

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

Benedikt Rauscher

2018-07-04

Abstract

Tumor heterogeneity is a major challenge to the treatment of colorectal cancer (CRC). Recently, a transcriptome-based classification was developed, segregating CRC into four consensus molecular subtypes (CMS) with distinct biological and clinical characteristics. Here, we applied the CMS classification on CRC cell lines to identify novel subtype-specific drug vulnerabilities. We combined publicly available transcriptome data from multiple resources to assign 159 CRC cell lines to CMS. By integrating results from large scale drug screens, we discovered that CMS1 cancer is highly vulnerable to the survivin suppressor YM155. We confirmed our results using an independent panel of CRC cell lines and demonstrate a 100-fold higher sensitivity of CMS1 lines. This vulnerability was specific to YM155 and not observed for commonly used chemotherapeutic agents. In CMS1 cancer, low concentrations of YM155 induced apoptosis and expression signatures associated with NFkappaB and ER stress mediated apoptosis signaling. Using a genome-wide CRISPR/Cas9 screen, we further discovered a novel role of genes involved in LDL-receptor recycling as modulators of YM155 response in CMS1 CRC. Our work shows that combining drug response data with CMS classification in cell lines can reveal specific vulnerabilities and propose YM155 as novel CMS1 specific drug.

Contents

| | | |
|-----|---|---|
| 1 | About. | 3 |
| 2 | Dependencies. | 3 |
| 3 | Meta data | 3 |
| 4 | Preprocessing microarray data | 4 |
| 4.1 | Adai cell line | 4 |
| 4.2 | Wagner cell line | 5 |
| 4.3 | Garnett cell line. | 6 |
| 4.4 | Barettina cell line | 7 |
| 4.5 | Medico cell line. | 8 |
| 5 | Aggregation of subtype data | 9 |

| | | |
|------|--|----|
| 6 | Comparison with other studies | 11 |
| 6.1 | Sveen et al | 11 |
| 6.2 | Linnekamp et al | 12 |
| 6.3 | Comparison | 12 |
| 7 | Characterization of predicted subtypes | 14 |
| 8 | CMS in PDX models | 17 |
| 9 | CMS-dependent drug response driven from public drug screening data | 18 |
| 9.1 | Identification of candidate substances | 18 |
| 9.2 | Validation of YM-155 | 21 |
| 9.3 | Validation of 5-FU | 22 |
| 9.4 | Validation of SN38 (Irinotecan) | 23 |
| 9.5 | Validation of cell line CMS | 24 |
| 10 | YM-155 induces apoptosis independent of Survivin | 25 |
| 10.1 | Gene and protein expression of Survivin. | 25 |
| 10.2 | Validation of Navitoclax | 26 |
| 10.3 | Differential expression after YM-155 inhibition | 27 |
| 11 | A CRISPR screen identifies resistance markers to YM155 treatment in CMS1 | 35 |
| 11.1 | Quality control | 36 |
| 11.2 | Reproducibility | 40 |
| 11.3 | Hit calling | 42 |
| 11.4 | Visualization of results | 43 |
| 12 | Session info | 46 |
| | References | 50 |

1 About

This document contains computer code to reproduce analyses and figures presented in the corresponding study manuscript.

2 Dependencies

We load a number of packages whose functions are needed throughout the analysis.

```
library(tidyverse)
library(affy)
library(GEOquery)
library(preprocessCore)
library(pheatmap)
library(openxlsx)
library(sva)
library(impute)
library(reshape2)
library(ggsignif)
library(CMSclassifier)
library(ggrepel)
library(perm)
library(patchwork)
library(readxl)
library(lumi)
library(limma)
library(fgsea)
library(GO.db)
library(Organism.dplyr)
library(edgeR)
library(patchwork)
```

3 Meta data

In the process of the analysis we will require meta data such as, for example, gene ID maps. In the following we load these meta data, which are provided within this package. We load a map that links gene IDs to microarray probe IDs for a number of different microarray platform used throughout the analysis. We further load a list of genes based on which classification of molecular subtypes is performed.

```
## gene probe ID map
data('probe_map', package='CMSYM1552018')
## reference genes for CMS classification
data('entrez_ref', package='CMSYM1552018')
## list of protein coding genes
data('pc', package='CMSYM1552018')
```

4 Preprocessing microarray data

In the next steps we load and process microarray data from different colorectal cancer cell lines for subtype classification.

4.1 Adai cell line

Adai cell line is data set featured in Oncomine (Rhodes et al., n.d.). Expression profiles for several colorectal cancer cell lines were generated using an Affymetrix HG-U133 Plus2 chip. We read a matrix that contains normalized expression values after background correction using RMA (R. A. Irizarry et al. 2003), cross-chip quantile normalization and expression level summary using medianpolish (all performed using this affy package (Gautier et al. 2004)).

```
data('adai_exprs', package='CMSYM1552018')
```

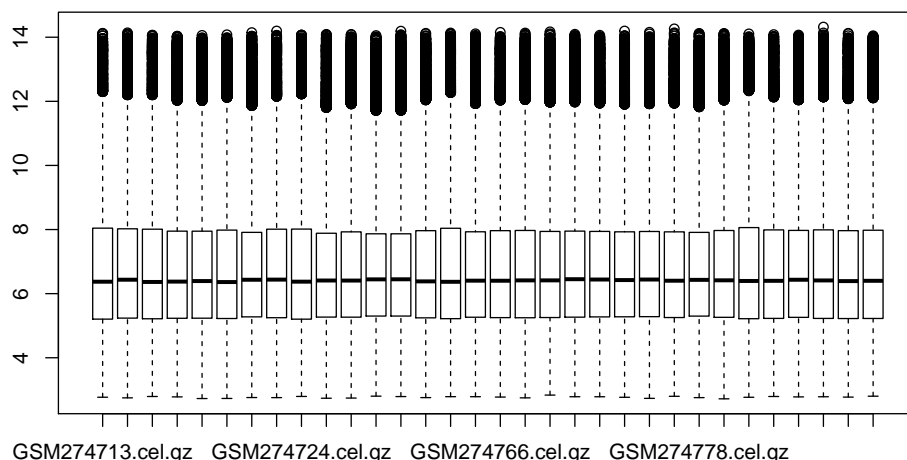


Figure 1: Normalized expression levels across Adai cell line samples

For classification of subtypes we need gene level expression values. We select for each gene the probe with the highest median expression value across samples hypothesizing that this probe corresponds to the 'main' transcript of the gene. We first generate a function for this that can then be applied to other data sets as well.

```
get_cms_gene_lvl2 <- function(eset, chip, filter_ref = T){
  ## make probe mapping
  gpl <- probe_map %>% dplyr::select(matches(chip), entrez) %>%
    drop_na() %>% distinct() %>%
    group_by_at(chip) %>% dplyr::slice(1) %>% ungroup() %>%
    as.data.frame() %>% column_to_rownames(chip)

  ## aggregate per entrez
  eset_aggr <- probesToEntrez(eset, gpl, entrez = 'entrez')

  ## convert to long format
  df_long <- eset_aggr %>% as_tibble(rownames='entrez') %>%
    gather(sample, expr, -entrez)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
if(filter_ref){
  df_long <- df_long %>% filter(entrez %in% entrez_ref)
}

## return results
return(df_long)
}
```

Next we can apply this function to extract the required data from the Adai expression set. The 'chip' argument to the extraction function refers to the column names of the probe_map object.

```
adai_for_cms <- get_cms_gene_lvl2(adai_exprs, 'Affy HG U133-PLUS-2 probeset') %>%
  mutate(dataset='Adai', platform = 'HG_U133_Plus2',
    sample=gsub('.cel.gz', '', sample))
```

Finally we need to know for the purpose of downstream analysis which sample corresponds to which CRC cell lines. We downloaded information about the samples from GEO which we load from an R data file and which we are going to use for this purpose.

```
data('adai_sample_anno', package='CMSYM1552018')
```

```
## add to adai_for_cms data
adai_for_cms <- adai_for_cms %>% left_join(adai_sample_anno)
```

Now we classify the CRC subtypes of the Adai cell line data set using the CMSclassifier (Guinney et al. 2015).

```
## annotation frame and expression matrix for classification
adai_mat <- adai_for_cms %>% dplyr::select(entrez, sample, expr) %>%
  group_by(sample, entrez) %>% summarise(expr=mean(expr)) %>% ungroup() %>%
  spread(sample, expr) %>% `rownames<-`(NULL) %>%
  data.frame() %>% column_to_rownames('entrez')
adai_anno <- adai_for_cms %>% distinct(sample, cellline)

## classify using RF
adai_cms <- classifyCMS.RF(as.data.frame(adai_mat))

## annotate samples with CMS predictions
adai_cms <- adai_anno %>%
  left_join(adai_cms %>% data.frame() %>% mutate(sample=rownames(.)) %>% tbl_df %>%
    `colnames<-`(gsub('.posteriorProb', '', gsub('RF.', '', colnames(.)))) %>%
    mutate(predictedCMS = ifelse(is.na(predictedCMS), 'np', predictedCMS))
```

4.2 Wagner cell line

The Wagner cell line dataset (Wagner et al. 2007) is another Oncomine dataset that includes colorectal cancer cell lines similar to the Adai data set. Again an Affymetrix HG-U133 Plus2 chip was used to characterize gene expression in cell lines. We proceed similar to above.

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
data('wagner_exprs', package='CMSYM1552018')
```

We select the required data from all probe level expression values and aggregate to gene level using the function defined for this purpose above.

```
wagner_for_cms <- get_cms_gene_lvl2(wagner_exprs, 'Affy HG U133-PLUS-2 probeset') %>%  
  mutate(dataset='Wagner', platform = 'HG_U133_Plus2') %>%  
  mutate(sample=gsub('.CEL.gz', '', sample)) %>%  
  mutate(sample=gsub('_\\d$', '', sample))
```

Again we annotate the cell line for each sample.

```
data('wagner_sample_info', package='CMSYM1552018')
```

```
## add to wagner_for_cms data  
wagner_for_cms <- wagner_for_cms %>% left_join(wagner_sample_info)
```

Now we classify the CRC subtypes of the Wagner cell line data set using the CMSclassifier.

```
## annotation frame and expression matrix for classification  
wagner_mat <- wagner_for_cms %>% dplyr::select(entrez, sample, expr) %>%  
  group_by(sample, entrez) %>% summarise(expr=mean(expr)) %>% ungroup() %>%  
  spread(sample, expr) %>% `rownames<-`(NULL) %>%  
  data.frame() %>% column_to_rownames('entrez')  
wagner_anno <- wagner_for_cms %>% distinct(sample, cellline)  
  
## classify using RF  
wagner_cms <- classifyCMS.RF(as.data.frame(wagner_mat))  
  
## annotate samples with CMS predictions  
wagner_cms <- wagner_anno %>%  
  left_join(wagner_cms %>% data.frame() %>% mutate(sample=rownames()) %>% tbl_df %>%  
    `colnames<-`(gsub('.posteriorProb', '', gsub('RF.', '', colnames())))) %>%  
  mutate(predictedCMS = ifelse(is.na(predictedCMS), 'np', predictedCMS))
```

4.3 Garnett cell line

For the next dataset (Garnett et al. 2012) we proceed as above.

```
data('garnett_exprs', package='CMSYM1552018')
```

We select the required data from all probe level expression values and aggregate to gene level using the function defined for this purpose above. A small difference is, that in comparison to the other data sets the Affymetrix HG-U133A platform was used for the experiments.

```
garnett_for_cms <- get_cms_gene_lvl2(garnett_exprs, 'Affy HG U133A probeset') %>%  
  mutate(dataset='Garnett', platform = 'HG_U133A')
```

Again we annotate the cell line for each sample.

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
data('garnett_sample_info', package='CMSYM1552018')
```

```
## add to wagner_for_cms data
garnett_for_cms <- garnett_for_cms %>% left_join(garnett_sample_info)
```

Now we classify the CRC subtypes of the Garnett cell line data set using the CMSclassifier.

```
## annotation frame and expression matrix for classification
garnett_mat <- garnett_for_cms %>% dplyr::select(entrez, sample, expr) %>%
  group_by(sample, entrez) %>% summarise(expr=mean(expr)) %>% ungroup() %>%
  spread(sample, expr) %>% `rownames<-`(NULL) %>%
  data.frame() %>% column_to_rownames('entrez')
garnett_anno <- garnett_for_cms %>% distinct(sample, cellline)

## classify using RF
garnett_cms <- classifyCMS.RF(as.data.frame(garnett_mat))

## annotate samples with CMS predictions
garnett_cms <- garnett_anno %>%
  left_join(garnett_cms %>% data.frame() %>% mutate(sample=gsub('^X', '', rownames(.))) %>% tbl_df %>%
    `colnames<-`(gsub('.posteriorProb', '', gsub('RF.', '', colnames(.)))) %>%
    mutate(predictedCMS = ifelse(is.na(predictedCMS), 'np', predictedCMS))
```

4.4 Barrett cell line

This cell line data set is the older microarray based expression data set of the Cancer Cell Line Encyclopedia (CCLE) (Barretina et al. 2012). We process as above. The HG U133 Plus2 platform was used.

```
data('ccle_exprs', package='CMSYM1552018')
```

We select the required data from all probe level expression values and aggregate to gene level.

```
ccle_for_cms <- get_cms_gene_lv12(ccle_exprs, 'Affy HG U133-PLUS-2 probeset') %>%
  mutate(dataset='CCLE', platform = 'HG_U133_Plus2',
    sample = gsub('.CEL', '', sample))
```

Again we annotate the cell line for each sample.

```
data('ccle_sample_info', package='CMSYM1552018')
```

```
## add to wagner_for_cms data
ccle_for_cms <- ccle_for_cms %>% mutate(sample = gsub('\\.', '-', sample)) %>%
  inner_join(ccle_sample_info)
```

Now we classify the CRC subtypes of the CCLE cell line data set using the CMSclassifier.

```
## annotation frame and expression matrix for classification
ccle_mat <- ccle_for_cms %>% dplyr::select(entrez, sample, expr) %>%
  group_by(sample, entrez) %>% summarise(expr=mean(expr)) %>% ungroup() %>%
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
spread(sample, expr) %>% `rownames<-` (NULL) %>%
data.frame() %>% column_to_rownames('entrez')
ccle_anno <- ccle_for_cms %>% distinct(sample, cellline) %>%
mutate(sample = gsub('-', '.', sample))

## classify using RF
ccle_cms <- classifyCMS.RF(as.data.frame(ccle_mat))

## annotate samples with CMS predictions
ccle_cms <- ccle_anno %>%
left_join(ccle_cms %>% data.frame() %>% mutate(sample=rownames(.)) %>% tbl_df %>%
`colnames<-`(gsub('.posteriorProb', '', gsub('RF.', '', colnames(.)))) %>%
mutate(predictedCMS = ifelse(is.na(predictedCMS), 'np', predictedCMS))
```

Recently there has also been a large RNAseq data set with colorectal cell lines released (Meyers et al. 2017). We can also try to classify based on this set to see if the results are coherent.

```
data('ccle_rnaseq_mat', package='CMSYM1552018')

## classify subtypes
ccle_rnaseq_cms <- classifyCMS.RF(as.data.frame(ccle_rnaseq_mat)) %>%
rownames_to_column('cellline') %>% tbl_df %>%
`colnames<-`(gsub('RF.', '', gsub('.posteriorProb', '',
colnames(.)))) %>%
mutate(sample=cellline) %>%
mutate(predictedCMS = ifelse(is.na(predictedCMS), 'np', predictedCMS))
```

4.5 Medico cell line

This data set was published by Medico and colleagues (Medico et al. 2015) is an especially rich resource of mRNA expression in colorectal cancer cell lines. In contrast to the data sets processed so far, an Illumina microarray platform was used to generate the data so the data cannot be processed exactly as the others. We used the GEOquery package to download the data from GEO (Edgar, Domrachev, and Lash 2002) and can now load the expression values. As we do not have access to the probe level data we will just quantile normalize log-transformed expression values.

```
data('medico_exprs', package='CMSYM1552018')

## quantile normalize
medico_norm <- normalize.quantiles(medico_exprs) %>% log2()
colnames(medico_norm) <- colnames(medico_exprs)
rownames(medico_norm) <- rownames(medico_exprs)
```

Now we can proceed as we did with the other samples.

```
medico_for_cms <- get_cms_gene_lvl2(medico_norm, 'Illumina Human HT 12 V4 probe') %>%
mutate(dataset='Medico', platform = 'Illumina_HT12')
```

Finally, we annotate the cell line corresponding to each sample.


```
data('medico_samples', package='CMSYM1552018')
```

```
## merge sample info with data
medico_for_cms <- medico_for_cms %>% inner_join(medico_samples)
```

Now we classify the CRC subtypes of the CCLE cell line data set using the CMSclassifier.

```
## annotation frame and expression matrix for classification
medico_mat <- medico_for_cms %>% dplyr::select(entrez, sample, expr) %>%
  group_by(sample, entrez) %>% summarise(expr=mean(expr)) %>% ungroup() %>%
  spread(sample, expr) %>% `rownames<-`(NULL) %>%
  data.frame() %>% column_to_rownames('entrez')
medico_anno <- medico_for_cms %>% distinct(sample, cellline) %>%
  mutate(sample = gsub('-', '.', sample))

## classify using RF
medico_cms <- classifyCMS.RF(as.data.frame(medico_mat))

## annotate samples with CMS predictions
medico_cms <- medico_anno %>%
  left_join(medico_cms %>% data.frame() %>% mutate(sample=rownames()) %>% tbl_df %>%
    `colnames<-`(gsub('.posteriorProb', '', gsub('RF.', '', colnames())))) %>%
  mutate(predictedCMS = ifelse(is.na(predictedCMS), 'np', predictedCMS))
```

