

# [Tutorial] Human Genome Annotation

## 1. Introduction

### 1.1. What is gene annotation?

Over the past years, we have learnt that there are a number of chromosomes and genes in our genome. Counting the number of chromosomes is fairly easy but students might find difficult to say how many genes we have in our genome. If you can get an answer for this, could you tell how many genes encode protein and how many do not?

To answer this question, we need to access the database for gene annotation. Gene annotation is the process of making nucleotide sequence meaningful - where genes are located? whether it is protein-coding or noncoding. If you would like to get an overview of gene annotation, please find this link.

One of well-known collaborative efforts in gene annotation is the GENCODE consortium. It is a part of the Encyclopedia of DNA Elements (The ENCODE project consortium) and aims to identify all gene features in the human genome using a combination of computational analysis, manual annotation, and experimental validation (Harrow et al. 2012). You might find another database for gene annotation, like RefSeq, CCDS, and need to understand differences between them.

In this tutorial, we will access to gene annotation from the GENCODE consortium and explore genes and functional elements in our genome.

### 1.2. Aims

What we will do with this dataset:

- Be familiar with gene annotation modality.
- Tidy data and create a table for your analysis.
- Apply tidyverse functions for data munging.

Please note that there is better solution for getting gene annotation in R if you use a biomart. Our tutorial is only designed to have a practice on tidyverse exercise.

## 2. Explore your data

### 2.1. Unboxing your dataset

This tutorial will use a gene annotation file from the GENCODE. You will need to download the file from the GENCODE. If you are using terminal, please download file using `wget`:

```
# Run from your terminal, not R console
# wget ftp://ftp.ebi.ac.uk/pub/databases/gencode/Gencode_human/release_31/gencode.v31.basic.annotation.

# Once you downloaded the file, you won't need to download it again. So please comment out the command
```

