# Descriptive statistics

```r
library(magrittr)
library(ggplot2)
library(gridExtra)
library(ggcorrplot)

source("../R/visualizations.R")
source("../R/feature_definitions.R")

training_set <- read.csv("../preprocessed_training_data.csv", row.names = 1, as.is = TRUE)
outcome <- read.csv("../training_outcomes.csv", row.names = 1)[,1]

stopifnot(row.names(training_set) == row.names(outcome))

features <- colnames(training_set)
```
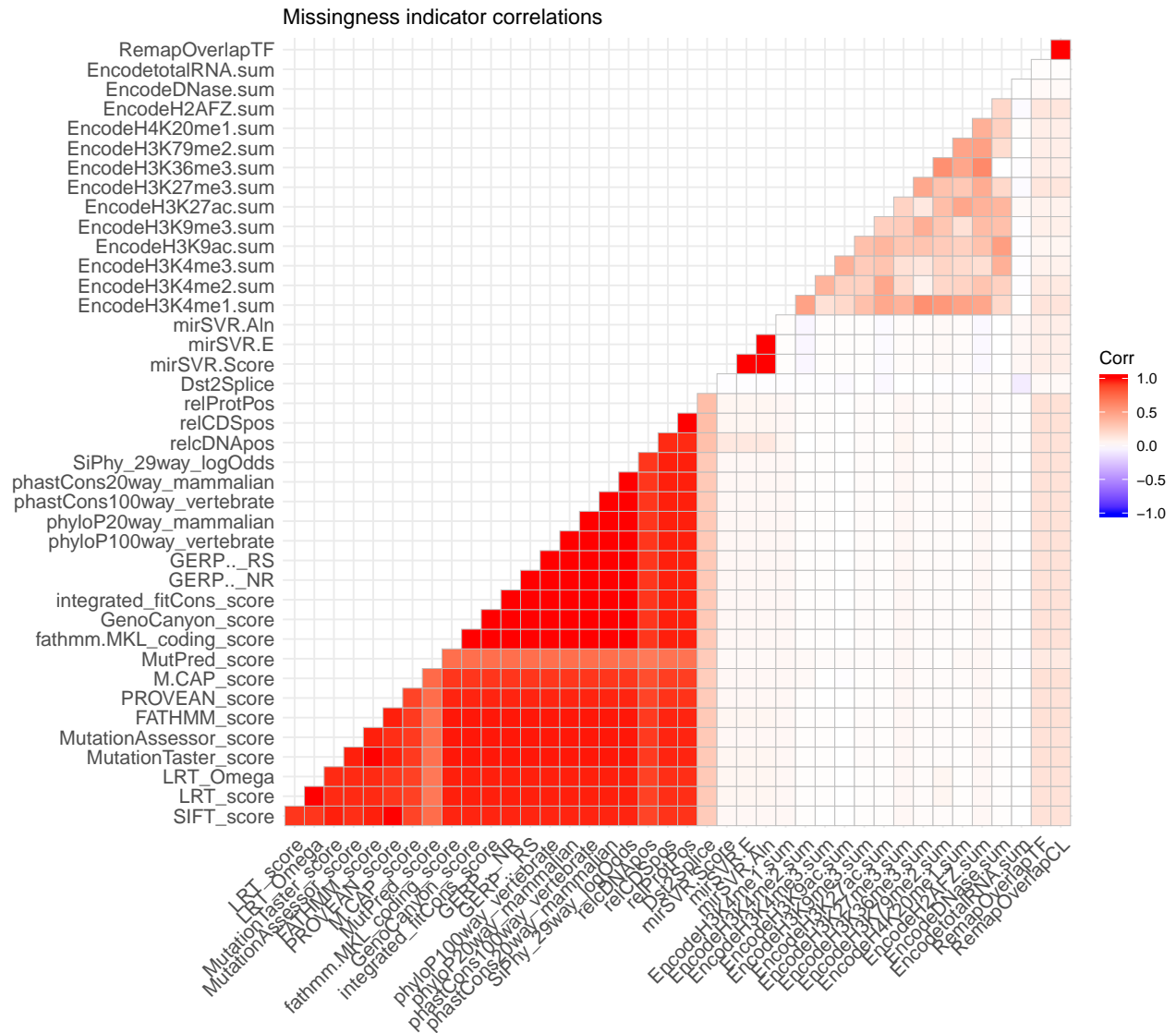
## Correlations

