

main_exons

March 8, 2022

1 eQTL boxplot

This is script ported from python to fix unknown plotting error.

```
[1]: suppressPackageStartupMessages({  
      library(tidyverse)  
      library(ggpubr)  
    })
```

1.1 Functions

```
[2]: feature = "exons"
```

1.1.1 Basic loading functions

```
[3]: get_residualized_df <- function(){  
      expr_file = paste0("/ceph/projects/v4_phase3_paper/analysis/eql_analysis/  
      ↪all/",  
                          feature, "/expression_gct/covariates/  
      ↪residualized_expression/_m/",  
                          feature, "_residualized_expression.csv")  
      return(data.table::fread(expr_file) %>% column_to_rownames("gene_id"))  
    }  
    memRES <- memoise::memoise(get_residualized_df)  
  
    get_pheno_df <- function(){  
      phenotype_file = paste0('/ceph/projects/v4_phase3_paper/inputs/',  
                              'phenotypes/_m/merged_phenotypes.csv')  
      return(data.table::fread(phenotype_file))  
    }  
    memPHENO <- memoise::memoise(get_pheno_df)  
  
    get_genotypes <- function(){  
      traw_file = paste0("/ceph/projects/brainseq/genotype/download/topmed/  
      ↪convert2plink/",  
                          "filter_maf_01/a_transpose/_m/LIBD_Brain_TopMed.traw")  
      traw = data.table::fread(traw_file) %>% rename_with(~ gsub('\\_.*', '', .x))  
      return(traw)
```

```

}
memSNPs <- memoise::memoise(get_genotypes)

```

1.1.2 eQTL and helpful functions

```

[4]: feature_map <- function(feature){
  return(list("genes"="Gene", "transcripts"= "Transcript",
    "exons"= "Exon", "junctions"= "Junction")[[feature]])
}

save_ggplots <- function(fn, p, w, h){
  for(ext in c('.pdf', '.png', '.svg')){
    ggsave(paste0(fn, ext), plot=p, width=w, height=h)
  }
}

get_caudeate_eqtls <- function(){
  mashr_file = "../summary_table/_m/BrainSeq_caudeate_eQTL.txt.gz"
  return(data.table::fread(mashr_file) %>%
    filter(Type == feature_map(feature)) %>%
    select(gene_id, variant_id, AA, EA))
}

memCAUDEATE <- memoise::memoise(get_caudeate_eqtls)

get_eqtl_df <- function(){
  eGenes_file = paste0("../summary_table/_m/", feature, "/lfsr_allpairs_ancestry.txt.gz")
  eGenes = data.table::fread(eGenes_file)
  return(eGenes)
}

memEQTL <- memoise::memoise(get_eqtl_df)

```

1.1.3 Basic eQTL plotting functions

```

[5]: get_genotype_annot <- function(){
  return(memSNPs() %>% select(CHR, SNP, POS, COUNTED, ALT))
}

get_snps_df <- function(){
  return(memSNPs() %>% select("SNP", starts_with("Br")))
}

letter_snp <- function(number, a0, a1){
  if(is.na(number)){ return(NA) }
  if( length(a0) == 1 & length(a1) == 1){
    seps = ""; collapse=""
  } else {
    seps = " "; collapse=NULL
  }
}