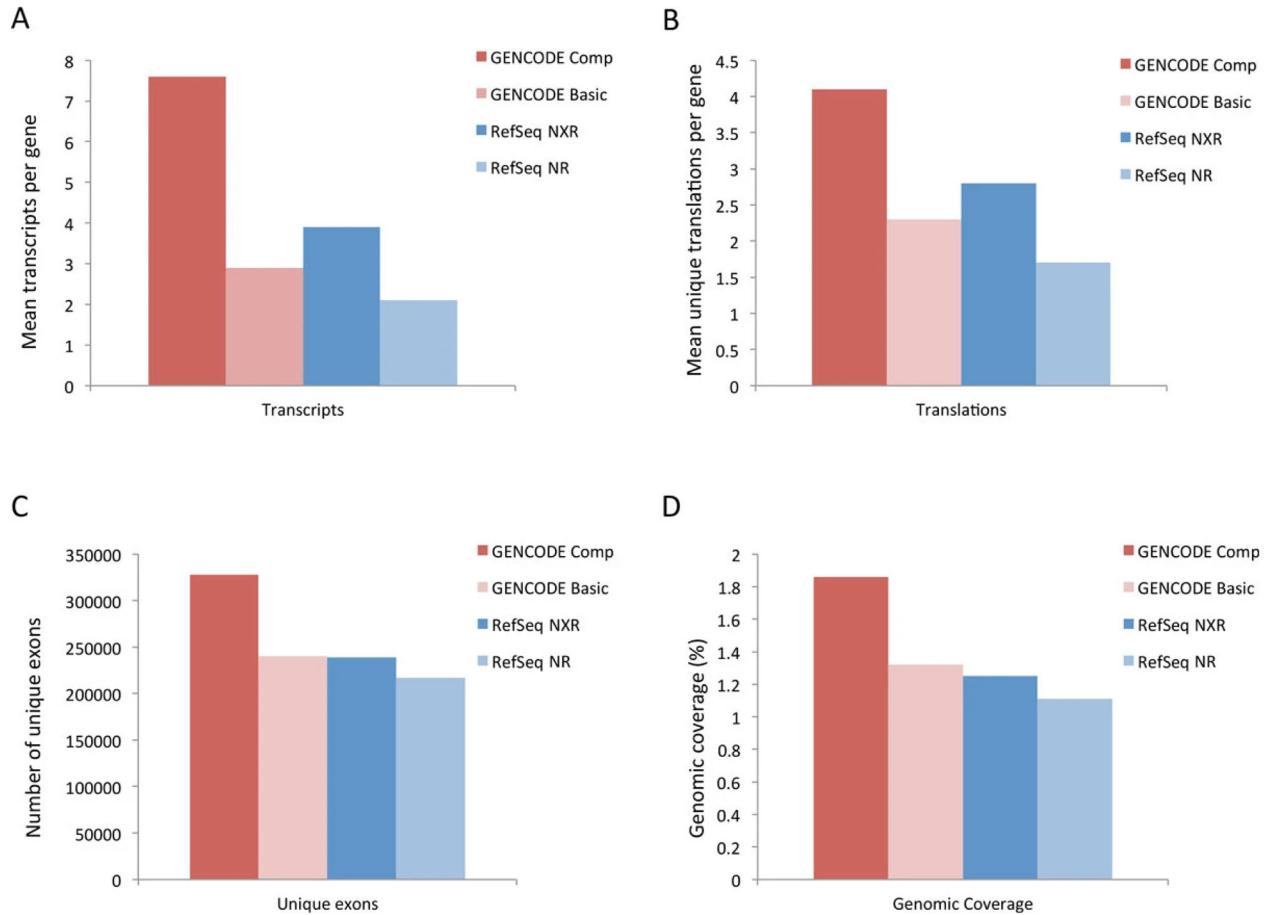


Figure 1: Comparison of GENCODE and RefSeq gene annotation and the impact of reference geneset on variant effect prediction (Frankish et al. 2015). A) Mean number of alternatively spliced transcripts per multi-exon protein-coding locus B) Mean number of unique CDS per multi-exon protein-coding locus C) Mean number of unique (non-redundant) exons per multi-exon protein-coding locus D) Percentage genomic coverage of unique (non-redundant) exons at multi-exon protein-coding loci.

Once you download the file, you can print out the first few lines using the following bash command (we will learn UNIX commands later):

```
# Run from your terminal, not R console
# gzcat gencode.v31.basic.annotation.gtf.gz | head -7
```

The file is the GTF file format, which you will find most commonly in gene annotation. Please read the file format thoroughly in the link above. For the tutorial, we need to load two packages. If the package is not installed in your system, please install it.

- **tidyverse**, a package you have learnt from the chapter 5.
- **readr**, a package provides a fast and friendly way to read. Since the file `gencode.v31.basic.annotation.gtf.gz` is pretty large, you will need some function to load data quickly into your workspace. `readr` is a part of `tidyverse`, so you can just load `tidyverse` to use `readr` functions.

Let's load the GTF file into your workspace. We will use `read_delim` function from the `readr` package. This is much faster loading than `read.delim` or `read.csv` from R base. However, please keep in mind that some parameters and output class for `read_delim` are slightly different from them.

```
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --

## v ggplot2 3.3.5     v purrr    0.3.4
## v tibble   3.1.4     v dplyr    1.0.7
## v tidyr    1.1.3     v stringr  1.4.0
## v readr    2.0.1     vforcats  0.5.1

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()   masks stats::lag()

read_delim('gencode.v31.basic.annotation.gtf.gz',
           delim = '\t', skip = 5, progress = F,
           col_names = F)

## Rows: 1756502 Columns: 9

## -- Column specification -----
## Delimiter: "\t"
## chr (7): X1, X2, X3, X6, X7, X8, X9
## dbl (2): X4, X5

##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.

## # A tibble: 1,756,502 x 9
##       X1      X2      X3          X4      X5  X6      X7      X8      X9
##       <chr> <chr> <chr>      <dbl> <dbl> <chr> <chr> <chr> <chr>
```

```

## 1 chr1 HAVANA gene      11869 14409 . + . "gene_id \"ENSG0000022~
## 2 chr1 HAVANA transcript 11869 14409 . + . "gene_id \"ENSG0000022~
## 3 chr1 HAVANA exon      11869 12227 . + . "gene_id \"ENSG0000022~
## 4 chr1 HAVANA exon      12613 12721 . + . "gene_id \"ENSG0000022~
## 5 chr1 HAVANA exon      13221 14409 . + . "gene_id \"ENSG0000022~
## 6 chr1 HAVANA transcript 12010 13670 . + . "gene_id \"ENSG0000022~
## 7 chr1 HAVANA exon      12010 12057 . + . "gene_id \"ENSG0000022~
## 8 chr1 HAVANA exon      12179 12227 . + . "gene_id \"ENSG0000022~
## 9 chr1 HAVANA exon      12613 12697 . + . "gene_id \"ENSG0000022~
## 10 chr1 HAVANA exon     12975 13052 . + . "gene_id \"ENSG0000022~
## # ... with 1,756,492 more rows

```

Can you find out what the parameters mean? Few things to note are:

- The GTF file contains the first few lines for comments (#). In general, the file contains description, provider, date, format.
- The GTF file does not have column names so you will need to assign ‘FALSE’ for col\_names.