5 Aggregation of subtype data

Now that we have classified all datasets for molecular subtypes we aggregate them into one table to choose the consensus molecular subtypes. We further annotate information about the microsatellite instability status of each cell line as it correlates well with molecular subtype 1 and might be an interesting biological covariate for downstream analyses.

```
data('msi_status', package='CMSYM1552018')

## combine classification results of individual datasets
all_cms <- adai_cms %>% bind_rows(wagner_cms) %>% bind_rows(garnett_cms) %>%
  bind_rows(ccle_cms) %>% bind_rows(medico_cms) %>% bind_rows(ccle_rnaseq_cms) %>%
  mutate(cellline = toupper(cellline)) %>% group_by(cellline) %>%
  summarise(cms_profile = paste(predictedCMS, collapse=', ')) %>%
  ungroup()

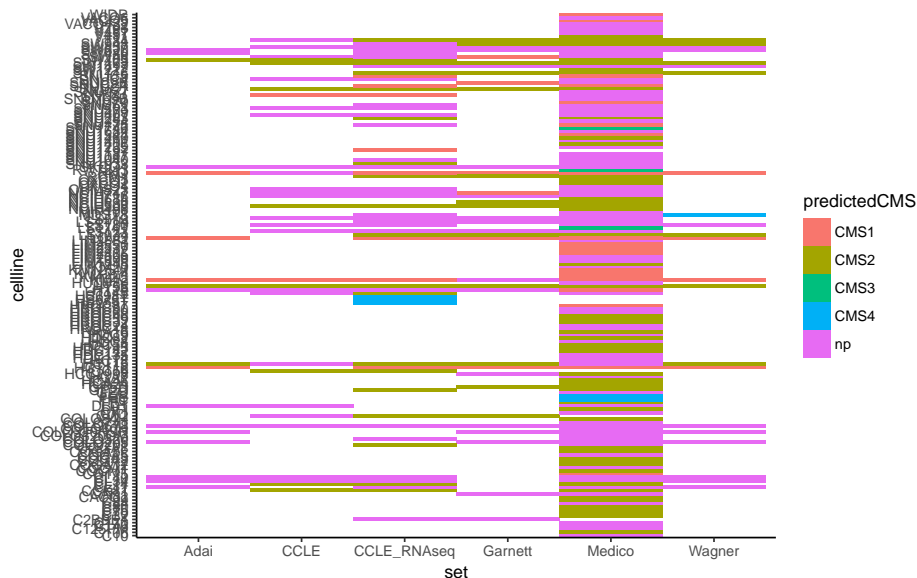
## add MSI info
all_cms <- all_cms %>% left_join(msi_status)
```

We can also generate a visualization of the predicted subtypes across data sets.

```
adai_cms %>% mutate(set = 'Adai') %>%
  bind_rows(wagner_cms %>% mutate(set = 'Wagner')) %>%
  bind_rows(garnett_cms %>% mutate(set = 'Garnett')) %>%
  bind_rows(ccle_cms %>% mutate(set = 'CCLE')) %>%
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
bind_rows(medico_cms %>% mutate(set = 'Medico')) %>%
bind_rows(ccle_rnaseq_cms %>% mutate(set = 'CCLE_RNAseq')) %>%
mutate(cellline = toupper(cellline)) %>%
dplyr::select(cellline, set, predictedCMS) %>% distinct() %>%
ggplot(aes(set, cellline, fill = predictedCMS)) +
geom_tile() + theme_classic()
```



There are no disagreements which is nice to see. It can, however, happen that based on one sample a sub-type could not be assigned where based on another sample classification was possible. In these cases we select the most frequent classification result. Based on this we firmly assign a subtype to each cell line. We are then interest in observing the distribution of predicted subtypes.

```
## by sorting alphabetically we can retrieve a
## subtype if possible as 'C' comes before 'n'
cms_cl <- all_cms %>%
  mutate(cms = map(cms_profile,
                    function(x) strsplit(x, ',') %>% unlist() %>%
                      table() %>% [. == max(.)] %>% names() %>%
                      sort() %>% .[1])) %>% unnest(cms)

## we exclude COL0741, which was shown to be a melanoma cell line
cms_cl <- cms_cl %>% filter(!cellline %in% c('COL0741', 'HCT15', 'C32'))
```

In agreement with other reports in tumour organoids we observe a strong tendency towards classification into CMS1 or CMS2, whereas CMS3 and CMS4 are only rarely observed.

Based on the paper we would assume that there is a significant relationship between the MSI/MSS status and the molecular subtype, where subtype 1 should be more commonly MSI.

```
cms_cl %>% filter(cms %in% c('CMS1', 'CMS2')) %>%
  dplyr::select(-c(cellline, cms_profile)) %>% table() %>% fisher.test()
```

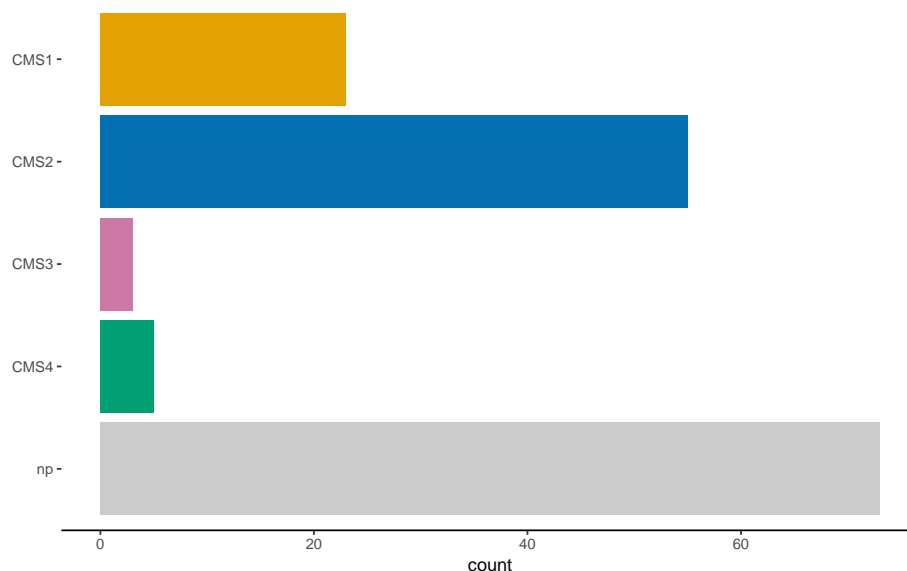


Figure 2: Distribution of CRC cell line subtypes

6 Comparison with other studies

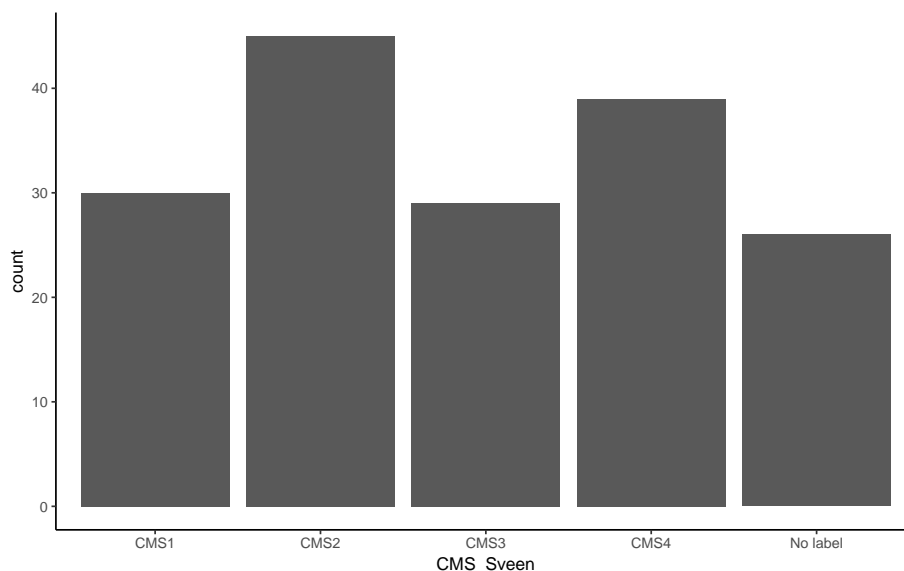
Recently, studies have been published where CMS classification was performed for CRC cell lines (Sveen et al. 2018) (Linnekamp et al. 2018).

6.1 Sveen et al

Here a specific CMS classifier for CRC cell lines was derived based on a smaller set (~450) of genes highly expressed in cell lines. We want to see how consistent these results are with our classification results where the CMS classifier was used as published based on primary tumour samples. We read the resulting Sveen et al classifications into R and compare them to the CMSclassifier class labels.

```
## load sveen cms data
data('sveen_cms', package='CMSYM1552018')

## bar plot of cms label distribution
sveen_cms %>% ggplot(aes(CMS_Sveen)) + geom_bar() + theme_classic()
```



6.2 Linnekamp et al

A supplementary table with subtype predictions is not available from the paper so we type in the predictions manually based on Figure 2.

```
lk_cms <- list(
  CMS1 = c('SNUC2B', 'LS411N', 'LOV0', 'HCT116', 'RK0', 'KM12', 'SW48'),
  CMS2 = c('CCK81', 'GP5D', 'SW1417', 'NCIH630', 'SW1463', 'SW948', 'T84',
            'SNUC1', 'RCM1', 'SW1116', 'HT55', 'LS1034'),
  CMS3 = c('LS123', 'OUMS23', 'NCIH716', 'COL0320HSR', 'HUTU80', 'MDST8', 'CAR1'),
  CMS4 = c('SW620', 'C2BBE1', 'COL0678', 'NCIH747', 'HCT15', 'CL11', 'SW837',
            'HT29', 'SKC01', 'HT115', 'HCC2998', 'COL0205', 'CW2')
)

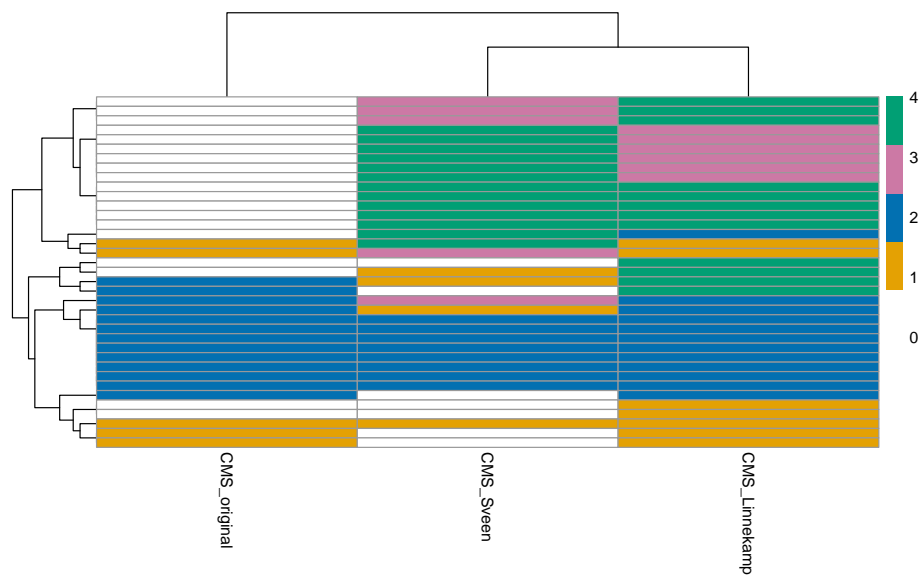
linne_cms <- map2(names(lk_cms), 1:length(lk_cms),
  function(x, y) tibble(CMS_Linnekamp=x, cellline=lk_cms[[y]])) %>%
  bind_rows()
```

6.3 Comparison

We compare the classification labels of our cell lines and the other cell line CMS predictions. We generate both a heat map and a bar plot to visualize overlap.

```
sveen_cms %>% full_join(cms_cl %>% dplyr::select(cellline, CMS_original=cms)) %>%
  full_join(linne_cms) %>% drop_na() %>% mutate_all(~gsub('CMS', '', .)) %>%
  mutate_all(~ifelse(. %in% c('np', 'No label'), 0, .)) %>%
  data.frame() %>% `rownames<-` (NULL) %>% column_to_rownames('cellline') %>%
  apply(2, as.integer) %>%
  pheatmap(color=c('#ffffff', '#e4a103', '#0270b1', '#cc79a5', '#019f74'))
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

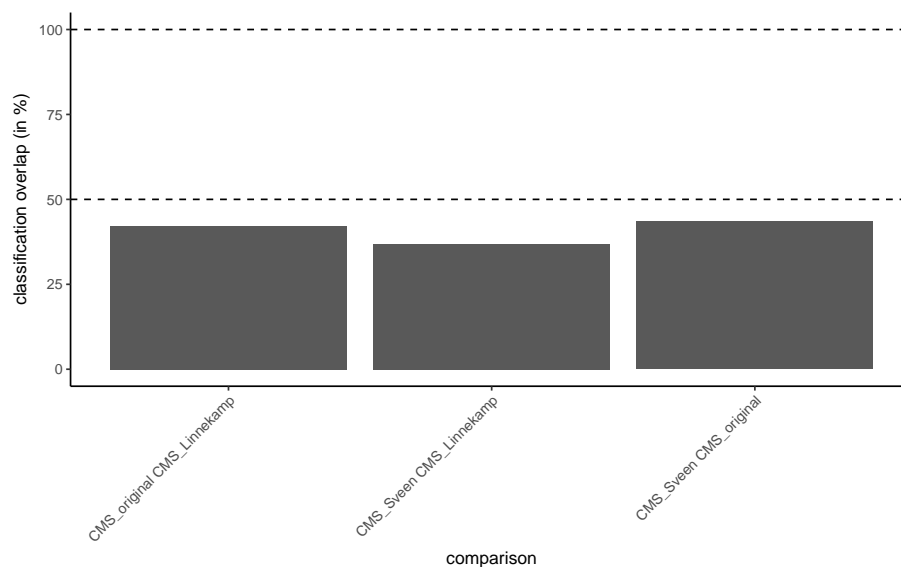


Percent of matching class labels for each pair of classifiers.

```
cms_comp <- sveen_cms %>% full_join(cms_cl %>%
  dplyr::select(cellline, CMS_original=cms)) %>%
  full_join(linne_cms) %>% mutate_all(~gsub('CMS', '', .)) %>%
  mutate_all(~ifelse(. %in% c('np', 'No label'), 0, .)) %>%
  dplyr::select(-cellline) %>% mutate_all('as.integer')

calc_ol <- function(x, y) {
  dat <- cms_comp[,c(x,y)] %>% drop_na() %>% data.frame()
  ol <- sum(dat[,1] == dat[,2]) / nrow(dat)
  return(ol)
}

combn(colnames(cms_comp), 2) %>% t %>% tbl_df %>%
  mutate(overlap = map2(V1, V2, function(x,y) calc_ol(x,y))) %>%
  unnest(overlap) %>% unite(comparison, V1, V2, sep=' ') %>%
  mutate(overlap = overlap * 100) %>%
  ggplot(aes(comparison, overlap)) + geom_bar(stat='identity') +
  theme_classic() + ylim(c(0,100)) +
  geom_hline(yintercept=100, linetype='dashed') +
  geom_hline(yintercept=50, linetype='dashed') +
  ylab('classification overlap (in %)') +
  theme(axis.text.x = element_text(angle=45, hjust=1))
```



7 Characterization of predicted subtypes

We next try to recreate the copy number and mutation plots shown in the paper by Guinney et al., to observe if cell lines behave similarly. To do this, we use data from the COSMIC database (Forbes et al. 2015). Mutations include all non-silent mutation events. Copy number changes of X and Y chromosomes are disregarded.