Plot correlation matrices of missingness indicators against missingness indicators, observed values against observed values, and missingness indicators against observed values.

```r
positive_data <- training_set[outcome == "positive", ]
negative_data <- training_set[outcome == "negative", ]

# Missingness indicator correlations
plot_missingness_correlations(training_set, numeric_features, "Missingness indicator correlations")
```
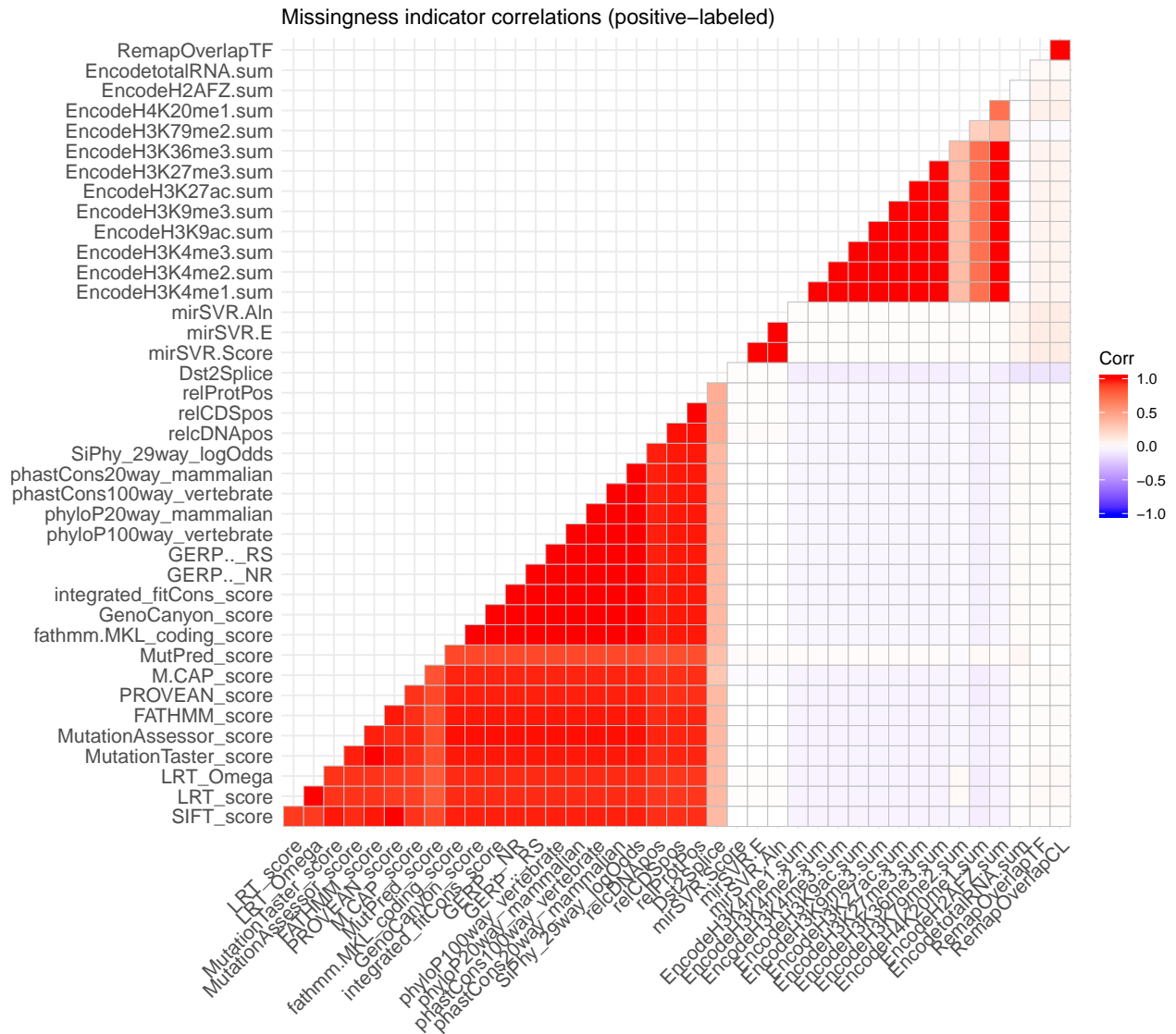
```
## Warning in cor(miss_data[, features]): the standard deviation is zero
```