This is sort of canonical way to load your dataset into R. However, we are using a GTF format, which is specific to gene annotation so we can use a package to specifically handle a GTF file.

Here I introduce the package rtracklayer. Let’s install the package first.

```

# if (!requireNamespace("BiocManager", quietly = TRUE))
#   install.packages("BiocManager")

#BiocManager::install("rtracklayer")

```

Then, now you can read the GTF file using this package. Then, you can check the class of the object d.

```

d = rtracklayer::import('gencode.v31.basic.annotation.gtf.gz')
class(d)

## [1] "GRanges"
## attr(,"package")
## [1] "GenomicRanges"

```

You will find out that this is GRanges class. This is from the package Genomic Range, specifically dealing with genomic datasets but we are not heading into this in this tutorial. So please find this information if you are serious on this.

We are converting d into a data frame as following:

```
d = d %>% as.data.frame()
```

Let’s overview few lines from the data frame, and explore what you get in this object.

```
head(d)
```

```

##   seqnames start   end width strand source      type score phase
## 1     chr1 11869 14409  2541      + HAVANA    gene     NA    NA
## 2     chr1 11869 14409  2541      + HAVANA transcript  NA    NA

```

```

## 3      chr1 11869 12227    359      + HAVANA      exon     NA     NA
## 4      chr1 12613 12721    109      + HAVANA      exon     NA     NA
## 5      chr1 13221 14409   1189      + HAVANA      exon     NA     NA
## 6      chr1 12010 13670   1661      + HAVANA transcript  NA     NA
##           gene_id                  gene_type gene_name level
## 1 ENSG00000223972.5 transcribed_unprocessed_pseudogene DDX11L1    2
## 2 ENSG00000223972.5 transcribed_unprocessed_pseudogene DDX11L1    2
## 3 ENSG00000223972.5 transcribed_unprocessed_pseudogene DDX11L1    2
## 4 ENSG00000223972.5 transcribed_unprocessed_pseudogene DDX11L1    2
## 5 ENSG00000223972.5 transcribed_unprocessed_pseudogene DDX11L1    2
## 6 ENSG00000223972.5 transcribed_unprocessed_pseudogene DDX11L1    2
##           hgnc_id      havana_gene transcript_id
## 1 HGNC:37102 OTTHUMG00000000961.2 <NA>
## 2 HGNC:37102 OTTHUMG00000000961.2 ENST00000456328.2
## 3 HGNC:37102 OTTHUMG00000000961.2 ENST00000456328.2
## 4 HGNC:37102 OTTHUMG00000000961.2 ENST00000456328.2
## 5 HGNC:37102 OTTHUMG00000000961.2 ENST00000456328.2
## 6 HGNC:37102 OTTHUMG00000000961.2 ENST00000450305.2
##           transcript_type transcript_name transcript_support_level
## 1                      <NA>          <NA>                 <NA>
## 2                      lncRNA        DDX11L1-202            1
## 3                      lncRNA        DDX11L1-202            1
## 4                      lncRNA        DDX11L1-202            1
## 5                      lncRNA        DDX11L1-202            1
## 6 transcribed_unprocessed_pseudogene        DDX11L1-201            NA
##           tag      havana_transcript exon_number      exon_id      ont
## 1 <NA>          <NA>          <NA>          <NA>          <NA>
## 2 basic          OTTHUMT00000362751.1 <NA>          <NA>          <NA>
## 3 basic          OTTHUMT00000362751.1    1 ENSE00002234944.1 <NA>
## 4 basic          OTTHUMT00000362751.1    2 ENSE00003582793.1 <NA>
## 5 basic          OTTHUMT00000362751.1    3 ENSE00002312635.1 <NA>
## 6 basic          OTTHUMT0000002844.2 <NA>          <NA> PGO:0000019
##           protein_id ccdsid
## 1          <NA> <NA>
## 2          <NA> <NA>
## 3          <NA> <NA>
## 4          <NA> <NA>
## 5          <NA> <NA>
## 6          <NA> <NA>

```

One thing you can find is that there is no columns in the data frame. Let's match which information is provided in columns. You can find the instruction page in the website ([link](#)).