```
## read and preprocess data exported from CCLE
data('ccle_cnv', package='CMSYM1552018')

## plot cna across groups (boxplot)
scna_plot <- ccle_cnv %>% mutate(cnv=round((2^value)*2, digits=0)) %>%
  filter(!cnv %in% 2, CHR %in% as.character(1:22)) %>%
  count(cms, cellline) %>%
  ggplot(aes(cms, n, fill=cms)) + geom_boxplot() + theme_classic() +
  geom_signif(comparisons = list(c('CMS1', 'CMS2'))) +
  xlab('') + ylab('SCNA count') +
  scale_fill_manual(values=c('#e4a103', '#0270b1',
                             '#019f74', '#ffffff')) +
  theme(legend.position='none')
```

We further aim to reproduce the mutational characteristics described in the original CMS study (Guinney et al. 2015).

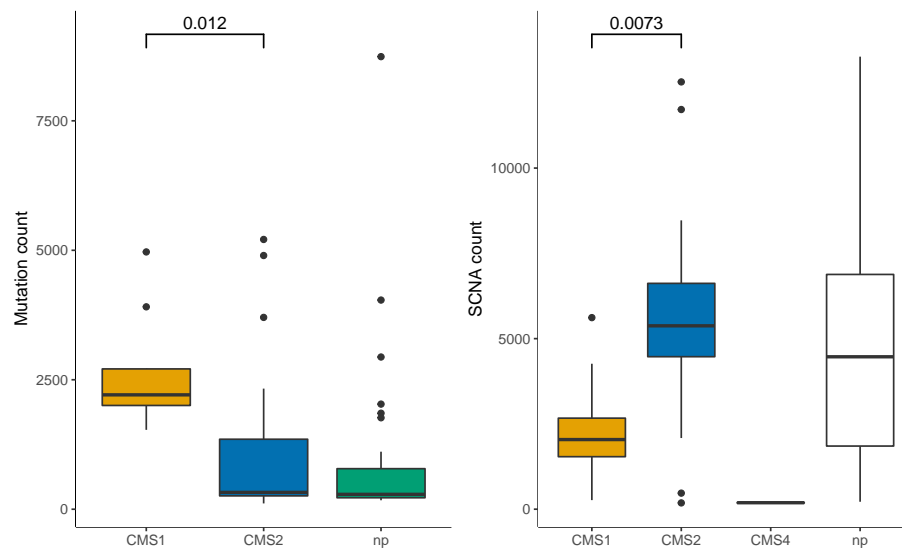
```
## read mutation data derived from COSMIC
data('cosmic_mut_crc', package='CMSYM1552018')

## boxplot of mutations across subtypes
mut_plot <- cosmic_mut_crc %>% dplyr::select(cellline, symbol) %>%
  distinct() %>% inner_join(cms_cl) %>%
  count(cellline, cms) %>% dplyr::select(`Molecular subtype`=cms, Count=n) %>%
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

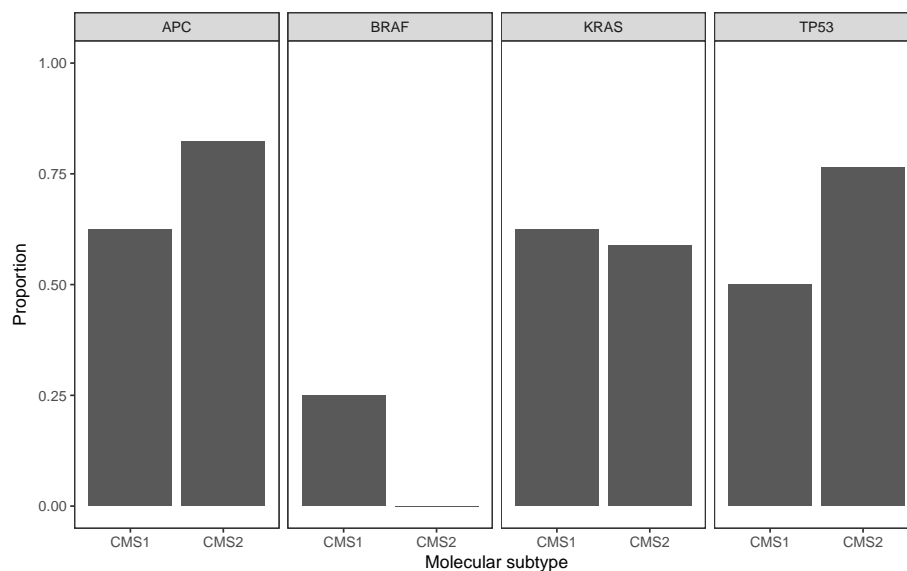
```
ggplot(aes(`Molecular subtype`, Count, fill=`Molecular subtype`)) +
  geom_boxplot() + theme_classic() +
  scale_fill_manual(values=c('#e4a103', '#0270b1',
                             '#019f74', '#ffffff')) +
  geom_signif(comparisons=list(c('CMS1', 'CMS2')) +
  xlab('') + ylab('Mutation count') +
  theme(legend.position='none')

## plot
mut_plot + scna_plot
```



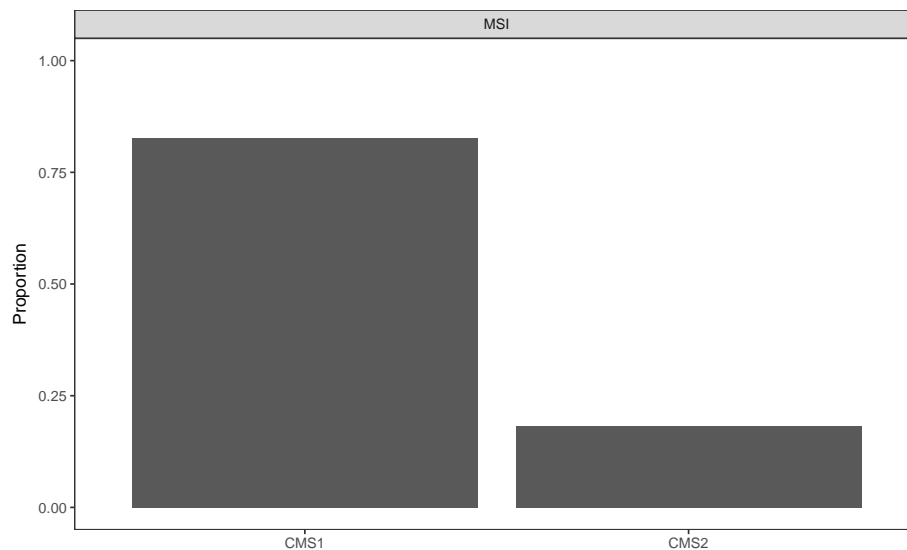
```
## barplot of mutation frequency for known oncogenes
cosmic_mut_crc %>% filter(symbol %in% c('KRAS', 'BRAF', 'APC', 'TP53')) %>%
  dplyr::select(symbol, cellline, `Mutation AA`) %>% distinct() %>%
  inner_join(cms_cl) %>% filter(cellline != 'HT115') %>%
  group_by(cellline, cms) %>% summarise(
    APC=ifelse('APC' %in% symbol, T, F),
    TP53=ifelse('TP53' %in% symbol, T, F),
    KRAS=ifelse('KRAS' %in% symbol, T, F),
    BRAF=ifelse(('BRAF' %in% symbol) &
      ('p.V600E' %in% `Mutation AA`), T, F)) %>% ungroup() %>%
  filter(cms %in% c('CMS1', 'CMS2')) %>%
  group_by(cms) %>% summarise(
    APC=length(which(APC))/n(),
    TP53=length(which(TP53))/n(),
    KRAS=length(which(KRAS))/n(),
    BRAF=length(which(BRAF))/n()) %>% ungroup() %>%
  gather(oncogene, `Mutation rate`, -cms) %>%
  dplyr::select(`Molecular subtype`=cms, everything()) %>%
  ggplot(aes(x=`Molecular subtype`, y=`Mutation rate`)) +
  geom_bar(stat='identity') + ylim(c(0,1)) + theme_bw() +
  facet_wrap(~oncogene, nrow=1) + ylab('Proportion') +
  theme(panel.grid=element_blank())
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



Finally, we compare microsatellite instability status.

```
cms_cl %>% filter(cms %in% c('CMS1', 'CMS2')) %>%  
  group_by(cms) %>%  
  summarise(Proportion=sum(msi_status == 'MSI', na.rm=T)/n()) %>%  
  ungroup() %>% mutate(type='MSI') %>%  
  ggplot(aes(cms, Proportion)) + geom_bar(stat='identity') +  
  theme_bw() + theme(panel.grid=element_blank()) +  
  facet_wrap(~type) + xlab('') + ylim(c(0,1))
```



8 CMS in PDX models

One might hypothesize that lack of CMS3 and CMS4 could be explained by a missing tumour microenvironment. Perhaps engrafted cell lines can change their subtype over time. To investigate this we downloaded expression data of engrafted human cell lines and observe if different subtype classifications are observed over time.

The supplementary materials in Hollingshead et al (Hollingshead et al. 2014) contain already normalized expression data (rma) from the U133+2 microarray platform. We should be able to use these data off the shelf to classify CMS.

```
data('pdx_rma', package='CMSYM1552018')
```

We select relevant probes, select only CRC cell line samples and perform CMS classification.

```
pdx_for_cms <- probe_map %>% filter(entrez %in% entrez_ref) %>%
  dplyr::select(probe = `Affy HG U133-PLUS-2 probeset`, entrez) %>%
  filter(entrez %in% entrez_ref) %>%
  drop_na() %>% distinct() %>%
  inner_join(pdx_rma %>% dplyr::select(-c(symbol, entrez))) %>%
  filter(cellline %in% c(cms_cl$cellline, 'HCT15'))

## select probe for each entrez
## make probe mapping
gpl <- probe_map %>% dplyr::select(matches('PLUS-2'), entrez) %>%
  drop_na() %>% distinct() %>%
  group_by_at(1) %>% dplyr::slice(1) %>% ungroup() %>%
  as.data.frame() %>%
  column_to_rownames('Affy HG U133-PLUS-2 probeset')

## expression matrix
em <- pdx_for_cms %>% acast(probe ~ sample, value.var = 'expr', fun.aggregate = mean)

## aggregate per entrez
eset_aggr <- probesToEntrez(em, gpl, entrez = 'entrez')

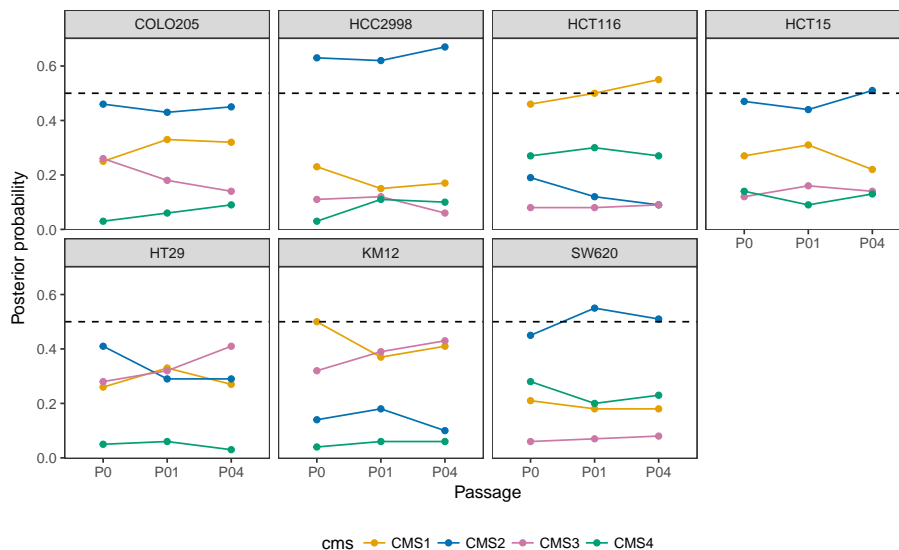
## subtypes for engrafted samples
pdx_cms <- eset_aggr %>%
  as.data.frame() %>% classifyCMS.RF() %>%
  as_tibble(rownames='sample') %>%
  inner_join(distinct(pdx_for_cms, sample, cellline, passage))
```

We visualize the results.

```
## data in format that is easy to plot with ggplot2
for_pdx_viz <- pdx_cms %>% dplyr::select(-sample) %>%
  gather(cms, prob, RF.CMS1.posteriorProb:RF.CMS4.posteriorProb) %>%
  mutate(cms = gsub('RF.', '', gsub('.posteriorProb', '', cms)),
         cms = factor(cms, levels=paste0('CMS', 1:4)),
         passage=factor(passage, levels=c('P0', 'P01', 'P04')))

## line chart with one line for each CMS.
```

```
for_pdx_viz %>%
  ggplot(aes(passage, prob, group=cms, colour=cms)) +
  geom_line() + geom_point() +
  geom_hline(yintercept = 0.5, linetype = 'dashed') +
  facet_wrap(~cellline, nrow = 2) +
  theme_bw() +
  theme(legend.position = 'bottom',
        panel.grid = element_blank()) +
  ylab('Posterior probability') + xlab('Passage') +
  scale_colour_manual(values=c('#e4a103', '#0270b1',
                              '#cc78a6', '#019f74'))
```



9 CMS-dependent drug response driven from public drug screening data

9.1 Identification of candidate substances

We read drug response AUC data from the 1001 cell lines study published in 2016 (Iorio et al. 2016). We select drug response data for CRC cell lines where subtype classifications are available.

```
## read auc data (Table S4; original publication)
data('drug_resp', package='CMSYM1552018')

## drug id-to-name mapping
drug_map <- drug_resp %>% dplyr::slice(1) %>% as.list %>% unlist
drug_map <- data.frame(drug_id=names(drug_map), name=drug_map) %>%
  tbl_df %>% drop_na()
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

Statistical testing for subtype specific drugs. For testing we use a permutation test as implemented in the 'permTS' R package with 10,000 Monte Carlo resamplings.

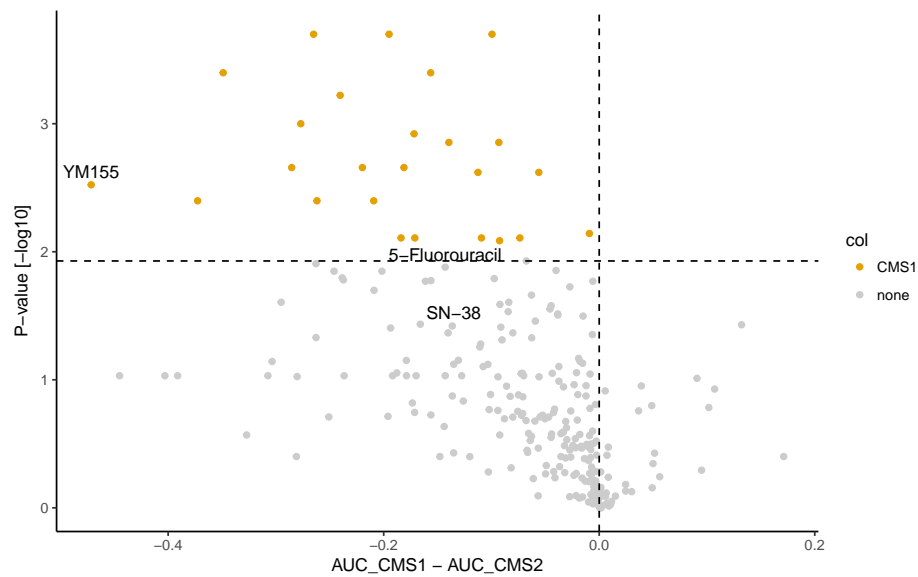
```
## reformat cell line names and filter sub-types
resp_auc <- drug_resp %>% dplyr::slice(-1) %>% tbl_df %>%
  dplyr::select(-Cell.line.cosmic.identifiers) %>%
  mutate(Sample.Names=gsub('-', '',
                             gsub('\\.', '', toupper(Sample.Names)))) %>%
  dplyr::select(cellline=Sample.Names, everything()) %>%
  filter(cellline %in% cms_cl %>% filter(cms %in% c('CMS1', 'CMS2')) %>% .$cellline))
## type cast drug effects to numeric
resp_auc <- cbind(resp_auc[,1], apply(resp_auc[,2:ncol(resp_auc)], 2, as.numeric)) %>% tbl_df
## annotate drugs and cms
resp_auc <- resp_auc %>% melt %>% tbl_df %>%
  dplyr::select(drug_id=variable, everything()) %>%
  inner_join(drug_map) %>% inner_join(cms_cl)

## permutation test
drug_results <- resp_auc %>% nest(-name) %>%
  mutate(test = map(data, ~ permTS(value ~ cms, data=.x, method="exact.mc",
                                   control=permControl(nmc=10^4)))) %>%
  mutate(res = map(test, ~ tibble(p.value = .x$p.value,
                                   delta_auc = .x$estimate))) %>%
  unnest(res) %>% mutate(fdr = p.adjust(p.value, method='BH')) %>%
  arrange(delta_auc)

## visualize as volcano plot
plot_drug_volcano_ym <- function(df){
  ## approximate 20% FDR cutoff
  fdr10 <- df %>% filter(fdr > 0.1) %>% arrange(fdr) %>%
    .$p.value %>% [1] %>% log10() %>% `*(-1)`

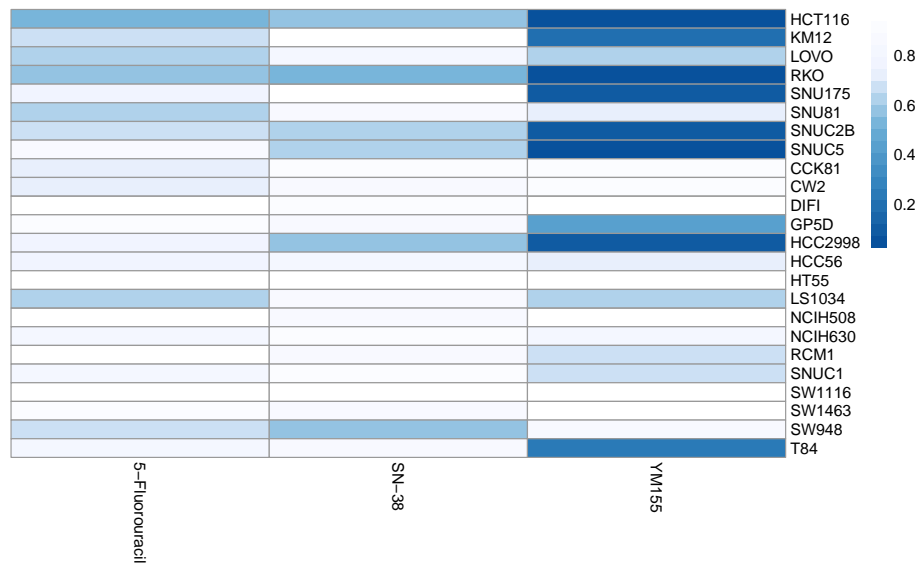
  ## plot volcano
  df %>% mutate(col=ifelse(fdr < 0.1 & delta_auc > 0, 'CMS2',
                          ifelse(fdr < 0.1 & delta_auc < 0, 'CMS1', 'none')) %>%
    mutate(label = ifelse(name %in% c('YM155', 'SN-38', '5-Fluorouracil'),
                          as.character(name), '')) %>%
  ggplot(aes(delta_auc, -log10(p.value), colour=col)) + geom_point() +
  geom_vline(xintercept=0, linetype = 'dashed') +
  geom_hline(yintercept = fdr10, linetype='dashed') +
  geom_text_repel(aes(label=label), nudge_y = 0.1, colour='black') +
  scale_colour_manual(values = c('#e4a103', '#cccccc')) +
  xlab('AUC_CMS1 - AUC_CMS2') + ylab('P-value [-log10]') +
  theme(legend.position = 'none') +
  theme_classic()
}
plot_drug_volcano_ym(drug_results)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



We further plot a small heat map of drug response in individual cell lines.