Missingness indicator correlations

```
plot_missingness_correlations(positive_data, numeric_features, "Missingness indicator correlations (pos
```

```
## Warning in cor(miss_data[, features]): the standard deviation is zero
```
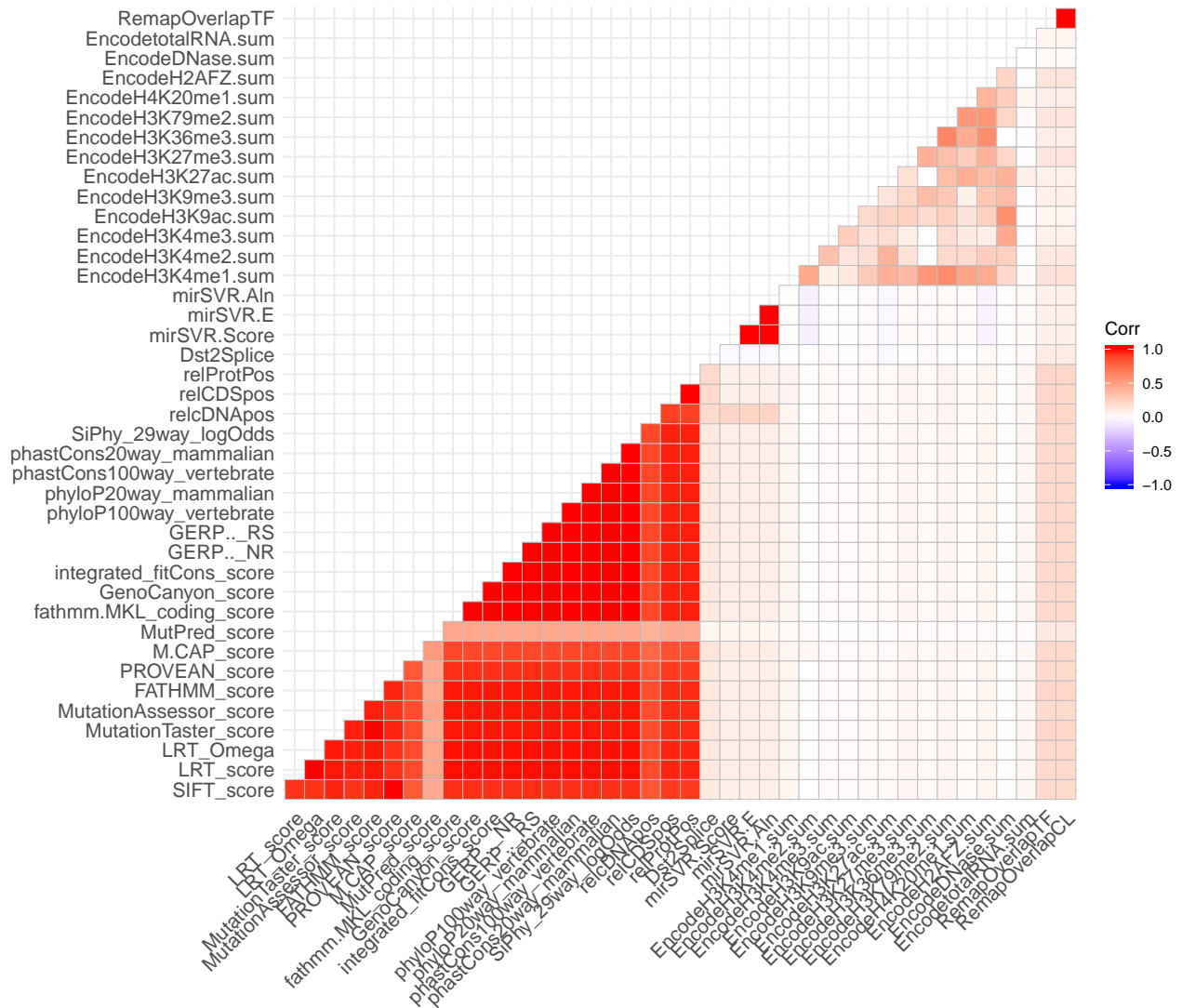
Missingness indicator correlations (positive-labeled)

```
plot_missingness_correlations(negative_data, numeric_features, "Missingness indicator correlations (nega
```

```
## Warning in cor(miss_data[, features]): the standard deviation is zero
```
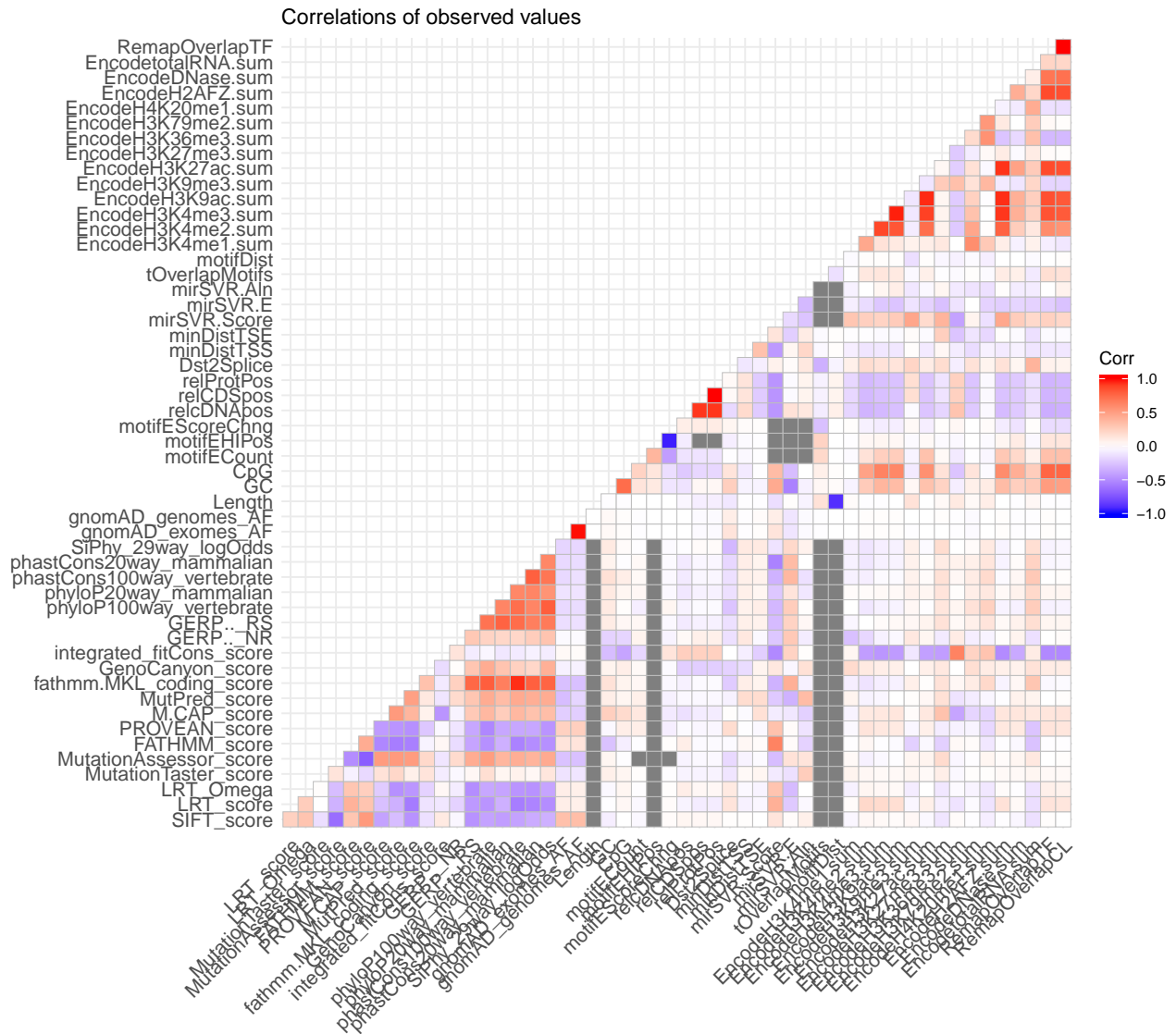
Missingness indicator correlations (negative–labeled)

```
# Observed value correlations
plot_observed_correlations(training_set, numeric_features, "Correlations of observed values")

## Warning in cor(data[, features], use = "pairwise.complete.obs"): the
## standard deviation is zero
```

Correlations of observed values

```
plot_observed_correlations(positive_data, numeric_features, "Correlations of observed values (positive-
```

```
## Warning in cor(data[, features], use = "pairwise.complete.obs"): the
## standard deviation is zero
```