```

```

    }
    return(paste(paste0(rep(a0, number), collapse = collapse),
                 paste0(rep(a1, (2-number)), collapse = collapse), sep=seps))
  }

get_snp_df <- function(variant_id, gene_id){
  zz = get_genotype_annot() %>% filter(SNP == variant_id)
  xx = get_snps_df() %>% filter(SNP == variant_id) %>%
    column_to_rownames("SNP") %>% t %>% as.data.frame %>%
    rownames_to_column("BrNum") %>% mutate(COUNTED=zz$COUNTED, ALT=zz$ALT)
  #>%
  rename("SNP"=all_of(variant_id))
  yy = memRES()[gene_id, ] %>% t %>% as.data.frame %>%
    rownames_to_column("BrNum") %>% inner_join(memPHENO(), by="BrNum")
  ## Annotated SNPs
  letters = c()
  for(ii in seq_along(xx$COUNTED)){
    a0 = xx$COUNTED[ii]; a1 = xx$ALT[ii]; number = xx$SNP[ii]
    letters <- append(letters, letter_snp(number, a0, a1))
  }
  xx = xx %>% mutate(LETTER=letters, ID=paste(SNP, LETTER, sep="\n"))
  df = inner_join(xx, yy, by="BrNum") %>% mutate_if(is.character, as.factor)
  return(df)
}
memDF <- memoise::memoise(get_snp_df)

get_gene_symbol <- function(gene_id){
  ensemblID = gsub("\\..*", "", gene_id)
  geneid = memMART() %>% filter(ensembl_gene_id == gsub("\\..*", "", gene_id))
  if(dim(geneid)[1] == 0){
    return("")
  } else {
    return(geneid$external_gene_name)
  }
}

```

```

[6]: plot_simple_eqtl <- function(fn, gene_id, variant_id, eqtl_annot, prefix,
  y0=NULL, y1=NULL){
  if(is.null(y0)){ y0 = quantile(memDF(variant_id, gene_id)[[gene_id]],
  probs=c(0.01))[[1]] - 0.2 }
  if(is.null(y1)){ y1 = quantile(memDF(variant_id, gene_id)[[gene_id]],
  probs=c(0.99))[[1]] + 0.2 }
  bxp = memDF(variant_id, gene_id) %>%
    ggboxplot(x="ID", y=gene_id, fill="Race", color="Race", add="jitter",
              xlab=variant_id, ylab="Residualized Expression", outlier.
  shape=NA,
              add.params=list(alpha=0.5), alpha=0.4, legend="bottom",

```

```

        palette="npg", ylim=c(y0,y1),
    ↪ggtheme=theme_pubr(base_size=20, border=TRUE)) +
        font("xy.title", face="bold") +
        ggtitle(paste(prefix, eqtl_annot, sep='\n')) +
        theme(plot.title = element_text(hjust = 0.5, face="bold"))
    print(bxp)
    save_ggplots(fn, bxp, 7, 7)
}

```

1.1.4 GWAS plots

```

[7]: get_gwas_snps <- function(){
    gwas_snp_file = paste0('/ceph/projects/v4_phase3_paper/inputs/sz_gwas/pgc3/
    ↪',
                                'map_phase3/_m/libd_hg38_pgc2sz_snps_p5e_minus8.tsv')
    gwas_df = data.table::fread(gwas_snp_file) %>% arrange(P)
    return(gwas_df)
}
memGWAS <- memoise::memoise(get_gwas_snps)

get_gwas_snp <- function(variant){
    return(memGWAS() %>% filter(our_snp_id == variant))
}

get_risk_allele <- function(variant){
    gwas_snp = get_gwas_snp(variant)
    if(gwas_snp$OR > 1){
        ra = gwas_snp$A1
    }else{
        ra = gwas_snp$A2
    }
    return(ra)
}

get_eqtl_gwas_df <- function(){
    return(memCAUDATE() %>% inner_join(memGWAS(),
    ↪by=c("variant_id"="our_snp_id")))
}

get_gwas_ordered_snp_df <- function(variant_id, gene_id,
    ↪pgc3_a1_same_as_our_counted, OR){
    df = memDF(variant_id, gene_id)
    if(!pgc3_a1_same_as_our_counted){ # Fix bug with matching alleles!
        if(OR < 1){ df = df %>% mutate(SNP = 2-SNP, ID=paste(SNP, LETTER,
    ↪sep="\n")) }
    } else {