```
resp_auc %>% filter(name %in% c('YM155', 'SN-38', '5-Fluorouracil'),
                        cms %in% c('CMS1', 'CMS2')) %>% arrange(cms, cellline) %>%
  mutate(cellline = factor(cellline, levels = unique(cellline))) %>%
  reshape2::acast(cellline ~ name, value.var = 'value') %>%
  pheatmap(cluster_rows = F, cluster_cols = F,
            color = colorRampPalette(rev(c('#ffffff', '#eff3ff', '#6baed6', '#3182bd', '#08519c')))(20))
```



9.2 Validation of YM-155

We load the data from proliferation assays that we performed to validate the YM155 effect experimentally in additional cell lines. We define a function that can plot a dose-response curve to use as a figure. We finally define a function that draws a jitter plot comparing area-under-the curve values of cell lines of CMS1 versus CMS2.

```
cell_lines <- c('HCT116', 'RK0', 'LIM2405', 'WiDr', 'LoVo', 'SNU-C2A',
               'CaCo2', 'HT55', 'GP2d', 'HCA-24', 'CL-14', 'SW403')
start_row <- c(3, 18, 33, 48, 63, 78, 93, 108, 123, 138, 153, 168)

cms1 <- c('RK0', 'HCT116', 'LoVo', 'SNU-C2A', 'WiDr', 'LIM2405')
cms2 <- c('HT55', 'CaCo', 'CL-14', 'SW403', 'GP2d', 'HCA-24')

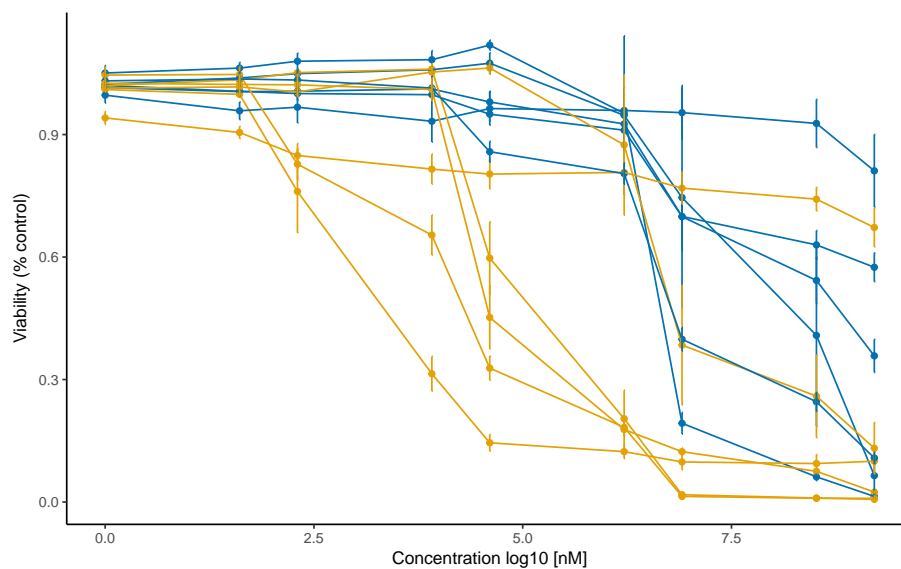
plot_prolif <- function(df){
  df %>% ggplot(aes(conc_log, vmean, colour=CMS, group=cellline)) +
    geom_point() + geom_line() +
    geom_errorbar(aes(ymin = ymin, ymax = ymax), width = 0) + theme_classic() +
    xlab('Concentration log10 [nM]') + ylab('Viability (% control)') +
    scale_colour_manual(values = c('#e4a103', '#0270b1')) +
    theme(legend.position = 'none')
}

compare_plot <- function(df){
  df %>% group_by(cellline, CMS) %>% arrange(conc_log) %>%
    mutate(conc_rank = 1:n()) %>%
    summarise(AUC = pracma::trapz(conc_rank, vmean)) %>% ungroup() %>%
    ggplot(aes(CMS, AUC)) + geom_jitter(width=0.2) +
    stat_summary(fun.y = 'mean', fun.ymin = 'mean',
                fun.ymax = 'mean', geom='crossbar',
                width = 0.5, colour = 'red') + theme_classic() +
    ggsignif::geom_signif(comparisons = list(c('CMS1', 'CMS2')),
                          test = 't.test',
                          test.args = list(var.equal = T))
}

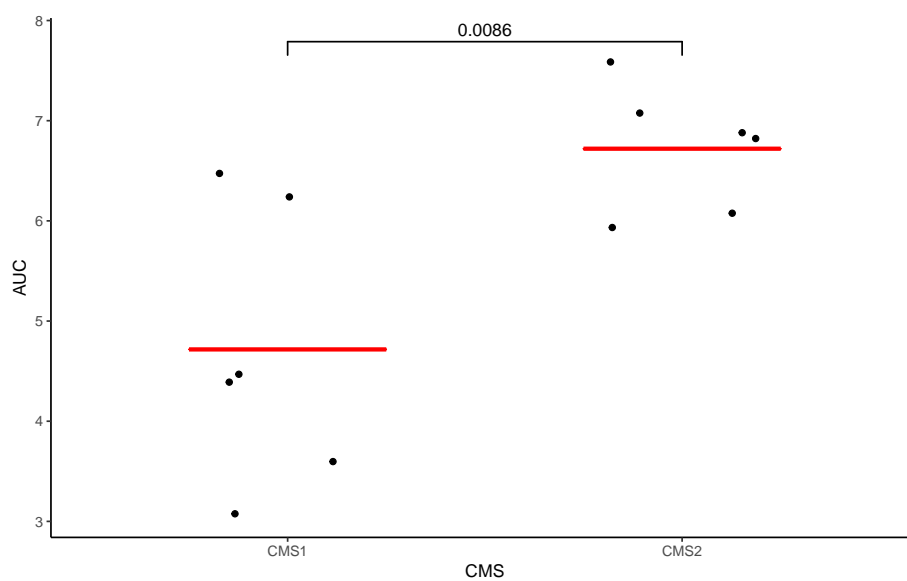
data('ym155', package='CMSYM1552018')

plot_prolif(ym155)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



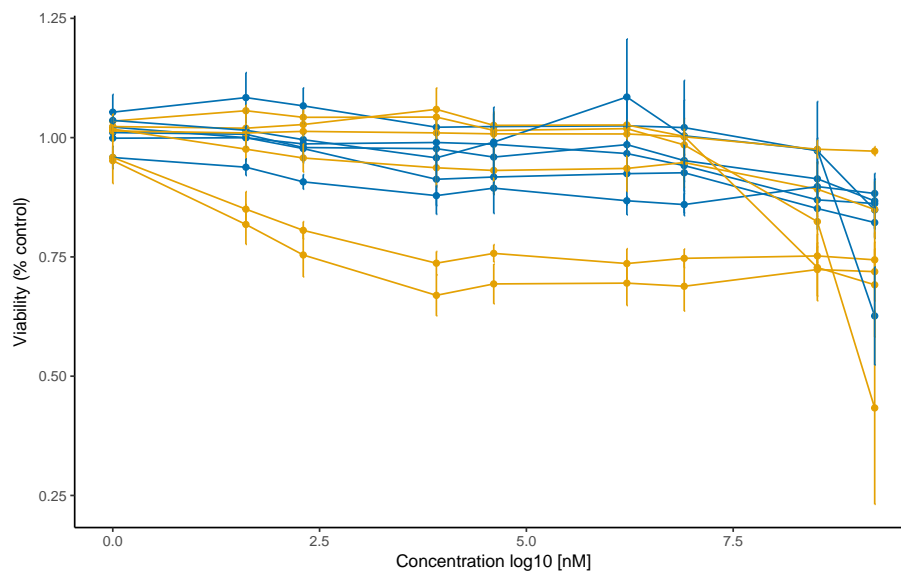
```
compare_plot(ym155)
```



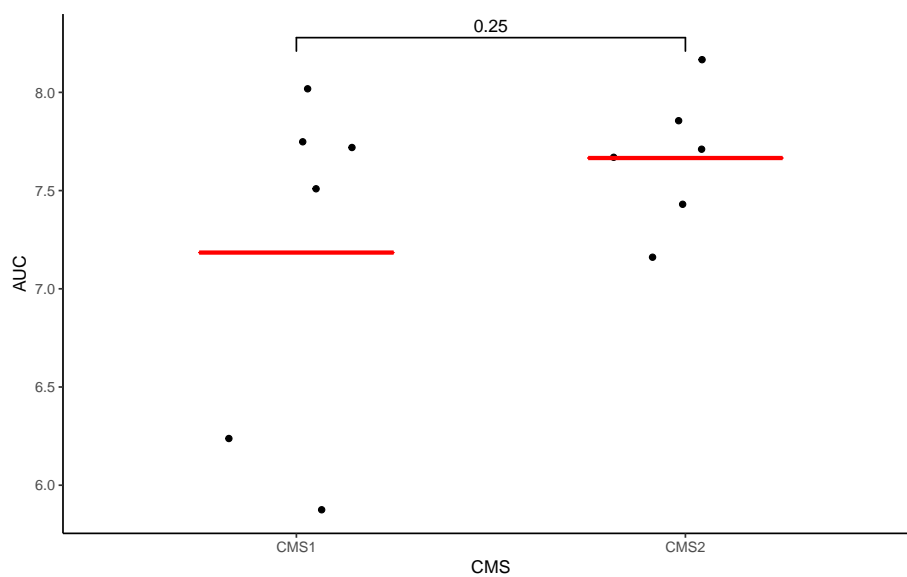
9.3 Validation of 5-FU

```
data('five_fu', package='CMSYM1552018')  
plot_prolif(five_fu)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



```
compare_plot(five_fu)
```

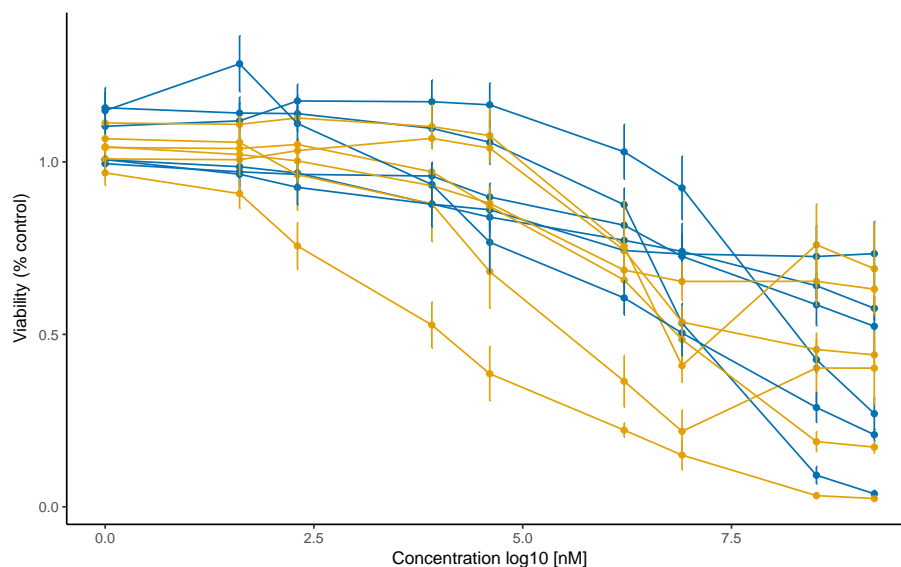


9.4 Validation of SN38 (Irinotecan)

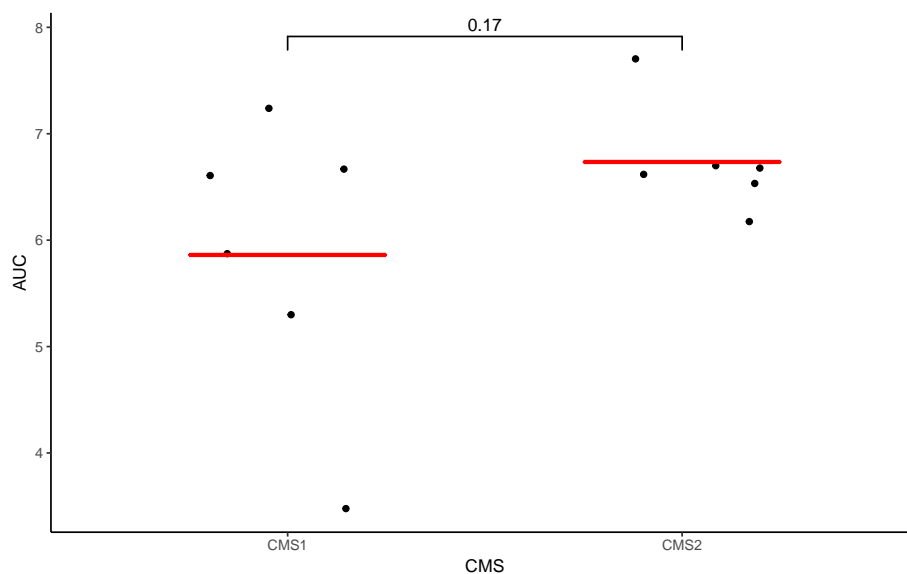
```
data('sn38', package='CMSYM1552018')
```

```
plot_prolif(sn38)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



```
compare_plot(sn38)
```



9.5 Validation of cell line CMS

In order to confirm that the cell lines we analyzed in our lab have the right CMS we performed microarray experiments using Affymetrix arrays. We load data that was normalized as described above.

```
data('exprs_cms_val', package='CMSYM1552018')
```

We analyze CMS status as above.

```
exprs_val <- get_cms_gene_lvl2(exprs_cms_val, 'Affy HG U133-PLUS-2 probeset') %>%
  mutate(dataset='CCLE', platform = 'HG_U133_Plus2',
```



```
sample = gsub('.CEL', '', sample))

cms_val <- acast(exprs_val, entrez ~ sample, value.var = 'expr') %>%
  as.data.frame() %>%
  classifyCMS.RF(minPosterior = 0.35)
```

The CMS of all cell lines can be confirmed if a slightly more relaxed posterior probability threshold of 0.35 is applied. There are no cell lines for which we observe a class label that disagrees with our previous results.

10 YM-155 induces apoptosis independent of Survivin

10.1 Gene and protein expression of Survivin

Protein expression data are derived from Frejno et al. (Frejno et al. 2017).