Correlations of observed values (positive–labeled)

```
plot_observed_correlations(negative_data, numeric_features, "Correlations of observed values (negative-
```

```
## Warning in cor(data[, features], use = "pairwise.complete.obs"): the
## standard deviation is zero
```
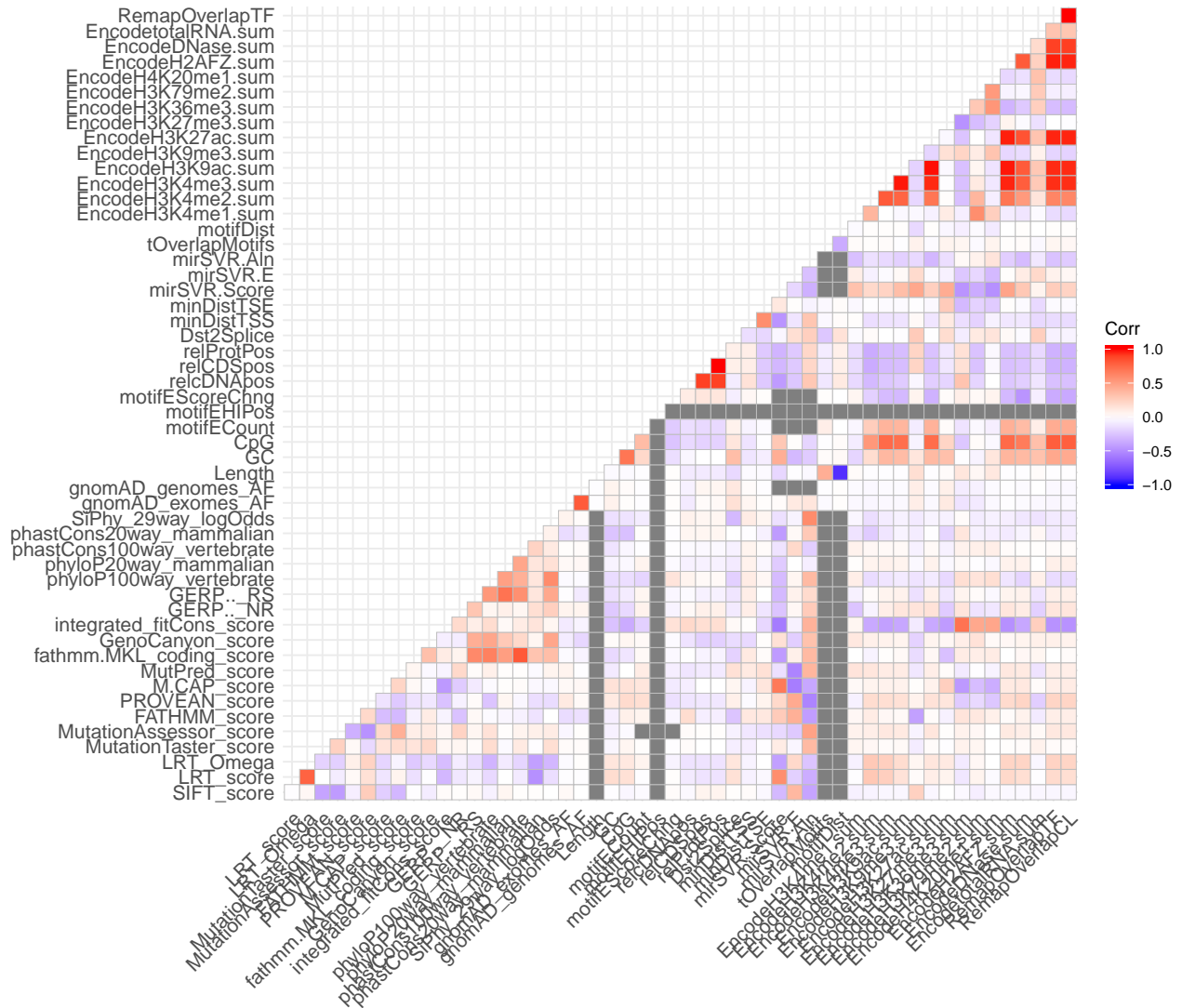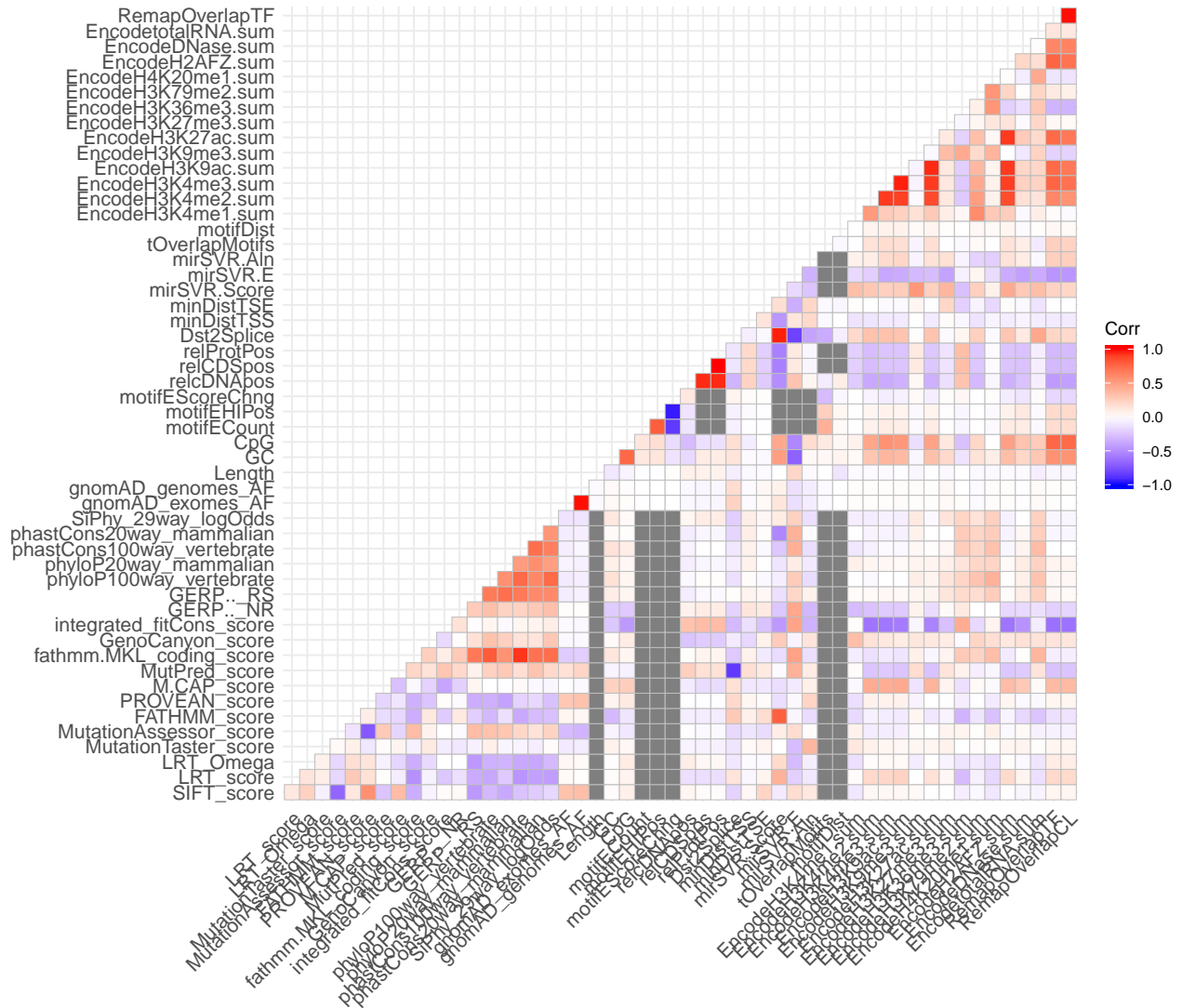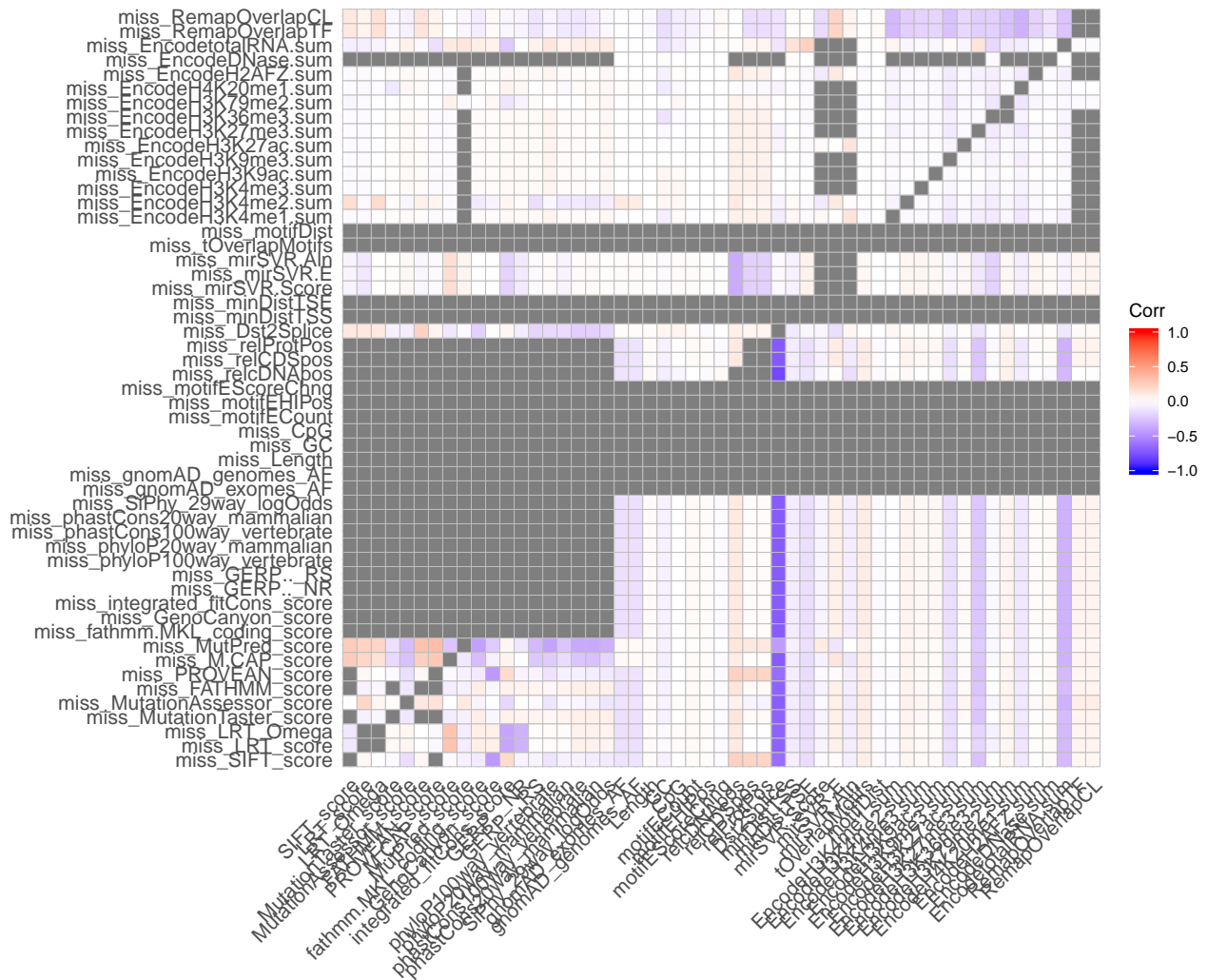
Correlations of observed values (negative−labeled)

```