```

```

        if(OR > 1){ df = df %>% mutate(SNP = 2-SNP, ID=paste(SNP, LETTER,
↪sep="\n")) }
    }
    return(df)
}

plot_gwas_eqtl <- function(fn, gene_id, variant_id, eqtl_annot,
↪pgc3_a1_same_as_our_counted,
                                OR, title){
    dt = get_gwas_ordered_snp_df(variant_id, gene_id,
↪pgc3_a1_same_as_our_counted, OR)
    bxp = dt %>% mutate_if(is.character, as.factor) %>%
        ggboxplot(x="ID", y=gene_id, fill="Race", color="Race", add="jitter",
                    xlab=variant_id, ylab="Residualized Expression", outlier.
↪shape=NA,
                    add.params=list(alpha=0.5), alpha=0.4, legend="bottom",
↪#ylim=c(y0,y1),
                    palette="npg", ggtheme=theme_pubr(base_size=20, border=TRUE))
    ↪+
        font("xy.title", face="bold") + ggtitle(title) +
        theme(plot.title = element_text(hjust = 0.5, face="bold"))
    print(bxp)
    save_ggplots(fn, bxp, 7, 8)
}

```

1.2 Plot eQTL

```

[8]: get_drd2_exon_annotation <- function(brainseq_exon_id){
    return(list(
        'e667152' = 'DRD2 exon 1S',
        'e667153' = 'DRD2 exon 2',
        'e667154' = 'DRD2 exon 3',
        'e667155' = 'DRD2 exon 4',
        'e667156' = 'DRD2 exon 5',
        'e667157' = 'DRD2 exon 7',
        'e667158' = 'DRD2 exon 8',
        'e667159' = 'DRD2 exon 1L',
        'e667164' = 'DRD2 exon 6'
    )[[brainseq_exon_id]])
}

```

1.2.1 DRD2 plot

```

[9]: drd2_short = "ENST00000346454.7"; drd2_long = "ENST00000362072.7"
exons_df = data.table::fread("../.../DE_gwas_eqtl_integration/_h/exons.
↪csv") %>%

```

```

    filter(transcript_id %in% c(drd2_short, drd2_long)) %>%
    select(gene_id, gene_name, exon_id, brainseq_exon_id) %>%
    distinct(brainseq_exon_id, .keep_all=TRUE)
all_drd2_exons = exons_df$brainseq_exon_id
all_drd2_exons

```

1. 'e667152' 2. 'e667153' 3. 'e667154' 4. 'e667155' 5. 'e667156' 6. 'e667157' 7. 'e667158' 8. 'e667159'
9. 'e667164'

```

[10]: drd2_df0 = memCAUDATE() %>% filter(gene_id %in% all_drd2_exons) %>%
      arrange(AA, EA) %>% group_by(gene_id) %>% slice(1) %>% arrange(AA, EA)
drd2_df0

```

| | gene_id | variant_id | AA | EA |
|---------------------|---------|---------------------|------------|------------|
| | <chr> | <chr> | <dbl> | <dbl> |
| A grouped_df: 2 × 4 | e667159 | chr11:113528769:T:C | 0.06756843 | 0.01091546 |
| | e667157 | chr11:113371155:C:T | 0.11607231 | 0.01183844 |

```

[11]: drd2_df = memEQTL() %>% filter(gene_id %in% all_drd2_exons) %>%
      arrange(EA, AA) %>% group_by(gene_id) %>% slice(1) %>% arrange(AA, EA)
drd2_df