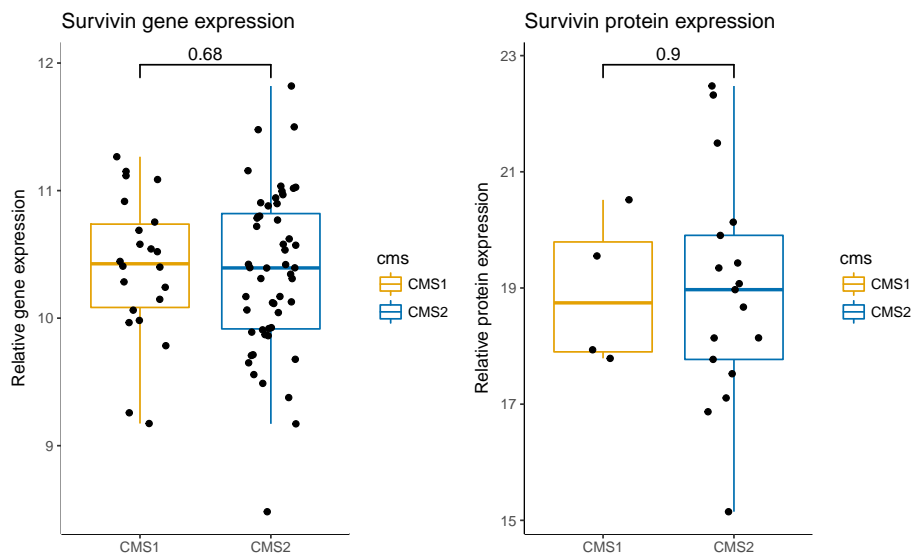
```
## gene expression of survivin across CMS1 and CMS2
survivin_gene <- medico_for_cms %>% filter(entrez == 332) %>% inner_join(cms_cl) %>%
  filter(cms %in% c('CMS1', 'CMS2')) %>%
  ggplot(aes(cms, expr)) + geom_boxplot(aes(colour=cms), fill = NA,
                                         outlier.colour = 'ffffff') +
  geom_jitter(width=0.2) +
  scale_colour_manual(values = c('#e4a103', '#0270b1')) +
  ylab('Relative gene expression') + xlab('') +
  ggsignif::geom_signif(comparisons = list(c('CMS1', 'CMS2'))) +
  ggtitle('Survivin gene expression') +
  theme(legend.position = 'none')

## protein expression Survivin vs YM155 response
data('protein_expr', package='CMSYM1552018')

## protein expression Survivin vs YM155 response
survivin_protein <- protein_expr %>% filter(symbol == 'BIRC5') %>%
  mutate(preexpr=as.numeric(preexpr)) %>%
  inner_join(cms_cl) %>% filter(cms %in% c('CMS1', 'CMS2')) %>%
  drop_na() %>% ggplot(aes(cms, preexpr)) +
  geom_boxplot(aes(colour=cms), outlier.colour = 'ffffff') +
  geom_jitter(width=0.2) +
  scale_colour_manual(values = c('#e4a103', '#0270b1')) +
  ylab('Relative protein expression') + xlab('') +
  ggsignif::geom_signif(comparisons = list(c('CMS1', 'CMS2'))) +
  ggtitle('Survivin protein expression') +
  theme(legend.position = 'none')

## plot both with patchwork
survivin_gene + theme_classic() +
  survivin_protein + theme_classic()
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

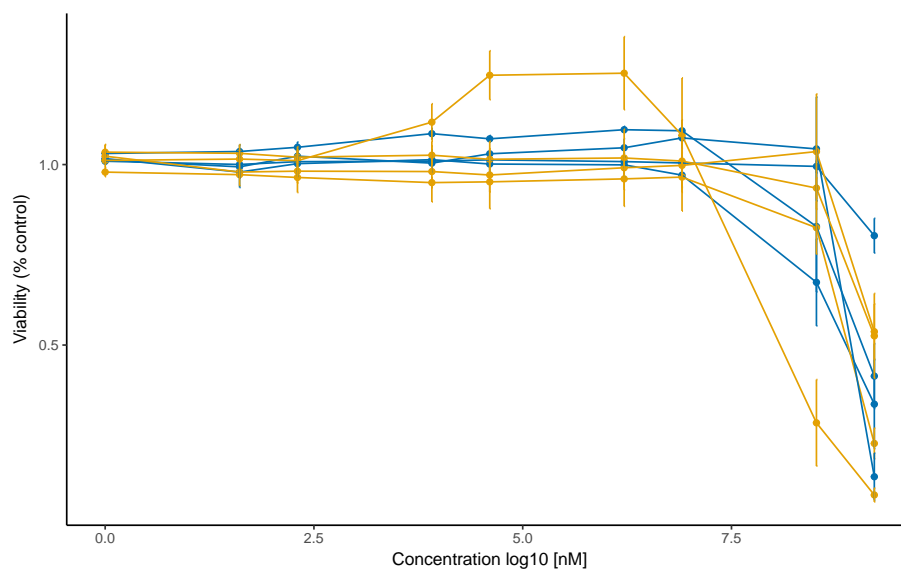


10.2 Validation of Navitoclax

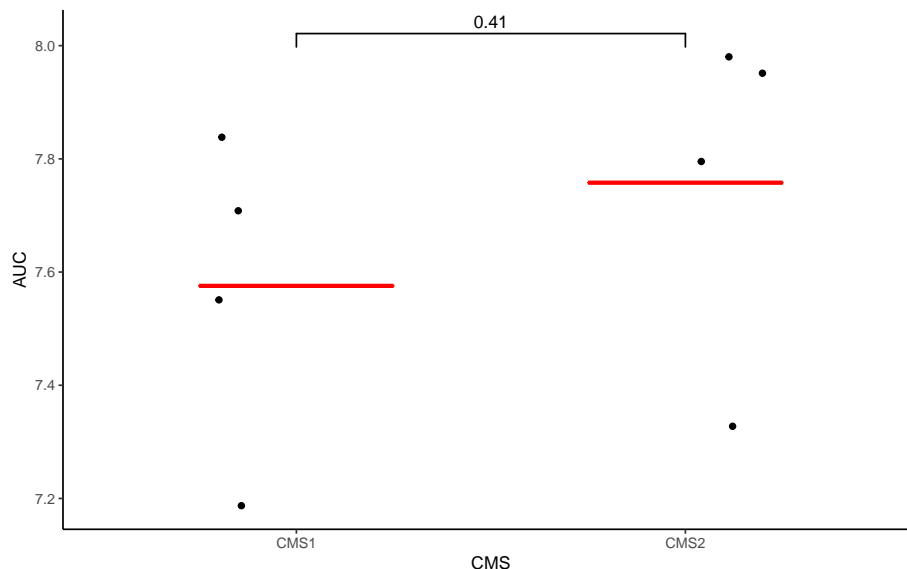
Proliferation assays were performed similar to above for YM155, 5-FU and SN-38. We load the data and plot dose response curves.

```
data('navitoclax', package='CMSYM1552018')
```

```
plot_prolif(navitoclax)
```



```
compare_plot(navitoclax)
```



10.3 Differential expression after YM-155 inhibition

To assess whether there are transcriptomic changes associated with YM-155 inhibition we treated cells with 100 ul of YM-155 and measured gene expression.

10.3.1 Loading array data

The data were generated using an Illumina BeadChip array. We use the lumi package (Du, Kibbe, and Lin 2008) to process the expression measurements. We start by loading the raw data in form of a lumi batch object.

```
data('lumi_raw', package='CMSYM1552018')
data('probe_mapping_ilmn', package='CMSYM1552018')
```

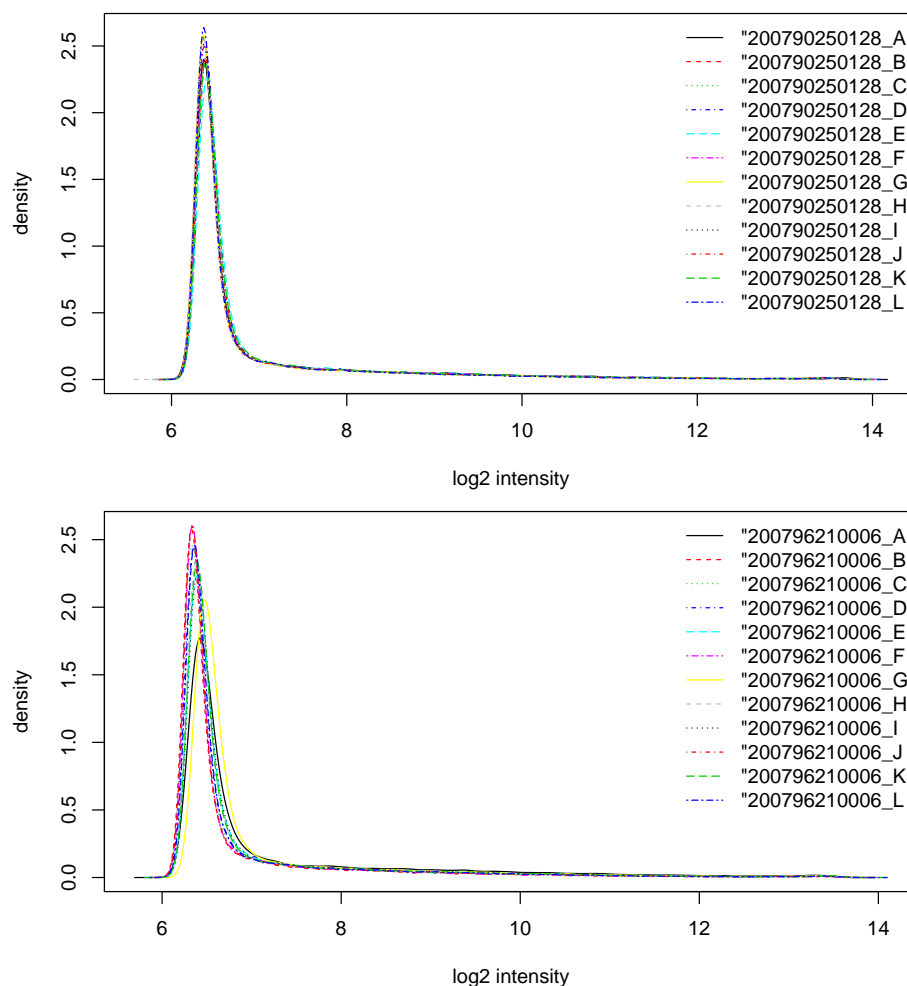
Lumi comes with a number of quality control features that we can use to determine microarray quality. First we report some summary statistics.

```
walk(lumi_raw, summary, 'QC')
```

Next we generate sample density plots.

```
walk(lumi_raw, density)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



10.3.2 Normalization

We use RMA and quantile normalization to normalize the data.

```
lumi_norm <- lumiExpresso(combine(lumi_raw[[1]], lumi_raw[[2]]))
```

10.3.3 Quality control

By PCA and clustering we want to understand whether the normalized data looks reasonable. First we annotate the sample columns for easy interpretation.

```
## load sample info to annotate samples
data('sample_sheet', package='CMSYM1552018')

## extract expression matrix
lumi_norm_exprs <- exprs(lumi_norm)

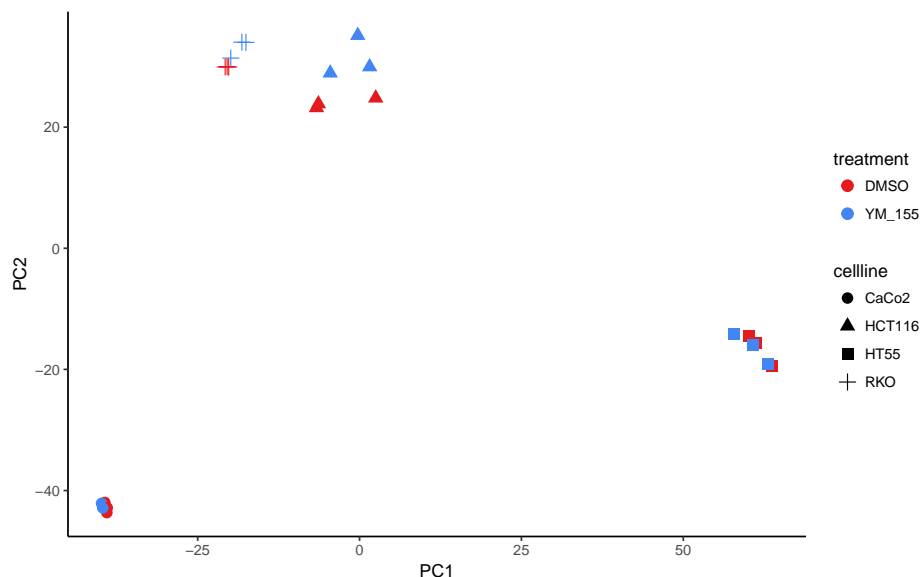
## set sample names
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
sn_ind <- match(paste0('', sample_sheet$sample_id), colnames(lumi_norm_exprs))
colnames(lumi_norm_exprs) <- sample_sheet$Sample_Name[sn_ind]
```

Then we run a PCA on the data and plot the first two principal components.

```
pca <- lumi_norm_exprs %>% t() %>% prcomp()
pca$x %>% as_tibble(rownames='sample') %>%
  separate(sample, c('cellline', 'replicate', 'treatment'), sep=' ', remove=F) %>%
  ggplot(aes(PC1, PC2, colour=treatment, shape=cellline)) +
  geom_point(size=3) + theme_classic() +
  scale_colour_manual(values = c('#E41A1C', '#4285f4'))
```



This looks as expected. CMS1 cell lines are closer together than to the CMS 2 cell lines. Differences between cell lines dominate the YM155 treatment which we would expect. Nevertheless differences can be seen between treated and untreated samples and these seem to be somewhat larger for the CMS1 lines.

10.3.4 Differential gene expression analysis

We now use limma (Ritchie et al. 2015) to find genes that are differentially regulated upon YM155 treatment in CMS1 and CMS2 cell lines.

```
## sample annotation for expression matrix columns
sample_anno <- sample_sheet %>%
  separate(Sample_Name, c('cellline', 'replicate', 'treatment'), sep=' ', remove=F) %>%
  mutate(cms = ifelse(cellline %in% c('CaCo2', 'HT55'), 'CMS2', 'CMS1')) %>%
  mutate(treatment_cms = ifelse(treatment == 'DMSO', 'DMSO',
                                ifelse((treatment == 'YM155') & (cms == 'CMS1'), 'YM155-CMS1',
                                       'YM155-CMS2'))))

## check that sample anno matches to expr matrix
identical(colnames(lumi_norm_exprs), sample_anno$Sample_Name)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
## define model matrix
mm <- model.matrix(~ cellline + treatment_cms, data = sample_anno)

fit <- lmFit(lumi_norm_exprs, mm)
fit_eb <- eBayes(fit)

## results for cms1
dge_res_cms1 <- topTable(fit_eb, coef = 5, n=Inf) %>%
  as_tibble(rownames='ProbeID') %>%
  left_join(probe_mapping_ilmn) %>%
  filter(!is.na(Symbol))

## results for cms2
dge_res_cms2 <- topTable(fit_eb, coef = 6, n=Inf) %>%
  as_tibble(rownames='ProbeID') %>%
  left_join(probe_mapping_ilmn) %>%
  filter(!is.na(Symbol))
```

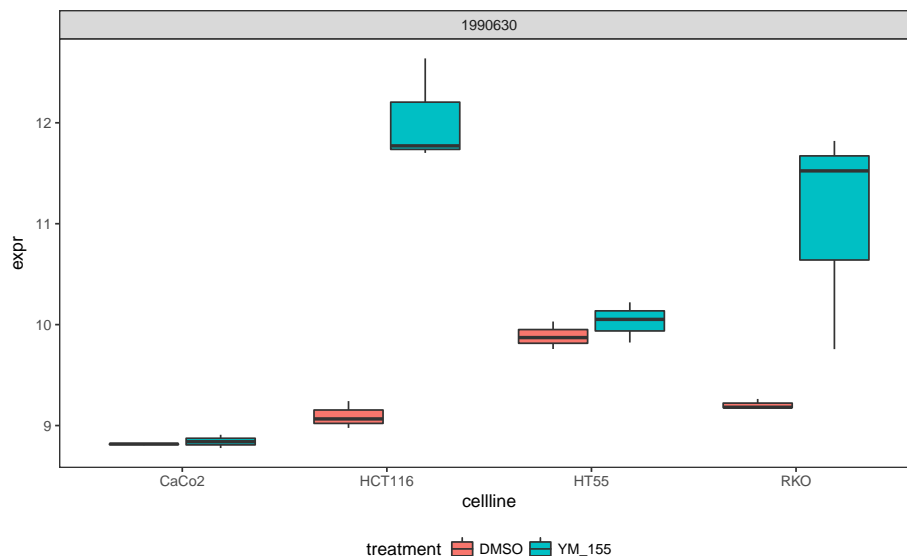
To explore the results we first generate a function that can plot expression levels for all probes targeting a gene of interest across different groups.

```
expr_long <- lumi_norm_exprs %>% as_tibble(rownames='ProbeID') %>%
  left_join(probe_mapping_ilmn) %>%
  gather(sample, expr, -c(ProbeID, Entrez_Gene_ID, Symbol)) %>%
  filter(!is.na(Symbol)) %>%
  separate(sample, c('cellline', 'replicate', 'treatment'), sep=' ', remove=F)

plot_gene_expr <- function(symbol){
  expr_long %>% filter(Symbol == symbol) %>%
    ggplot(aes(cellline, expr, fill = treatment)) +
    geom_boxplot() + facet_wrap(~ProbeID) +
    theme_bw() + theme(panel.grid = element_blank(), legend.position = 'bottom')
}

## Plot gene expression of TRIB3
## - a key player in ER-stress induced apoptosis
plot_gene_expr('TRIB3')
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



We find no differentially expressed genes in CMS2 samples why nice signal is present for CMS1. It seems from the data that especially ER-stress induced apoptosis and NFkappaB signaling are induced by the YM-155 treatment. We can confirm this by gene set enrichment analysis (Sergushichev 2016).

```
## gene ontology annotation
go_ids <- as_tibble(as.data.frame(GOTERM)) %>%
  filter(ontology == 'BP') %>% distinct(go_id, Term)
src <- src_organism("TxDb.Hsapiens.UCSC.hg38.knownGene")
go_map <- tbl(src, 'id_go') %>% filter(ontology == 'BP') %>%
  distinct(entrez, go) %>% collect(n=Inf) %>%
  inner_join(go_ids %>% dplyr::select(go = go_id, term = Term)) %>%
  dplyr::select(-go) %>% split(. $term) %>% map(~ .x$entrez)

## ranked list
ranks <- setNames(dge_res_cms1$t, dge_res_cms1$Entrez_Gene_ID) %>%
  sort(decreasing=T)
ranks2 <- setNames(dge_res_cms2$t, dge_res_cms2$Entrez_Gene_ID) %>%
  sort(decreasing=T)
fgsea_res_go <- fgsea(go_map, ranks, nperm=1e5, maxSize=500) %>%
  as_tibble() %>% arrange(pval)
fgsea_res_cms2 <- fgsea(go_map, ranks2, nperm=1e4, maxSize=500) %>%
  as_tibble() %>% arrange(pval)
```

We define a custom function that can draw barcode plots to visualize gene set enrichment.