# Missingness vs. observed correlations
plot_missingness_vs_observed_correlations(training_set, numeric_features, "Missingness correlations vs.
```

```
## Warning in cor(data, miss_data, use = "pairwise.complete.obs"): the
## standard deviation is zero
```
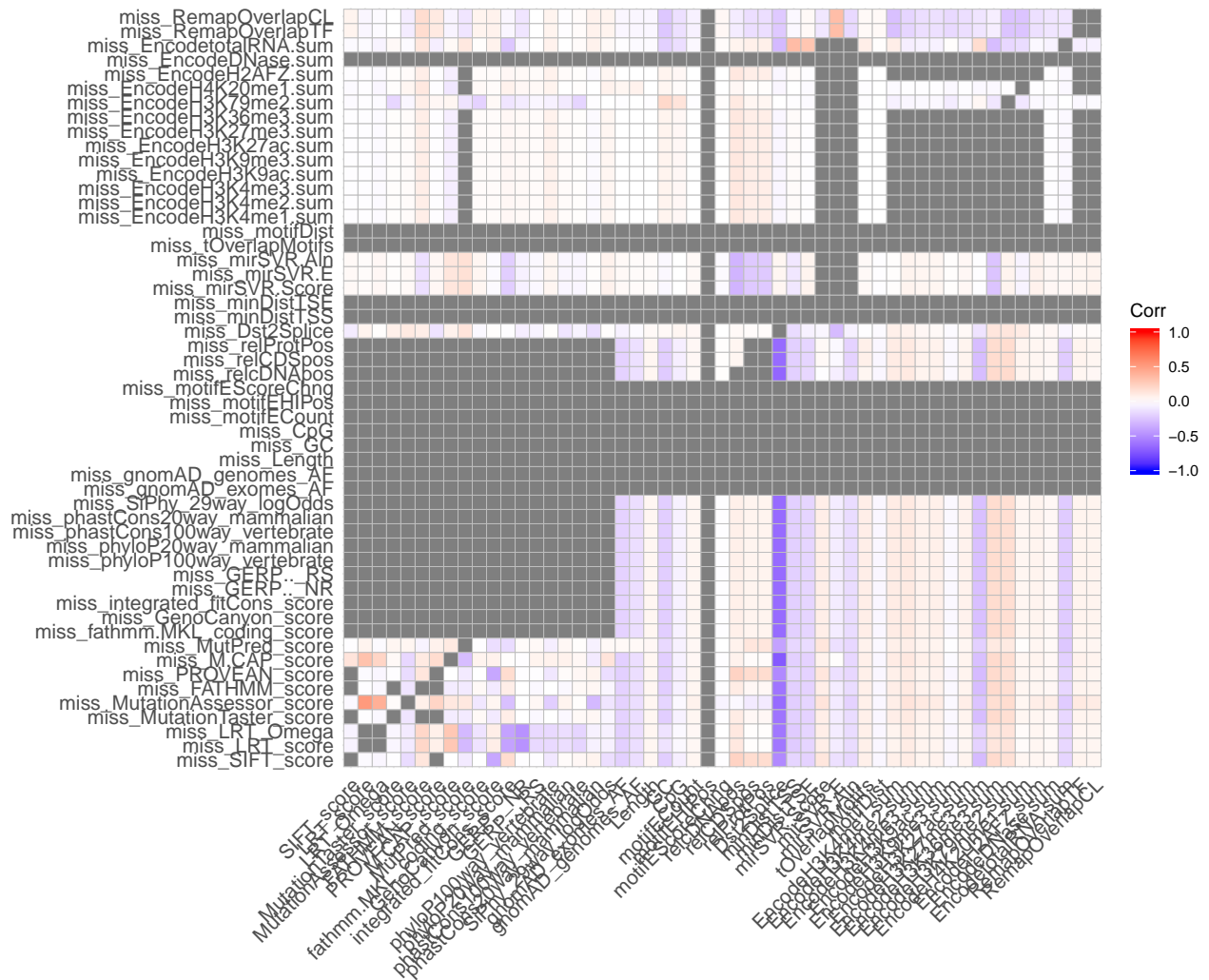
Missingness correlations vs. observed values

```
plot_missingness_vs_observed_correlations(positive_data, numeric_features, "Missingness correlations vs
```

```
## Warning in cor(data, miss_data, use = "pairwise.complete.obs"): the
## standard deviation is zero
```

