *Nucl Acids Res* 2012, **40**:e133.
13. Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK: **Limma powers differential expression analyses for RNA-sequencing and microarray studies.** *Nucl Acids Res* 2015, **43**:e47.
14. Dongen S van, Enright AJ: **Metric distances derived from cosine similarity and Pearson and Spearman correlations.** *arXiv:12083145 [cs, stat]* 2012.
15. **Color universal design (CUD) - - how to make figures and presentations that are friendly to colorblind people** [<http://jfly.iam.u-tokyo.ac.jp/color/> (<http://jfly.iam.u-tokyo.ac.jp/color/>)]