```
pretty_barcode_plot <- function(stat_vector, sig){
  require(fgsea)
  ## genes in signature
  sig_genes <- sig[[1]]
  if(NA %in% stat_vector){
    warning('Removing NAs from ranked gene list')
    stat_vector <- stat_vector[!is.na(stat_vector)]
  }
}
```

```
## generate barcode plot
bc_plot <- plotEnrichment(sig_genes, stat_vector)

## remove unwanted layers
bc_plot$layers <- list()

## add barcode at the bottom
lowest_pos <- min(bc_plot$data[,2])
dash_length <- abs(purrr::reduce(range(bc_plot$data[,2]), `~`)*0.1)
middle <- which.min(abs(stat_vector))

## p-value and negative enrichment score
pval <- limma::cameraPR(stat_vector, sig[[1]])$PValue

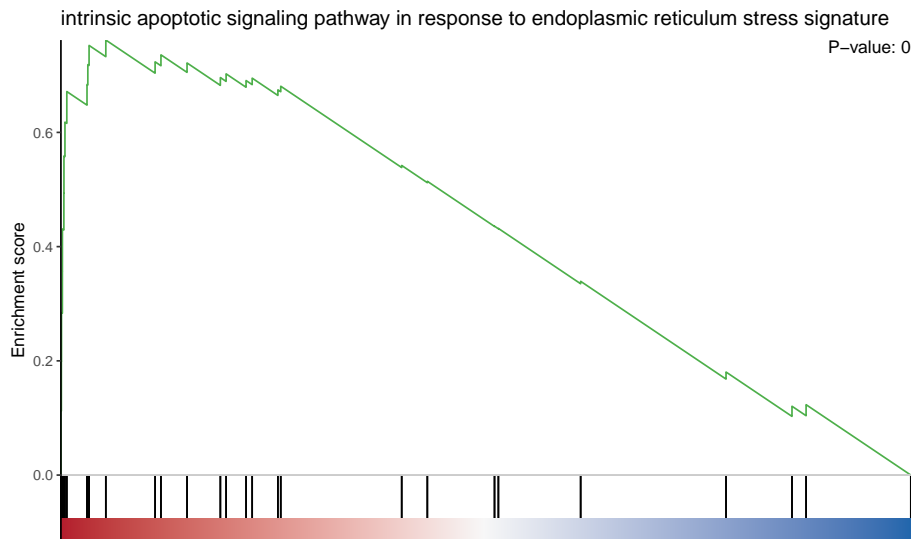
bc_plot_custom <- bc_plot + geom_segment(aes(x=x, xend=x), y=lowest_pos,
                                           yend=lowest_pos-dash_length) +
  geom_line(colour='#4daf4a') +
  geom_hline(yintercept=lowest_pos, colour='#ccccc') +
  geom_hline(yintercept=0, colour='#ccccc') + xlab('') +
  theme_classic() +
  geom_tile(data=tibble(rank=1:length(stat_vector),
                        y=lowest_pos-(1.25*dash_length)),
            aes(x=rank, y=y, fill=rank),
            width=1,
            height=0.5*dash_length) +
  scale_fill_gradient2(low = '#b2182b', high='#2166ac',
                       mid='#f7f7f7', midpoint = middle) +
  scale_x_continuous(expand = c(0, 0)) +
  scale_y_continuous(expand = c(0, 0)) +
  theme(panel.grid=element_blank(),
        axis.text.x=element_blank(),
        axis.ticks.x = element_blank(),
        legend.position = 'none') +
  ggtitle(paste(names(sig)[1], 'signature')) +
  ylab('Enrichment score') +
  annotate('text', x = Inf, y = Inf, hjust=1, vjust = 1,
          label = paste('P-value:', round(pval, 4)))

return(bc_plot_custom)
}
```

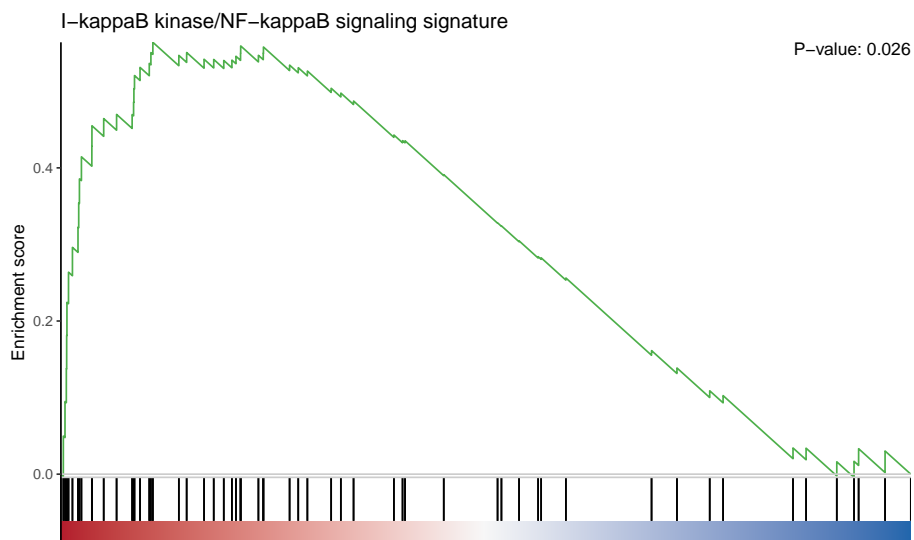
We draw barcode plots for the mentioned processes of interest. We further draw a volcano plot that can globally visualize expression changes upon YM-155 treatment selecting the probe with the strongest fold change to represent each gene.

```
## CMS1
pretty_barcode_plot(ranks, go_map['intrinsic apoptotic signaling pathway in response to endoplasmic reticulum'])
```


Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



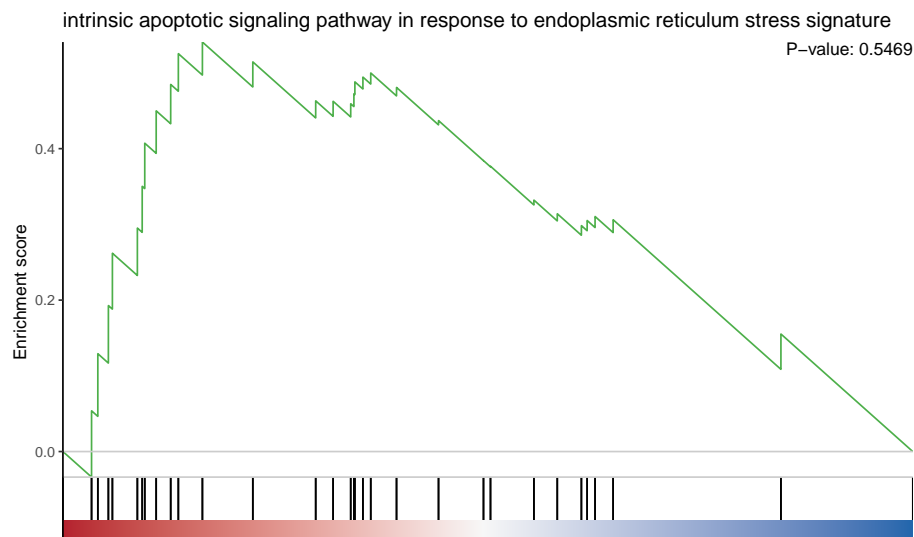
```
pretty_barcode_plot(ranks, go_map['I-kappaB kinase/NF-kappaB signaling'])
```



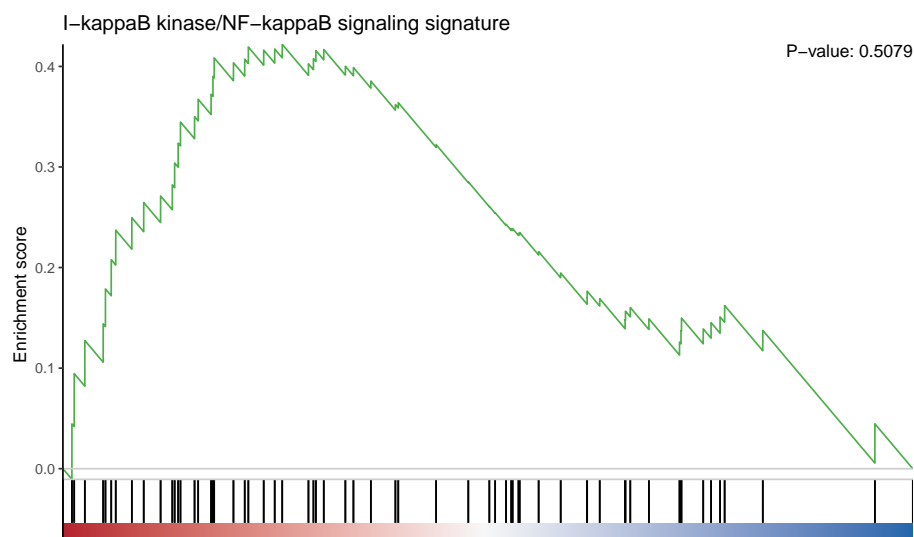
```
## CMS2
```

```
pretty_barcode_plot(ranks2, go_map['intrinsic apoptotic signaling pathway in response to endoplasmic reticulum stress'])
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



```
pretty_barcode_plot(ranks2, go_map['I-kappaB kinase/NF-kappaB signaling'])
```



```
bc_plot <- function(df, plot_names=T){
  if(plot_names){
    chop <- c('TIRB3', 'DDIT3', 'TNFRSF10B', 'ATF4', 'PPP1R15A', 'TRIB3')
    nfkb <- c('NFKB1', 'IKKB', 'NFKB2', 'NKIRAS2')
  } else {
    genes <- c()
  }

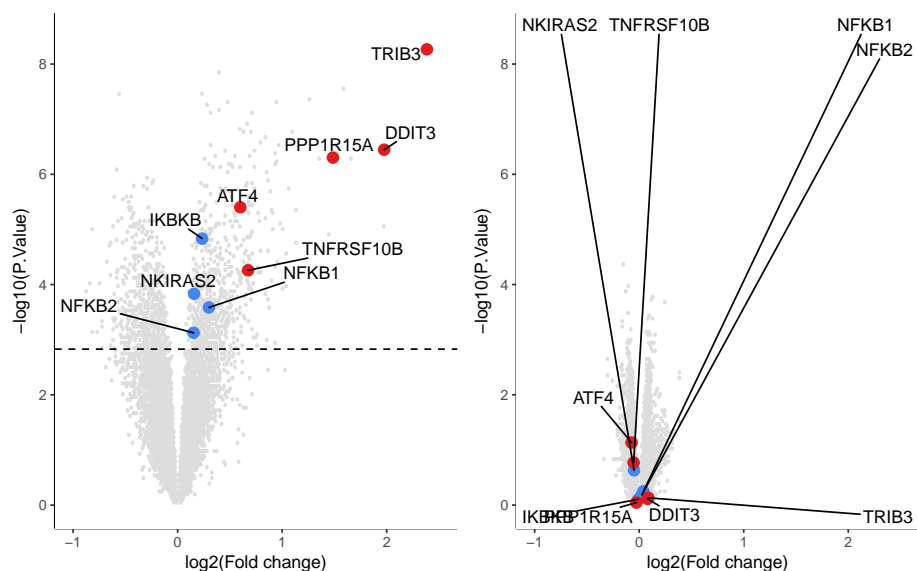
  df <- df %>% group_by(Symbol) %>% arrange(desc(abs(t))) %>%
    dplyr::slice(1) %>% ungroup() %>%
    mutate(label = ifelse(Symbol %in% c(chop, nfkb), Symbol, ''),
           highlight = ifelse(label %in% chop, 'CHOP',
                              ifelse(label %in% nfkb, 'NFKB', 'none')))
```

```
fdr5 <- df %>% filter(adj.P.Val < 0.05) %>% arrange(desc(P.Value)) %>%
  pull(P.Value) %>% head(1) %>% log10() %>% `*`(-1)

df %>% ggplot(aes(logFC, -log10(P.Value), label = label)) +
  geom_hex(data = subset(df, highlight == 'none'), colour = '#ddddd', fill='#ddddd', bins=150, size=2) +
  geom_point(data = subset(df, highlight == 'NFKB'), colour= '#4285f4', size=3) +
  geom_point(data = subset(df, highlight == 'CHOP'), colour= '#e41a1c', size=3) +
  ggrepel::geom_text_repel() +
  geom_hline(yintercept = fdr5, linetype = 'dashed') +
  theme_classic() + theme(legend.position = 'none') +
  xlab('log2(Fold change)') + ylim(c(0, 8.5)) + xlim(c(-1, 2.5))
}

p1 <- bc_plot(dge_res_cms1)
p2 <- bc_plot(dge_res_cms2)

p1 + p2
```



11 A CRISPR screen identifies resistance markers to YM155 treatment in CMS1

We generated read counts from the initial sequencing (fastq) data using [CRISPRanalyzerR](#) (Winter et al. 2017). We now load these counts into R to analyze them further.

```
data('counts', package='CMSYM1552018')
```

11.1 Quality control

11.1.1 Sequencing depth

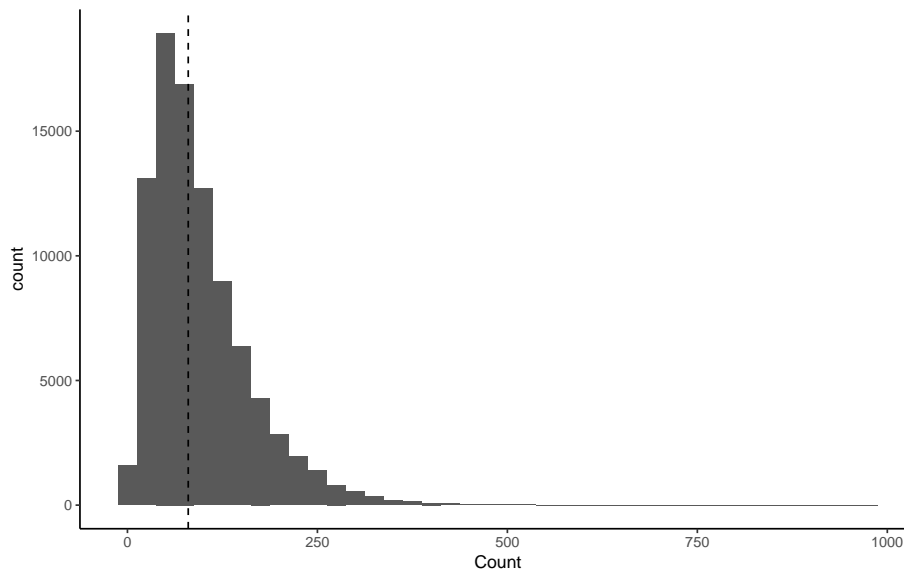
First we would like to assess sequencing depth for each sample and compare it to previous screens.

11.1.1.1 Time 0 sample

We start with the initial time point (T0).

```
## our T0 sample
median_our <- counts %>% filter(sample == 'd0screenlibrary') %>%
  .$Count %>% median()

counts %>% filter(sample == 'd0screenlibrary') %>%
  ggplot() + geom_histogram(aes(Count), bins=40) +
  geom_vline(xintercept=median_our, linetype='dashed') +
  theme_classic()
```

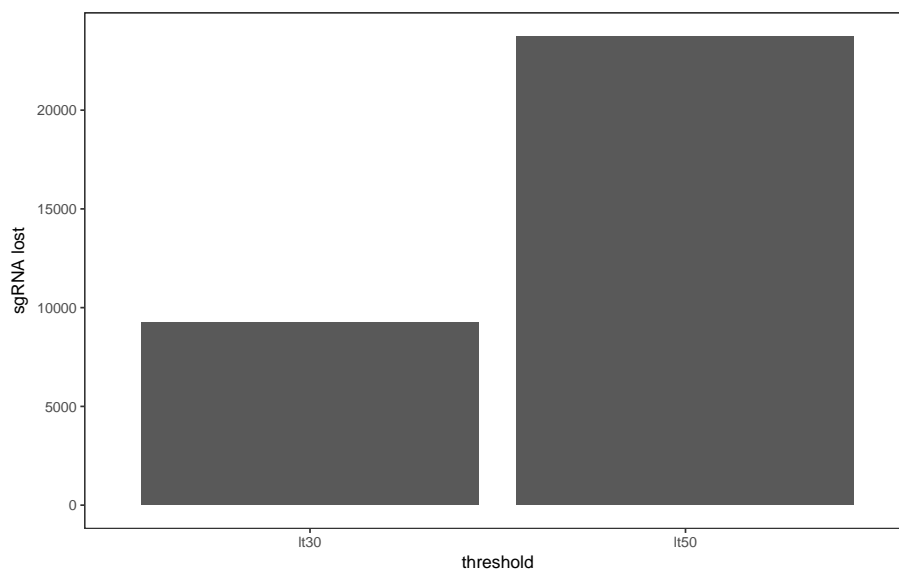


The sequencing depth should be sufficient to continue with the analysis of the screen.

```
counts %>%
  filter(sample %in% c('rc_initial', 'HCTwtT0', 'd0screenlibrary')) %>%
  mutate(lt30 = Count < 30, lt50 = Count < 50) %>%
  summarise(lt30 = sum(lt30), lt50 = sum(lt50)) %>%
  gather(threshold, `sgRNA lost`, lt30, lt50) %>%
  ggplot(aes(threshold, `sgRNA lost`)) +
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
geom_bar(stat='identity') +  
theme_bw() + theme(panel.grid = element_blank())
```



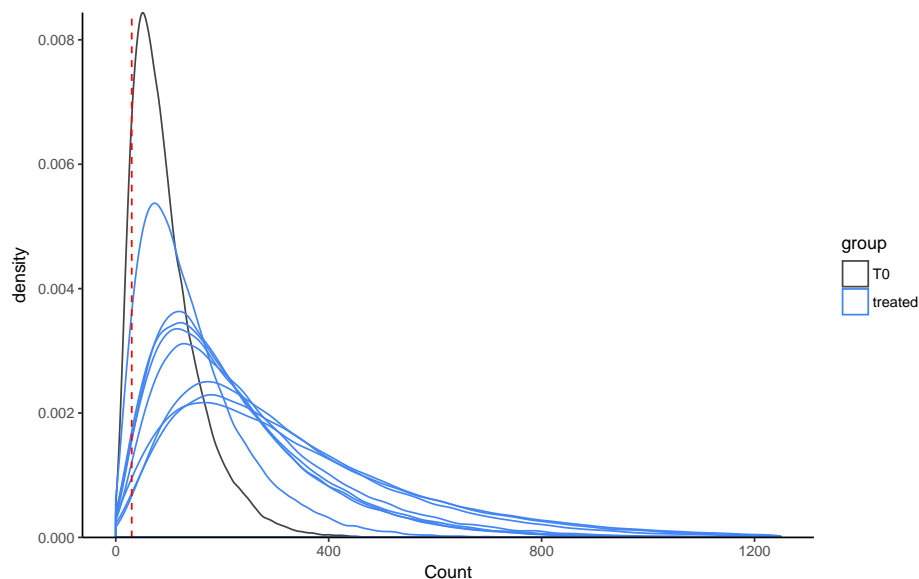
To avoid losing too many sgRNAs due to insufficient coverage we will continue with a less stringent exclusion criterium (< 30), which results in only a loss of ~10% of the sgRNAs.