Missingness correlations vs. observed values (positive–labeled)

```
plot_missingness_vs_observed_correlations(negative_data, numeric_features, "Missingness correlations vs
```

```
## Warning in cor(data, miss_data, use = "pairwise.complete.obs"): the
## standard deviation is zero
```
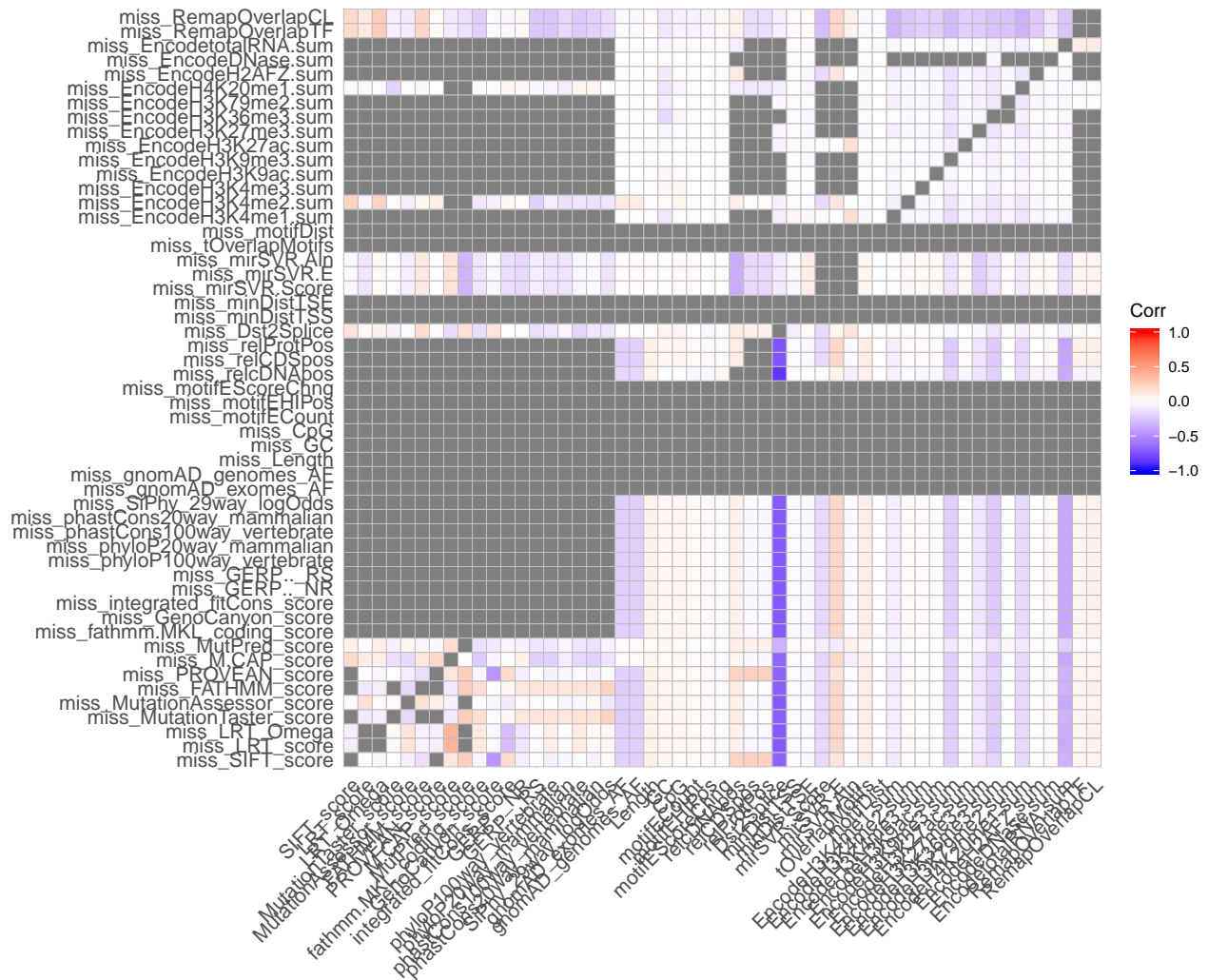
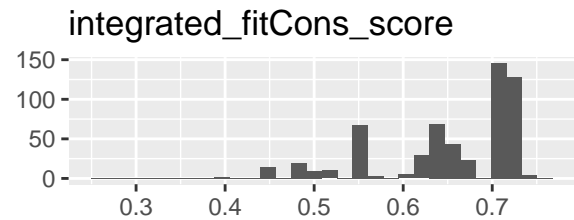Missingness correlations vs. observed values (negative–labeled)