```

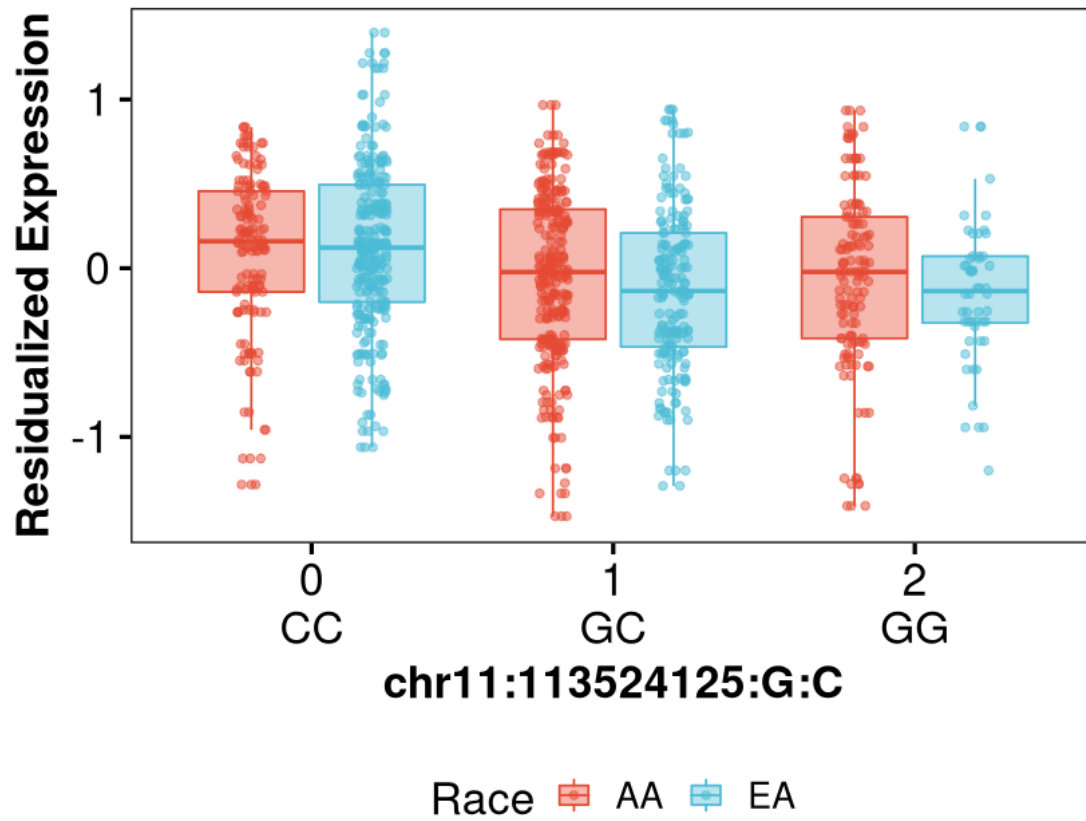
| | effect | gene_id | variant_id | AA | EA |
|---------------------|-----------------------------|---------|---------------------|------------|-----------|
| | <chr> | <chr> | <chr> | <dbl> | <dbl> |
| A grouped_df: 4 × 5 | e667159_chr11:113524125:G:C | e667159 | chr11:113524125:G:C | 0.07281645 | 0.0045421 |
| | e667158_chr11:112955463:T:G | e667158 | chr11:112955463:T:G | 0.07634220 | 0.0859227 |
| | e667157_chr11:113358943:T:C | e667157 | chr11:113358943:T:C | 0.13383811 | 0.0039085 |
| | e667153_chr11:113446345:G:A | e667153 | chr11:113446345:G:A | 0.47320254 | 0.1383174 |