11.1.1.2 All samples

We plot histograms for read count distributions of all samples.

```
counts %>%  
  mutate(group = ifelse(grepl('d0', sample), 'T0', 'treated')) %>%  
  ## exclude outlier counts so we can see the distribution  
  filter(Count < 1250) %>%  
  ggplot() + geom_density(aes(Count, group=sample, colour=group)) +  
  geom_vline(xintercept = 30, colour = 'red', linetype='dashed') +  
  scale_colour_manual(values=c('#444444', '#4285f4')) +  
  scale_y_continuous(expand = c(0,0)) +  
  theme_classic()
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



Summary statistics to sum up the sequencing.

```
counts %>%
  group_by(sample) %>%
  summarise(total = sum(Count), mean = mean(Count),
            median = median(Count)) %>% ungroup()

#> # A tibble: 9 x 4
#>   sample                total mean median
#>   <chr>                  <int> <dbl> <dbl>
#> 1 d0screenlibrary      8727438  95.6    80
#> 2 Screen1DMSOR1        20080262 220.    182
#> 3 Screen1DMSOR2        20872124 229.    189
#> 4 Screen1NavitoclaxR1  31602078 346.    282
#> 5 Screen1NavitoclaxR2  13021252 143.    118
#> 6 Screen2DMSOR1        22875111 250.    208
#> 7 Screen2DMSOR2        19900023 218.    180
#> 8 Screen2YM155R1       31287008 343.    284
#> 9 Screen2YM155R2       29227176 320.    267
```

11.1.2 Essential genes

We want to get an idea of how strong the phenotypes in our screen are by comparing core-essential and non-essential screens. We first use the edgeR TMM-normalization to normalize samples and then calculate fold changes for each sample, comparing later time points to the T0 samples.

```
## make sure that the T0 sample is the first in the alphabet
data_split <- list()
data_split[[1]] <- counts %>%
  separate(sgRNA, c('symbol', 'sequence'), sep='_') %>%
  mutate(sample = ifelse(grepl('d0|T0|initial', sample), 'aa', sample)) %>%
  unite(sgRNA, symbol, sequence)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
## a function to calculate fold changes for each screen (replicates individually)
calc_fc <- function(df){
  ## matrix of counts
  cmat <- df %>% distinct() %>%
    spread(sample, Count) %>%
    data.frame() %>% `rownames<-`(NULL) %>%
    column_to_rownames('sgRNA')

  ## normalize with edgeR
  cmat <- calcNormFactors(DGEList(cmat)) %>% cpm(log=T)

  ## calculate fold changes and convert back to long df
  data.frame(cmat[, -1] - cmat[, 1]) %>% rownames_to_column('sgRNA') %>%
    tbl_df %>% gather(sample, log2fc, -sgRNA) %>%
    separate(sgRNA, c('symbol', 'sequence'), sep='_')
}

## apply function to calculate fc for each screen
data_norm <- data_split %>% map(calc_fc)
## exclude low coverage sgRNAs
data_norm <- data_norm %>% bind_rows() %>%
  inner_join(distinct(counts, sample)) %>%
  anti_join(data_split %>% bind_rows() %>%
    filter(sample == 'aa', Count < 30) %>%
    distinct(sgRNA) %>%
    separate(sgRNA, c('symbol', 'sequence'), sep='_'))
```

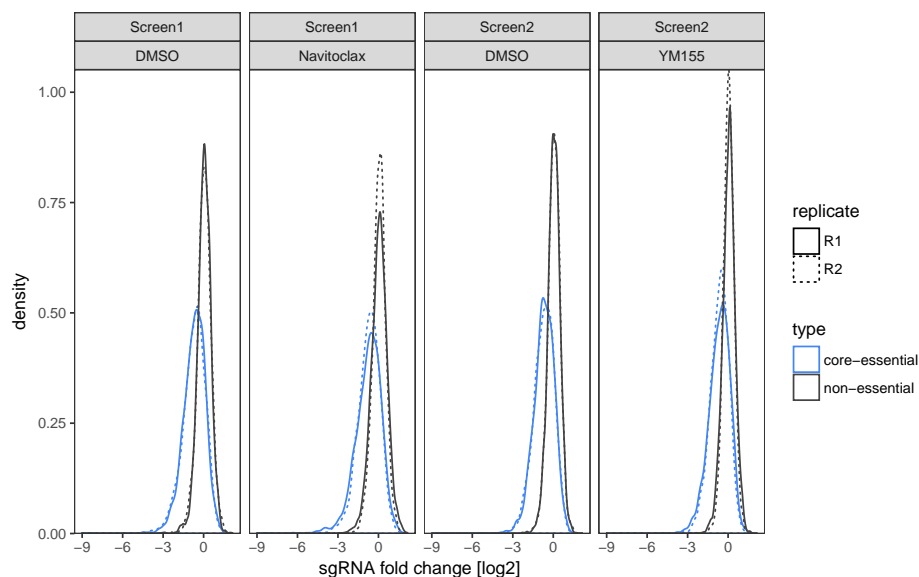
Next we load the core- and non-essential gene reference set defined in Hart et al., 2017.

```
## core-essential genes
data('ceg', package='CMSYM1552018')
## non-essential genes
data('neg', package='CMSYM1552018')
```

Now we can compare fold changes between core-essential and non-essential genes for each sample of our screen (Hart et al. 2017).

```
data_norm %>%
  mutate(type = ifelse(symbol %in% ceg, 'core-essential',
    ifelse(symbol %in% neg, 'non-essential', 'none'))) %>%
  filter(type != 'none') %>%
  extract(sample, c('screen', 'treatment', 'replicate'),
    regex = '(Screen\\d)(\\.+)(R\\d$)', remove=F) %>%
  ggplot(aes(log2fc, colour = type, linetype = replicate)) +
  geom_density() + facet_wrap(~screen + treatment, nrow=1) +
  scale_colour_manual(values = c('#4285f4', '#444444')) +
  scale_y_continuous(expand = c(0,0)) +
  xlab('sgRNA fold change [log2]') +
  theme(legend.position = 'bottom') +
  theme_bw() + theme(panel.grid = element_blank())
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



We can conclude that there are recognizable phenotypes, separating distributions of sgRNAs targeting core- and non-essential genes, respectively. This indicates that the experiment worked as intended.

11.2 Reproducibility

We next check how reproducible our phenotypes are across replicates. We plot scatter plots and annotate Pearson and Spearman correlation coefficients on log2-fold changes.

```
## disentangle sample column
ym_screen <- data_norm %>%
  extract(sample, c('screen', 'treatment', 'replicate'),
    regex='(Screen\\d)(DMSO|Navitoclax|YM155)(R\\d)')

reproducibility_plot <- function(df, sc, tr){
  df <- df %>% filter(screen == sc, treatment==tr)
  cors <- df %>% spread(replicate, log2fc) %>%
    summarise(pcc = cor(R1, R2, method='pearson'),
      scc = cor(R1, R2, method='spearman')) %>%
    unlist() %>% round(2)

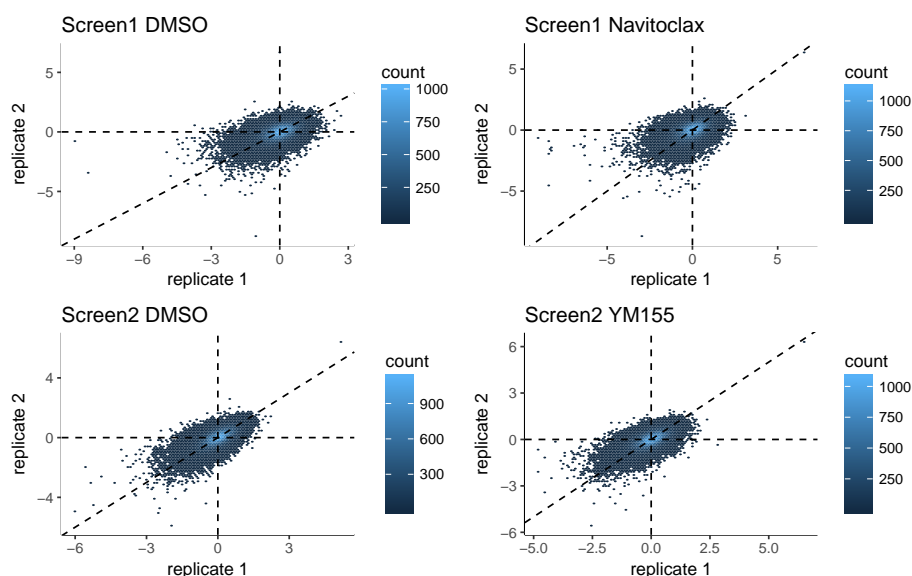
  df <- df %>% unite(experiment, screen, treatment) %>%
    spread(replicate, log2fc)

  df %>% ggplot(aes(R1, R2)) + geom_hex(bins=100) +
    geom_abline(linetype='dashed') +
    geom_hline(yintercept=0, linetype='dashed') +
    geom_vline(xintercept=0, linetype='dashed') +
    theme(legend.position='none') +
    xlab('replicate 1') + ylab('replicate 2') + ggtitle(paste(sc, tr)) +
    theme_classic()
}
```


Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
## generate plots using above function
rep_plots <- ym_screen %>% distinct(screen, treatment) %>%
  rowwise() %>% do(p = reproducibility_plot(ym_screen, .$screen, .$treatment)) %>%
  ungroup() %>% .$p

## draw to canvas
purrr::reduce(rep_plots, `+`)
```

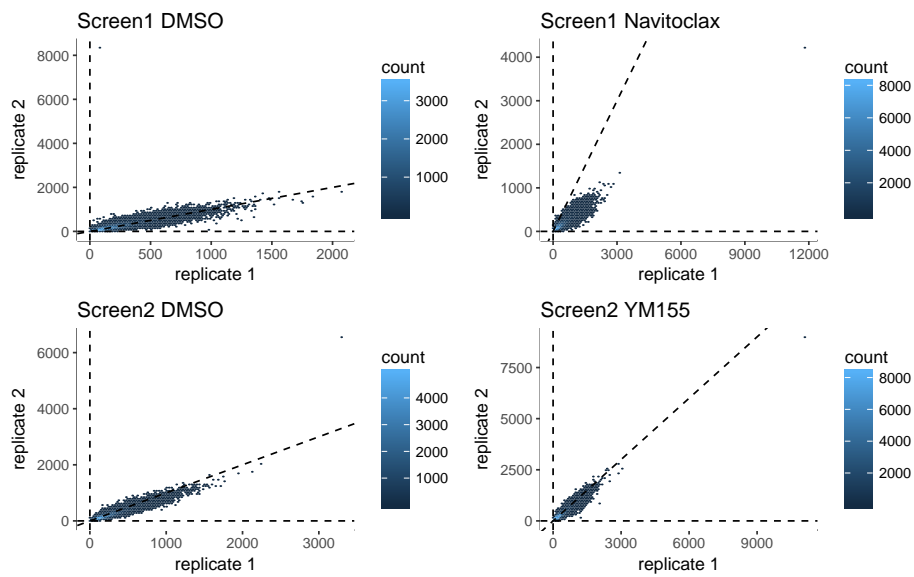


Correlation coefficients are in the same ballpark as previous studies ([based on a GenomeCRISPR meta analysis](#)). It might be worth pointing out that correlations based on log2 fold changes tend to always be considerably lower than coefficients based on normalized read counts that are also often presented. In our experience, however, fold change based correlation coefficients tend to be more 'honest' measures of reproducibility. For the sake of completeness we can also show similar scatter plots based on raw counts.

```
our_counts <- counts %>% filter(sample != 'd0screenlibrary') %>%
  extract(sample, c('screen', 'treatment', 'replicate'),
    regex='(Screen\\d)(.*) (R\\d)$') %>%
  ## we need to do this so the reproducibility plot will work, still counts though
  dplyr::select(log2fc=Count, everything())

our_counts %>% distinct(screen, treatment) %>%
  drop_na() %>% rowwise() %>%
  do(p = reproducibility_plot(our_counts, .$screen, .$treatment)) %>%
  ungroup() %>% .$p %>% purrr::reduce(`+`)
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



11.3 Hit calling

Next we would like to call hits between treated and untreated samples. For each sgRNA in the screen we would like to compare its abundance in the DMSO-treated sample compared to the sample that was treated with the drug of interest. For this purpose we use the MAGeCK software (version 0.5.7) (W. Li et al. 2014) that implements a negative binomial model to compare abundance of raw counts.

11.3.1 YM155

First we need to create an input file for mageck. This is a tab separated file that contains and sgRNA column, a gene column and additional columns for each sample.

```
## YM155
ym155_raw <- counts %>% separate(sgRNA, c('symbol', 'sequence'), sep='_') %>%
  filter(sequence %in% data_norm$sequence) %>%
  filter(sample != 'd0screenlibrary',
         grepl('Screen2', sample)) %>%
  dplyr::select(sequence, symbol, everything()) %>%
  spread(sample, Count)

## write mageck input file
ym155_raw %>% write_tsv('YM155.txt')
```

Next we run mageck on the sample file we just generated.

```
system('mageck test -k YM155.txt -t 2,3 -c 0,1 -n YM155')
```

We then load the results back into R.

```
data('ym_results', package='CMSYM1552018')
```

11.3.2 Navitoclax

For the Navitoclax screen, as above, we need to create an input file for mageck. This is a tab separated file that contains an sgRNA column, a gene column and additional columns for each sample.

```
## YM155
navito_raw <- counts %>% separate(sgRNA, c('symbol', 'sequence'), sep='_') %>%
  filter(sequence %in% data_norm$sequence) %>%
  filter(sample != 'd0screenlibrary',
         grepl('Screen1', sample)) %>%
  dplyr::select(sequence, symbol, everything()) %>%
  spread(sample, Count)

## write mageck input file
navito_raw %>% write_tsv('Navitoclax.txt')
```

Next we run MAGeCK on the sample file we just generated.

```
system('mageck test -k Navitoclax.txt -t 2,3 -c 0,1 -n Navitoclax')
```

We can then read the results file back into R.

```
data('navito_results', package='CMSYM1552018')
```

11.4 Visualization of results

We can now visualize the results in the form of two volcano plots. Mageck performs two sets of tests: one for positive and one for negative enrichment. Separate plots have to be generated for each of those.