## Feature value distributions

Next, plot distributions of each feature. Are they normal or linear?

```r
feature_distribution_plots <- lapply(numeric_features,
                                     function(column) {
                                       ggplot2::quickplot(
                                         na.omit(training_set[,column]),
                                         main = column,
                                         xlab = "",
                                         bins = 30
                                       )
                                     })

marrangeGrob(
  ncol = 2, nrow = 3,
  grobs = feature_distribution_plots
)
```

## SIFT_score

## MutationTaster_score

## LRT_score

## MutationAssessor_score

## LRT_Omega

## FATHMM_score

## PROVEAN_score

## fathmm.MKL_coding_score

## M.CAP_score

## GenoCanyon_score

## MutPred_score

## integrated_fitCons_score

## GERP.._NR

## phyloP20way_mammalian

## GERP.._RS

## phastCons100way_vertebrate

## phyloP100way_vertebrate

## phastCons20way_mammalian

## SiPhy_29way_logOdds

## Length

## gnomAD_exomes_AF

## GC

## gnomAD_genomes_AF

## CpG

## motifECount

## relcDNApos

## motifEHIPos

## relCDSpos

## motifEScoreChng

## relProtPos

## Dst2Splice

## mirSVR.Score

## minDistTSS

## mirSVR.E

## minDistTSE

## mirSVR.Aln

## tOverlapMotifs

## EncodeH3K4me2.sum

## motifDist

## EncodeH3K4me3.sum

## EncodeH3K4me1.sum

## EncodeH3K9ac.sum

## EncodeH3K9me3.sum

## EncodeH3K36me3.sum

## EncodeH3K27ac.sum

## EncodeH3K79me2.sum

## EncodeH3K27me3.sum

## EncodeH4K20me1.sum

## EncodeH2AFZ.sum

## RemapOverlapTF

## EncodeDNase.sum

## RemapOverlapCL

## EncodetotalRNA.sum

They are not, and thus it might be worth considering data transformations. In the case of random forest, however, monotone transformations should have no effect.

## Categorical level occurence counts

Print (one-dimensional) contingency tables, i.e. occurence counts of each level of categorical variables.

```r
for (cat_feat in categorical_features) {
  table(training_set[, cat_feat, drop = FALSE], dnn = cat_feat, useNA = "always") %>% as.data.frame %>%
}
```

```
##   LRT_pred Freq
## 1        D  382
## 2        N  162
## 3        U    4
## 4     <NA> 2044
##   Dst2SplType Freq
## 1    ACCEPTOR  115
## 2       DONOR  139
## 3        <NA> 2338
##       Consequence.x Freq
## 1         3PRIME_UTR   22
## 2         5PRIME_UTR   17
## 3         DOWNSTREAM  458
## 4           INTRONIC  794
## 5     NON_SYNONYMOUS  580
## 6   NONCODING_CHANGE   13
## 7        SPLICE_SITE   73
## 8           UPSTREAM  635
```

15

```
## 9              <NA>    0
```

## Heatmap of feature missingness against consequence

It is likely that missing values are more or less common in some variables depending on the predicted consequence. This can be visualized by a heatmap:

```
missing_value_sum_per_consequence <- lapply(training_set[, c(numeric_features, categorical_features), d
                                            function(column) {
                                              sapply(
                                                split(is.na(column), training_set$Consequence.x),
                                                sum
                                              )
                                            })
missing_value_sum_per_consequence %<>% data.frame %>% as.matrix
heatmap(missing_value_sum_per_consequence)
```



Non-synonymous variants have much less missingness in certain variables and more in others (as expected).