```

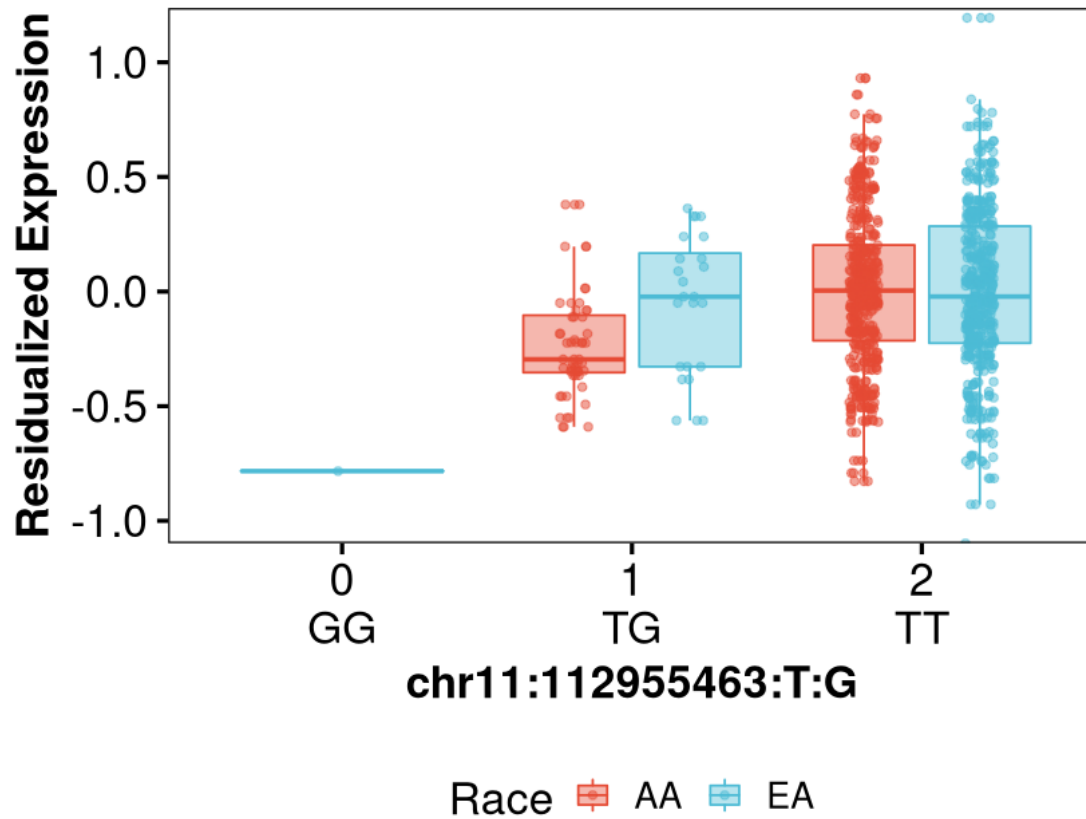
[12]: for(x in seq_along(drd2_df$gene_id)){
  anno = get_drd2_exon_annotation(drd2_df$gene_id[x])
  en = gsub(" ", "_", anno)
  ids = exons_df %>% filter(brainseq_exon_id == drd2_df$gene_id[x]) %>%
  ↪select("exon_id")
  fn = paste("drd2_eqtl", en, sep="_")
  eqtl_annot = paste(paste("eQTL (AA) lfsr:", signif(drd2_df$AA[x], 2)),
                    paste("eQTL (EA) lfsr:", signif(drd2_df$EA[x], 2)),
  ↪sep='\n')
  prefix = paste(anno, ids$exon_id, sep='\n')
  plot_simple_eqtl(fn, drd2_df$gene_id[x], drd2_df$variant_id[x], eqtl_annot,
  ↪prefix)
  #print(eqtl_annot)
}