11.4.1 YM155

We start with the negative enrichment.

```
mageck_volcano <- function(df, type='negative', highlight=c()){
  ## hit colour
  hit_col <- ifelse(type == 'negative', '#4285f4', '#e41a1c')

  ## fdr 20 cutoff
  fdr20 <- df %>% filter(fdr > 0.2) %>%
    arrange(p.value) %>% .$p.value %>% .[1] %>% log10() %>% `*(-1)`

  ## draw volcano
  df %>% mutate(label = ifelse(symbol %in% highlight, symbol, ''),
                 colour = ifelse(fdr < 0.2 & abs(log2fc) > 0.25,
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
      'hit', 'none'),
      type = type) %>%
  filter(ifelse(type == 'negative', log2fc < 0, log2fc > 0)) %>%
  ggplot(aes(log2fc, -log10(p.value))) +
  geom_point(aes(colour = colour)) +
  ggrepel::geom_text_repel(aes(label=label)) +
  geom_vline(xintercept=0, linetype = 'dashed') +
  geom_hline(yintercept = fdr20, linetype = 'dashed') +
  scale_colour_manual(values = c(hit_col, '#cccccc')) +
  theme(legend.position = 'none') +
  xlab('Fold change [log2]') + ylab('P-value [-log10]') +
  theme_classic()
}

## list of candidates selected form the Screen
candidates_neg <- c('SMAGP')
candidates_pos <- c('SLC35F2', 'CCDC22', 'SNX17', 'KIAA1033')

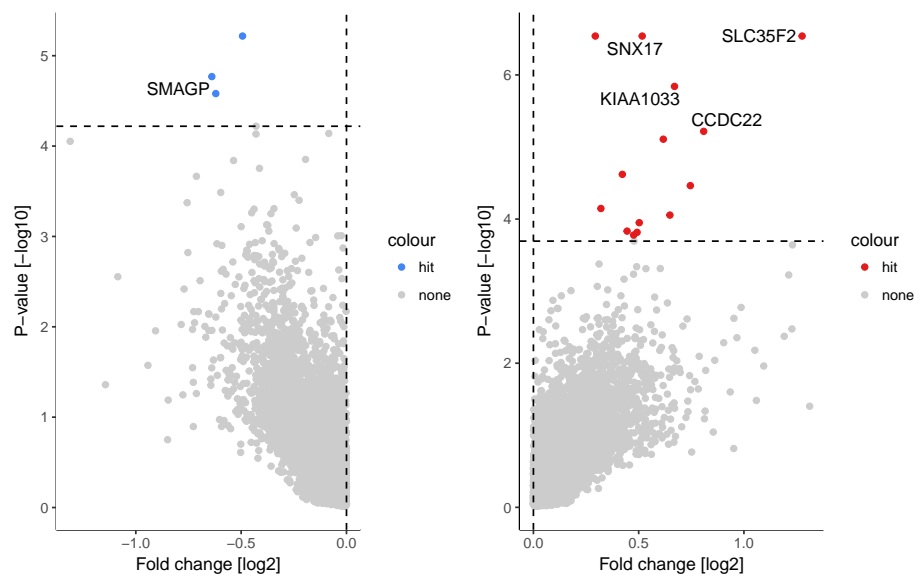
## generate plot
ym_negative <- ym_results %>%
  dplyr::select(symbol=id, log2fc = `neg|lfc`,
                p.value = `neg|p-value`, fdr = `neg|fdr`) %>%
  mageck_volcano(type = 'negative', highlight = candidates_neg)
```

We now add a plot for the positive enrichment and combine them using `patchwork`.

```
## generate plot
ym_positive <- ym_results %>%
  dplyr::select(symbol=id, log2fc = `pos|lfc`,
                p.value = `pos|p-value`, fdr = `pos|fdr`) %>%
  mageck_volcano(type = 'positive', highlight = candidates_pos)

## combine
ym_negative + ym_positive
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155



11.4.2 Navitoclax

We proceed as above with the Navitoclax Screen.

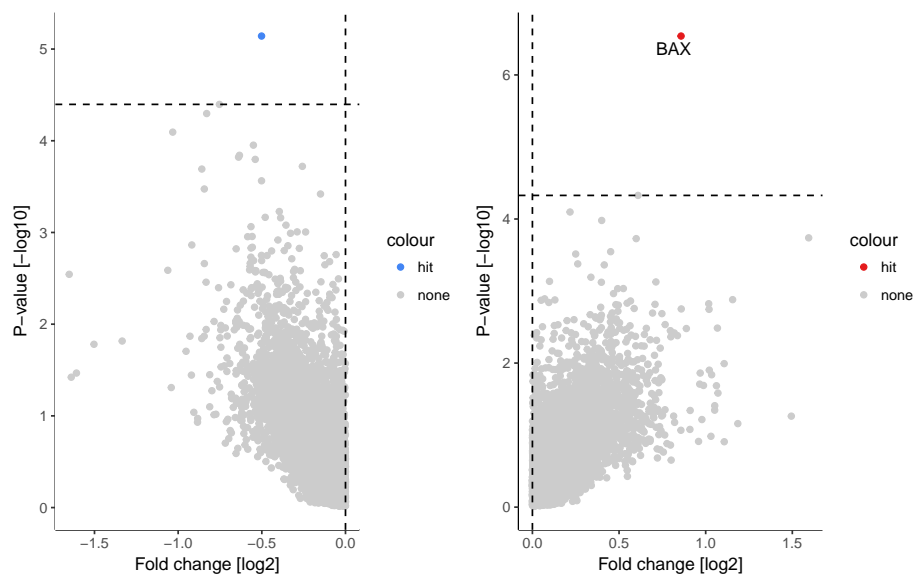
```
candidates_pos <- c('BAX')

navito_positive <- navito_results %>%
  dplyr::select(symbol=id, log2fc = `pos|lfc`,
                p.value = `pos|p-value`, fdr = `pos|fdr`) %>%
  mageck_volcano(type = 'positive', highlight = candidates_pos)
```

```
candidates_neg <- c()

navito_negative <- navito_results %>%
  dplyr::select(symbol=id, log2fc = `neg|lfc`,
                p.value = `neg|p-value`, fdr = `neg|fdr`) %>%
  mageck_volcano(type = 'negative', highlight = candidates_neg)

## plot
navito_negative + navito_positive
```



12 Session info

```
sessionInfo()
#> R version 3.4.1 (2017-06-30)
#> Platform: x86_64-apple-darwin15.6.0 (64-bit)
#> Running under: macOS Sierra 10.12.6
#>
#> Matrix products: default
#> BLAS: /Library/Frameworks/R.framework/Versions/3.4/Resources/lib/libRblas.0.dylib
#> LAPACK: /Library/Frameworks/R.framework/Versions/3.4/Resources/lib/libRlapack.dylib
#>
#> locale:
#> [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
#>
#> attached base packages:
#> [1] stats4      parallel    stats       graphics    grDevices   utils       datasets
#> [8] methods     base
#>
#> other attached packages:
#> [1] hexbin_1.27.2      bindrcpp_0.2.2      edgeR_3.20.9
#> [4] Organism.dplyr_1.6.2 AnnotationFilter_1.2.0 G0.db_3.5.0
#> [7] AnnotationDbi_1.40.0 IRanges_2.12.0      S4Vectors_0.16.0
#> [10] fgsea_1.4.1        Rcpp_0.12.16        limma_3.34.9
#> [13] lumi_2.30.0        readxl_1.1.0        patchwork_0.0.1
#> [16] perm_1.0-0.0        ggrepel_0.7.0       CMSclassifier_1.0.0
#> [19] randomForest_4.6-14 ggsignif_0.4.0       reshape2_1.4.3
#> [22] impute_1.52.0       sva_3.26.0          BiocParallel_1.12.0
#> [25] genefilter_1.60.0   mgcv_1.8-23         nlme_3.1-137
#> [28] openxlsx_4.0.17     pheatmap_1.0.8      preprocessCore_1.40.0
#> [31] GEOquery_2.46.15    affy_1.56.0         Biobase_2.38.0
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
#> [34] BiocGenerics_0.24.0      forcats_0.3.0           stringr_1.3.0
#> [37] dplyr_0.7.4             purrr_0.2.4            readr_1.1.1
#> [40] tidyr_0.8.0             tibble_1.4.2           ggplot2_2.2.1.9000
#> [43] tidyverse_1.2.1         BiocStyle_2.6.1
#>
#> loaded via a namespace (and not attached):
#> [1] backports_1.1.2
#> [2] fastmatch_1.1-0
#> [3] BiocFileCache_1.2.3
#> [4] plyr_1.8.4
#> [5] lazyeval_0.2.1
#> [6] splines_3.4.1
#> [7] GenomeInfoDb_1.14.0
#> [8] digest_0.6.15
#> [9] foreach_1.4.4
#> [10] BiocInstaller_1.28.0
#> [11] htmltools_0.3.6
#> [12] magrittr_1.5
#> [13] memoise_1.1.0
#> [14] Biostrings_2.46.0
#> [15] annotate_1.56.2
#> [16] modelr_0.1.1
#> [17] matrixStats_0.53.1
#> [18] siggenes_1.52.0
#> [19] prettyunits_1.0.2
#> [20] colorspace_1.3-2
#> [21] rappdirs_0.3.1
#> [22] blob_1.1.1
#> [23] rvest_0.3.2
#> [24] haven_1.1.1
#> [25] xfun_0.1
#> [26] crayon_1.3.4
#> [27] RCurl_1.95-4.10
#> [28] jsonlite_1.5
#> [29] bindr_0.1.1
#> [30] survival_2.42-3
#> [31] iterators_1.0.9
#> [32] glue_1.2.0
#> [33] registry_0.5
#> [34] gtable_0.2.0
#> [35] zlibbioc_1.24.0
#> [36] XVector_0.18.0
#> [37] TxDb.Hsapiens.UCSC.hg38.knownGene_3.4.0
#> [38] DelayedArray_0.4.1
#> [39] scales_0.5.0.9000
#> [40] DBI_1.0.0
#> [41] rngtools_1.2.4
#> [42] xtable_1.8-2
#> [43] progress_1.1.2
#> [44] bumphunter_1.20.0
#> [45] foreign_0.8-70
```

Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
#> [46] bit_1.1-12
#> [47] mclust_5.4
#> [48] httr_1.3.1
#> [49] RColorBrewer_1.1-2
#> [50] pkgconfig_2.0.1
#> [51] reshape_0.8.7
#> [52] XML_3.98-1.11
#> [53] dbplyr_1.2.1
#> [54] utf8_1.1.3
#> [55] locfit_1.5-9.1
#> [56] labeling_0.3
#> [57] tidyselect_0.2.4
#> [58] rlang_0.2.0.9001
#> [59] munsell_0.4.3
#> [60] cellranger_1.1.0
#> [61] tools_3.4.1
#> [62] cli_1.0.0
#> [63] RSQLite_2.1.1
#> [64] broom_0.4.4
#> [65] evaluate_0.10.1
#> [66] yaml_2.1.19
#> [67] org.Hs.eg.db_3.5.0
#> [68] knitr_1.20
#> [69] bit64_0.9-7
#> [70] beanplot_1.2
#> [71] methylumi_2.24.1
#> [72] doRNG_1.6.6
#> [73] nor1mix_1.2-3
#> [74] pracma_2.1.4
#> [75] xml2_1.2.0
#> [76] biomaRt_2.34.2
#> [77] compiler_3.4.1
#> [78] rstudioapi_0.7
#> [79] affyio_1.48.0
#> [80] stringi_1.2.2
#> [81] GenomicFeatures_1.30.3
#> [82] minfi_1.24.0
#> [83] lattice_0.20-35
#> [84] Matrix_1.2-14
#> [85] psych_1.8.4
#> [86] multtest_2.34.0
#> [87] pillar_1.2.2
#> [88] data.table_1.10.4-3
#> [89] bitops_1.0-6
#> [90] rtracklayer_1.38.3
#> [91] GenomicRanges_1.30.3
#> [92] R6_2.2.2
#> [93] bookdown_0.7
#> [94] RMySQL_0.10.14
#> [95] gridExtra_2.3
#> [96] KernSmooth_2.23-15
```


Multi-omics integration identifies a selective vulnerability of colorectal cancer subtypes to YM155

```
#> [97] nleqslv_3.3.1
#> [98] codetools_0.2-15
#> [99] MASS_7.3-50
#> [100] assertthat_0.2.0
#> [101] SummarizedExperiment_1.8.1
#> [102] openssl_1.0.1
#> [103] pkgmaker_0.22
#> [104] rprojroot_1.3-2
#> [105] withr_2.1.2
#> [106] GenomicAlignments_1.14.2
#> [107] Rsamtools_1.30.0
#> [108] mnormt_1.5-5
#> [109] GenomeInfoDbData_1.0.0
#> [110] hms_0.4.2
#> [111] quadprog_1.5-5
#> [112] grid_3.4.1
#> [113] base64_2.0
#> [114] rmarkdown_1.9
#> [115] illuminaio_0.20.0
#> [116] lubridate_1.7.4
```

References

- Barretina, Jordi, Giordano Caponigro, Nicolas Stransky, Kavitha Venkatesan, Adam A Margolin, Sungjoon Kim, Christopher J Wilson, et al. 2012. "The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity." *Nature* 483 (7391): 603–7. doi:[10.1038/nature11003](https://doi.org/10.1038/nature11003).
- Du, Pan, Warren A Kibbe, and Simon M Lin. 2008. "lumi: a pipeline for processing Illumina microarray." *Bioinformatics (Oxford, England)* 24 (13): 1547–8. doi:[10.1093/bioinformatics/btn224](https://doi.org/10.1093/bioinformatics/btn224).
- Edgar, Ron, Michael Domrachev, and Alex E Lash. 2002. "Gene Expression Omnibus: NCBI gene expression and hybridization array data repository." *Nucleic Acids Research* 30 (1): 207–10. <http://www.ncbi.nlm.nih.gov/pubmed/11752295> <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC99122>.
- Forbes, Simon A, David Beare, Prasad Gunasekaran, Kenric Leung, Nidhi Bindal, Harry Boutselakis, Minjie Ding, et al. 2015. "COSMIC: exploring the world's knowledge of somatic mutations in human cancer." *Nucleic Acids Research* 43 (Database issue): D805–11. doi:[10.1093/nar/gku1075](https://doi.org/10.1093/nar/gku1075).
- Frejno, Martin, Riccardo Zenezini Chiozzi, Mathias Wilhelm, Heiner Koch, Runsheng Zheng, Susan Klaeger, Benjamin Ruprecht, et al. 2017. "Pharmacoproteomic characterisation of human colon and rectal cancer." *Molecular Systems Biology* 13 (11): 951. doi:[10.15252/msb.20177701](https://doi.org/10.15252/msb.20177701).
- Garnett, Mathew J., Elena J. Edelman, Sonja J. Heidorn, Chris D. Greenman, Anahita Dastur, King Wai Lau, Patricia Greninger, et al. 2012. "Systematic identification of genomic markers of drug sensitivity in cancer cells." *Nature* 483 (7391): 570–75. doi:[10.1038/nature11005](https://doi.org/10.1038/nature11005).
- Gautier, L., L. Cope, B. M. Bolstad, and R. A. Irizarry. 2004. "affy-analysis of Affymetrix GeneChip data at the probe level." *Bioinformatics* 20 (3): 307–15. doi:[10.1093/bioinformatics/btg405](https://doi.org/10.1093/bioinformatics/btg405).
- Guinney, Justin, Rodrigo Dienstmann, Xin Wang, Aurélien de Reyniès, Andreas Schlicker, Charlotte Sonesson, Laetitia Marisa, et al. 2015. "The consensus molecular subtypes of colorectal cancer." *Nature Medicine* 21 (11): 1350–6. doi:[10.1038/nm.3967](https://doi.org/10.1038/nm.3967).
- Hart, Traver, Amy Hin Yan Tong, Katie Chan, Jolanda Van Leeuwen, Ashwin Seetharaman, Michael Aregger, Megha Chandrashekhar, et al. 2017. "Evaluation and Design of Genome-Wide CRISPR/SpCas9 Knockout Screens." *G3 (Bethesda, Md.)* 7 (8): 2719–27. doi:[10.1534/g3.117.041277](https://doi.org/10.1534/g3.117.041277).
- Hollingshead, Melinda G, Luke H Stockwin, Sergio Y Alcoser, Dianne L Newton, Benjamin C Orsburn, Carrie A Bonomi, Suzanne D Borgel, et al. 2014. "Gene expression profiling of 49 human tumor xenografts from in vitro culture through multiple in vivo passages—strategies for data mining in support of therapeutic studies." *BMC Genomics* 15 (1): 393. doi:[10.1186/1471-2164-15-393](https://doi.org/10.1186/1471-2164-15-393).
- Iorio, Francesco, Theo A. Knijnenburg, Daniel J. Vis, Graham R. Bignell, Michael P. Menden, Michael Schubert, Nanne Aben, et al. 2016. "A Landscape of Pharmacogenomic Interactions in Cancer." *Cell* 166 (3): 740–54. doi:[10.1016/j.cell.2016.06.017](https://doi.org/10.1016/j.cell.2016.06.017).
- Irizarry, R. A., Bridget Hobbs, Francois Collin, Yasmin D Beazer-Barclay, Kristen J Antonellis, Uwe Scherf, and Terence P Speed. 2003. "Exploration, normalization, and summaries of high density oligonucleotide array probe level data." *Biostatistics* 4 (2): 249–64. doi:[10.1093/biostatistics/4.2.249](https://doi.org/10.1093/biostatistics/4.2.249).
- Li, Wei, Han Xu, Tengfei Xiao, Le Cong, Michael I Love, Feng Zhang, Rafael A Irizarry, Jun S Liu, Myles Brown, and X Shirley Liu. 2014. "MAGeCK enables robust identification of essential genes from genome-scale CRISPR/Cas9 knockout screens." *Genome Biology* 15 (12): 554. doi:[10.1186/s13059-014-0554-4](https://doi.org/10.1186/s13059-014-0554-4).
- Linnekamp, Janneke F., Sander R. van Hooff, Pramudita R. Prasetyanti, Raju Kandimalla, Joyce Y. Buikhuizen, Evelyn Forster, Prachanthi Ramesh, et al. 2018. "Consensus molecular