## Compute number of observed missingness patterns

```
missingness_patterns <- training_set[, c(numeric_features, categorical_features)] %>% is.na
unique_missingness_patterns <- missingness_patterns %>% unique
num_missingness_patterns <- unique_missingness_patterns %>% nrow
print(paste(num_missingness_patterns, "out of", 2^length(c(numeric_features, categorical_features)), "po
```

```
## [1] "118 out of 72057594037927936 possible missingness patterns."
```

```r
missingness_pattern_factor <- apply(missingness_patterns, MARGIN = 1, function(x) paste0(as.integer(x),
rows_per_missingness_pattern <- table(missingness_pattern_factor)
rows_per_missingness_pattern <- rows_per_missingness_pattern %>% as.data.frame
rows_per_missingness_pattern[order(rows_per_missingness_pattern$Freq, decreasing = TRUE),]
```

```
##                            missingness_pattern_factor Freq
## 85  1111111111111111111000000001111001110000000000000000110 1042
## 86  1111111111111111111000000001111001110000000000000011110  629
## 7   0000000000000000000000000000001001110000000000000000010  151
## 18  0000000010000000000000000000001001110000000000000000010   86
## 81  1111111111111111111000000001111000000000000000000000110   65
## 3   0000000000000000000000000000000001110000000000000000000   54
## 74  1111111111111111111000000001110001110000000000000000100   44
## 87  1111111111111111111000000001111001110000000000000100110   32
## 8   0000000000000000000000000000001001110000000000000011010   28
## 32  0000000110000000000000000000001001110000000000000000010   24
## 12  0000000010000000000000000000001001110000000000000000000   21
## 19  0000000010000000000000000000001001110000000000000011010   21
## 4   0000000000000000000000000000000001110000000000000011000   19
## 75  1111111111111111111000000001110001110000000000000011100   19
## 88  1111111111111111111000000001111001110000000000000111110   19
## 71  1111111111111111111000000001111001110000000000000000110   18
## 6   0000000000000000000000000001000000000000000000000000010   16
## 62  1111111111111111111000000000000001110000000000000000100   14
## 33  0000000110000000000000000000001001110000000000000011010   12
## 68  1111111111111111111000000001110000000000000000000000110   12
## 96  1111111111111111111000000001111001110000000100000011110   12
## 17  0000000010000000000000000000001000000000000000000000010   11
## 52  1000001010000000000000000000001001110000000000000000010   11
## 69  1111111111111111111000000001110000000000000000000011110    9
## 82  1111111111111111111000000001111000000000000000000011110    9
## 9   0000000000000000000000000000001001110000000000000100010    8
## 45  0110000000000000000000000000001001110000000000000000110    8
## 46  0110000000000000000000000000001001110000000000000011110    7
## 77  1111111111111111111000000001110001110000000000000111100    7
## 25  0000000010000000000000000000000001110000000000000000000    6
## 48  0110000010000000000000000000001001110000000000000000110    6
## 50  1000001010000000000000000000000001110000000000000000000    6
## 63  1111111111111111111000000000000001110000000000000011100    6
## 76  1111111111111111111000000001110001110000000000000100100    6
## 26  0000000010000000000000000000000001110000000000000011000    5
## 56  1001011000000000000000000000001001110000000000000000010    5
## 72  1111111111111111111000000001111001110000000000000011110    5
## 89  1111111111111111111000000001111001110000000000010011110    5
## 13  0000000010000000000000000000001001110000000000000011000    4
## 39  0000100010000000000000000000001001110000000000000000010    4
## 57  1001011000000000000000000000001001110000000000000011010    4
## 66  1111111111111111111000000001100011100000000000000000100    4
## 92  1111111111111111111000000001111001110000000001000000110    4
## 106 1111111111111111111000000001111001110010000000000011110    4
## 5   0000000000000000000000000000000001110000000000000100000    3
## 11  0000000010000000000000000000000000000000000000000000000    3
## 28  0000000110000000000000000000000001110000000000000000000    3
## 38  0000100000000000000000000000001001110000000000000000010    3
```

```
## 41   0000100110000000000000000000001001110000000000000000000010   3
## 51   1000001010000000000000000000001000000000000000000000000010   3
## 53   1000001110000000000000000000001001110000000000000000000010   3
## 93   1111111111111111111000000001111001110000000000001000011110   3
## 99   1111111111111111111000000001111001110000001000000000011110   3
## 110  1111111111111111111000000001111001110010000011111001110 3
## 1    0000000000000000000000000000000000000000000000000000000000   2
## 2    0000000000000000000000000000000000000000000000000000011000   2
## 20   0000000010000000000000000000001001110000000000000100010   2
## 23   0000000100000000000000000000000000000000000000000000000000   2
## 35   0000000110000000000000000000001001110001000000000011010   2
## 42   0110000000000000000000000000001110000000000000000000100   2
## 43   0110000000000000000000000000001110000000000000000011100   2
## 44   0110000000000000000000000000001110000000001000000100   2
## 58   1001011010000000000000000000001110000000000000011000   2
## 61   1111111111111111111000000000000000000000000000000000100   2
## 64   1111111111111111111000000000001001110000000000000000110   2
## 90   1111111111111111111000000001111001110000000000100000110   2
## 105  1111111111111111111000000001111001110001100000000011110   2
## 113  1111111111111111111000000001111001110010001011101010011110   2
## 114  1111111111111111111000000001111001110011000000000011110   2
## 115  1111111111111111111000000001111001110011000001000011110   2
## 10   0000000000000000000000000000001001110000000000000111010   1
## 14   0000000010000000000000000000000111000000000000000100000   1
## 15   0000000010000000000000000000000111000000000000000111000   1
## 16   0000000010000000000000000000000111001111111111110011000   1
## 21   0000000010000000000000000001001110000000000000111010   1
## 22   0000000010000000000000000001001110000000000100011010   1
## 24   0000000010000000000000000000000000000000000000000011000   1
## 27   0000000010000000000000000001000000000000000000000010   1
## 29   0000000110000000000000000000000111000000000000000011000   1
## 30   0000000110000000000000000000000111000000000000000100000   1
## 31   0000000110000000000000000001000000000000000000000010   1
## 34   0000000110000000000000000001001110000000000100011010   1
## 36   0000100000000000000000000000000111000000000000000011000   1
## 37   0000100000000000000000000000001000000000000000000000010   1
## 40   0000100100000000000000000000001001110000000000000000010   1
## 47   0110000000000000000000000000001001110000000001000000110   1
## 49   0110000110000000000000000001001110000000000000000110   1
## 54   1000101010000000000000000000001110000000000000000000   1
## 55   1001011000000000000000000000001110000000000000011000   1
## 59   1001011010000000000000000001001110000000000000000010   1
## 60   1001111010000000000000000001001110000000000000000010   1
## 65   1111111111111111111000000000011000000000000000000000100   1
## 67   1111111111111111111000000000011000111000000000000011100   1
## 70   1111111111111111111000000000111000000010000000010011110   1
## 73   1111111111111111111000000000111001110000000000000100110   1
## 78   1111111111111111111000000001110001110000000001000000100   1
## 79   1111111111111111111000000001110001110010000010000011100   1
## 80   1111111111111111111000000001110001110011000100010011100   1
## 83   1111111111111111111000000001111000000001000000010011110   1
## 84   1111111111111111111000000001111000000110001000010011110   1
## 91   1111111111111111111000000001111001110000000000100011110   1
## 94   1111111111111111111000000001111001110000000001100011110   1
```

```
## 95  111111111111111111110000000011110011100000000001110011110      1
## 97  111111111111111111100000000011110011100000000100010011110      1
## 98  111111111111111111100000000011110011100000001001010011110      1
## 100 111111111111111111100000000011110011100000100000000011110      1
## 101 111111111111111111100000000011110011100001001110100011110      1
## 102 111111111111111111100000000011110011100010000000000011110      1
## 103 111111111111111111100000000011110011100010000010000011110      1
## 104 111111111111111111100000000011110011100010010000000011110      1
## 107 111111111111111111100000000011110011100100000001000011110      1
## 108 111111111111111111100000000011110011100100000011100011110      1
## 109 111111111111111111100000000011110011100100000011110011110      1
## 111 111111111111111111100000000011110011100100001001100111110      1
## 112 111111111111111111100000000011110011100100010011010011110      1
## 116 111111111111111111100000000011110011100110000100000011110      1
## 117 111111111111111111100000000011110011100110001101110011110      1
## 118 111111111111111111100000000011110011100111111101111011110      1
```