```

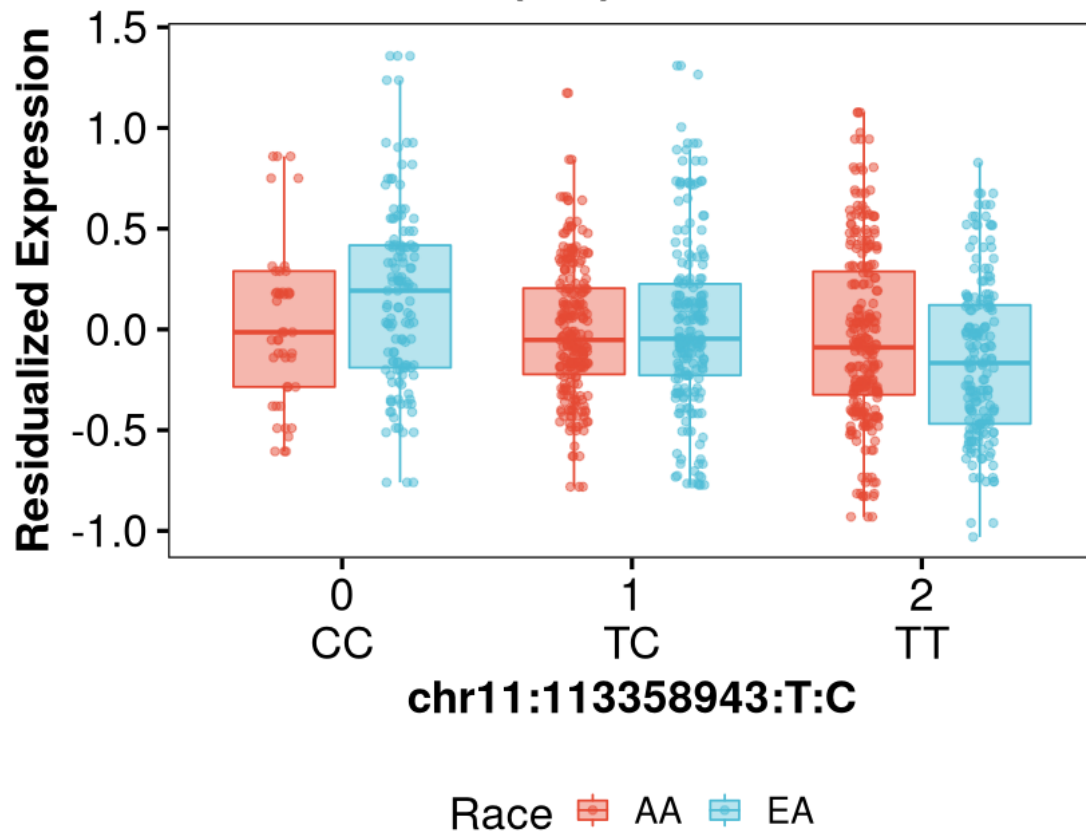
DRD2 exon 1L
ENSE00002282098.1
eQTL (AA) Ifsr: 0.073
eQTL (EA) Ifsr: 0.0045



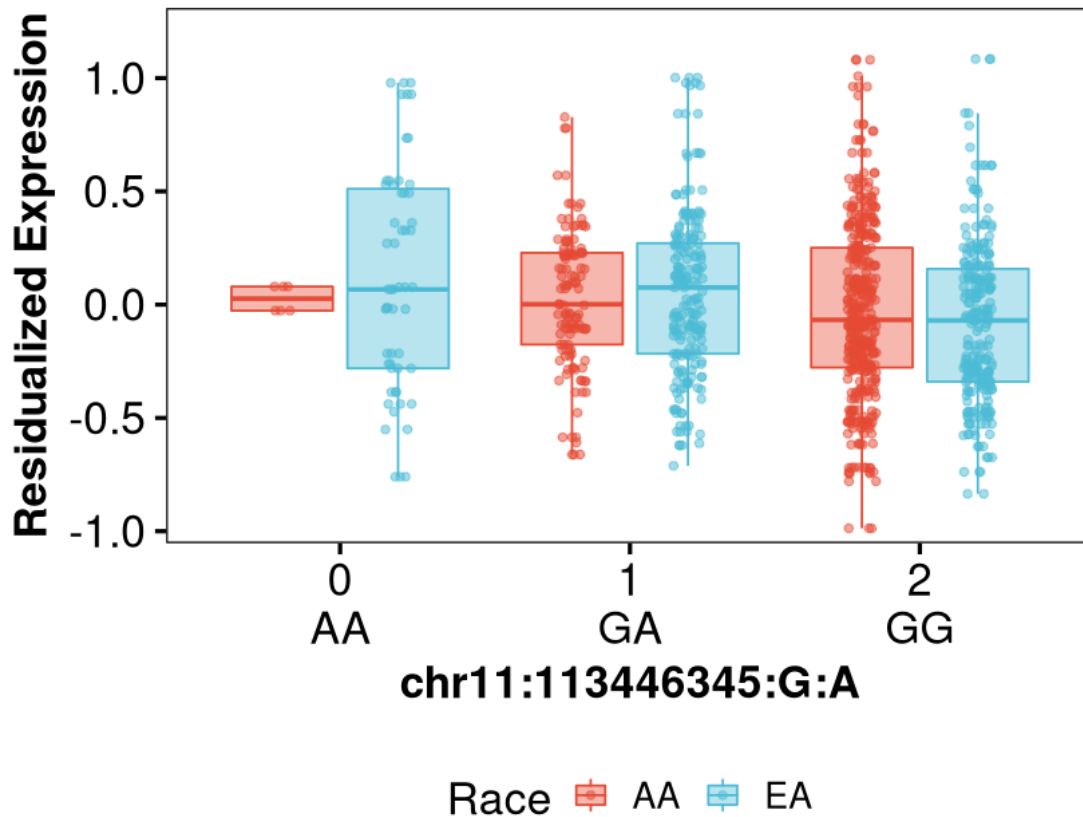
DRD2 exon 8
ENSE00002210277.1
eQTL (AA) lfsr: 0.076
eQTL (EA) lfsr: 0.086



DRD2 exon 7
ENSE00001247012.2
eQTL (AA) lfsr: 0.13
eQTL (EA) lfsr: 0.0039



DRD2 exon 2
ENSE00003665261.1
eQTL (AA) lfsr: 0.47
eQTL (EA) lfsr: 0.14



1.2.2 GWAS association

```
[13]: eGenes_gwas = get_eqtl_gwas_df() %>% filter(gene_id %in% all_drd2_exons) %>%
      arrange(EA, AA, P) %>% group_by(gene_id) %>% slice(1) %>% arrange(EA, AA, P)
eGenes_gwas
```

Warning message in cbind(parts\$left, ellip_h, parts\$right, deparse.level = 0L):
 "number of rows of result is not a multiple of vector length (arg 2)"
 Warning message in cbind(parts\$left, ellip_h, parts\$right, deparse.level = 0L):
 "number of rows of result is not a multiple of vector length (arg 2)"
 Warning message in cbind(parts\$left, ellip_h, parts\$right, deparse.level = 0L):
 "number of rows of result is not a multiple of vector length (arg 2)"
 Warning message in cbind(parts\$left, ellip_h, parts\$right, deparse.level = 0L):

"number of rows of result is not a multiple of vector length (arg 2)"

A grouped_df: 0 × 28

| gene_id | variant_id | AA | EA | V1 | chrN | cm | pos | our_counted |
|---------|------------|-------|-------|-------|-------|-------|-------|-------------|
| <chr> | <chr> | <dbl> | <dbl> | <int> | <int> | <int> | <int> | <chr> |

1.3 Session Info

```
[14]: Sys.time()
      proc.time()
      options(width = 120)
      sessioninfo::session_info()
```

```
[1] "2022-03-08 19:36:23 EST"
```

| user | system | elapsed |
|----------|---------|----------|
| 4868.976 | 268.288 | 1938.636 |

\$platform **\$version** 'R version 4.1.2 (2021-11-01)'

\$os 'Arch Linux'

\$system 'x86_64, linux-gnu'

\$ui 'X11'

\$language '(EN)'

\$collate 'en_US.UTF-8'

\$ctype 'en_US.UTF-8'

\$tz 'America/New_York'

\$date '2022-03-08'

\$pandoc '2.14.1 @ /usr/bin/pandoc'

| | package <chr> | ondiskversion <chr> | loadedversion <chr> | path <chr> |
|-------------------------------------|------------------|------------------------|------------------------|---------------|
| | abind | 1.4.5 | 1.4-5 | /home/jbe |
| | assertthat | 0.2.1 | 0.2.1 | /home/jbe |
| | backports | 1.4.1 | 1.4.1 | /home/jbe |
| | base64enc | 0.1.3 | 0.1-3 | /home/jbe |
| | broom | 0.7.12 | 0.7.12 | /home/jbe |
| | cachem | 1.0.6 | 1.0.6 | /home/jbe |
| | car | 3.0.12 | 3.0-12 | /home/jbe |
| | carData | 3.0.5 | 3.0-5 | /home/jbe |
| | cellranger | 1.1.0 | 1.1.0 | /home/jbe |
| | cli | 3.2.0 | 3.2.0 | /home/jbe |
| | colorspace | 2.0.3 | 2.0-3 | /home/jbe |
| | crayon | 1.5.0 | 1.5.0 | /home/jbe |
| | data.table | 1.14.2 | 1.14.2 | /home/jbe |
| | DBI | 1.1.2 | 1.1.2 | /home/jbe |
| | dbplyr | 2.1.1 | 2.1.1 | /home/jbe |
| | digest | 0.6.29 | 0.6.29 | /home/jbe |
| | dplyr | 1.0.8 | 1.0.8 | /home/jbe |
| | ellipsis | 0.3.2 | 0.3.2 | /home/jbe |
| | evaluate | 0.15 | 0.15 | /home/jbe |
| | fansi | 1.0.2 | 1.0.2 | /home/jbe |
| | farver | 2.1.0 | 2.1.0 | /home/jbe |
| | fastmap | 1.1.0 | 1.1.0 | /home/jbe |
| | forcats | 0.5.1 | 0.5.1 | /home/jbe |
| | fs | 1.5.2 | 1.5.2 | /home/jbe |
| | generics | 0.1.2 | 0.1.2 | /home/jbe |
| | ggplot2 | 3.3.5 | 3.3.5 | /home/jbe |
| | ggpubr | 0.4.0 | 0.4.0 | /home/jbe |
| | ggsci | 2.9 | 2.9 | /home/jbe |
| | ggsignif | 0.6.3 | 0.6.3 | /home/jbe |
| \$packages A packages_info: 78 x 11 | glue | 1.6.1 | 1.6.1 | /home/jbe |
| | purrr | 0.3.4 | 0.3.4 | /home/jbe |
| | R.methodsS3 | 1.8.1 | 1.8.1 | /home/jbe |
| | R.oo | 1.24.0 | 1.24.0 | /home/jbe |
| | R.utils | 2.11.0 | 2.11.0 | /home/jbe |
| | R6 | 2.5.1 | 2.5.1 | /home/jbe |
| | Rcpp | 1.0.8 | 1.0.8 | /home/jbe |
| | readr | 2.1.2 | 2.1.2 | /home/jbe |
| | readxl | 1.3.1 | 1.3.1 | /home/jbe |
| | repr | 1.1.4 | 1.1.4 | /home/jbe |
| | reprex | 2.0.1 | 2.0.1 | /home/jbe |
| | rlang | 1.0.1 | 1.0.1 | /home/jbe |
| | rstatix | 0.7.0 | 0.7.0 | /home/jbe |
| | rstudioapi | 0.13 | 0.13 | /home/jbe |
| | rvest | 1.0.2 | 1.0.2 | /home/jbe |
| | scales | 1.1.1 | 1.1.1 | /home/jbe |
| | sessioninfo | 1.2.2 | 1.2.2 | /home/jbe |
| | stringi | 1.7.6 | 1.7.6 | /home/jbe |
| | stringr | 1.4.0 | 1.4.0 | /home/jbe |
| | svglite | 2.1.0 | 2.1.0 | /home/jbe |
| | systemfonts | 1.0.4 | 1.0.4 | /home/jbe |