Based on this, you can assign a name for 9 columns. One thing to remember is you should not use space for the column name. Spacing in the column name is actually working but not a good habit for your code. So please replace a space with **underscore** in the column name.

```

# Assign column names according to the GENCODE instruction.
cols = c('chrom', 'source', 'feature_type', 'start', 'end', 'score', 'strand', 'phase', 'info')

```

Now you can set up the column names into the `col_names` parameter, and load the file into a data frame.

```

f = read_delim('gencode.v31.basic.annotation.gtf.gz',
               delim='\t', skip = 5,
               progress = F,
               col_names = cols)

## Rows: 1756502 Columns: 9

## -- Column specification -----
## Delimiter: "\t"
## chr (7): chrom, source, feature_type, score, strand, phase, info
## dbl (2): start, end

##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.

```

You can find the column names are now all set.

```
head(f)
```

```

## # A tibble: 6 x 9
##   chrom source feature_type start   end score strand phase info
##   <chr>  <chr>    <chr>     <dbl> <dbl> <chr>  <chr>  <chr>
## 1 chr1   HAVANA gene        11869 14409 .      +      .      "gene_id \\"ENSG00000~
## 2 chr1   HAVANA transcript  11869 14409 .      +      .      "gene_id \\"ENSG00000~
## 3 chr1   HAVANA exon       11869 12227 .      +      .      "gene_id \\"ENSG00000~
## 4 chr1   HAVANA exon       12613 12721 .      +      .      "gene_id \\"ENSG00000~
## 5 chr1   HAVANA exon       13221 14409 .      +      .      "gene_id \\"ENSG00000~
## 6 chr1   HAVANA transcript 12010 13670 .      +      .      "gene_id \\"ENSG00000~

```

When you loaded the file, you see the message about the data class. You might want to overview this data.

```
summary(f)
```

```

##      chrom          source      feature_type      start
##  Length:1756502  Length:1756502  Length:1756502  Min.   : 577
##  Class :character Class :character Class :character  1st Qu.: 32101517
##  Mode  :character Mode  :character Mode  :character Median : 61732754
##                                         Mean   : 75288563
##                                         3rd Qu.:111760181
##                                         Max.   :248936581
##      end            score      strand      phase
##  Min.   :     647  Length:1756502  Length:1756502  Length:1756502
##  1st Qu.: 32107331 Class :character Class :character Class :character
##  Median : 61738373 Mode  :character Mode  :character Mode  :character
##  Mean   : 75292632
##  3rd Qu.:111763007
##  Max.   :248937043
##      info
##  Length:1756502
##  Class :character

```

```
##  Mode :character
##
##
```

## 2.2. How many feature types in the GENCODE dataset?

As instructed in the GENCODE website, the GENCODE dataset provides a range of annotations for the feature type. You can check feature types using \_\_\_\_\_ function.

```
f %>% group_by(feature_type) %>% count(feature_type)
```

```
## # A tibble: 8 x 2
## # Groups:   feature_type [8]
##   feature_type     n
##   <chr>        <int>
## 1 CDS          567862
## 2 exon         744835
## 3 gene          60603
## 4 Selenocysteine    96
## 5 start_codon      57886
## 6 stop_codon       57775
## 7 transcript        108243
## 8 UTR            159202
```

```
# table(d$feature_type)
```

How many feature types provided in the GENCODE? And how many items stored for each feature type? Please write down the number of feature types from the dataset. Also, if you are not familiar with these types, it would be good to put one or two sentences that can describe each type).

```
#8 feature types
#CDS is coding sequence whose sequence can be translated into proteins
#exon is the part of a gene, also the counterpart of intron, that is not spliced out and remains in mature mRNA
#gene is a unit that contains the information that could make a functional protein.
#Selenocysteine is the 21st proteinogenic amino acid.
#start_codon is the sequence in an RNA or a template DNA that signals the start of protein synthesis
#stop_codon is a sequence within mRNA that signals a halt to protein synthesis
#transcript is the transcribed mRNA from its corresponding DNA.
#UTR is the untranslated region on each side of a coding sequence on a strand of mRNA.
```

## 2.3. How many genes we have?

Let's count the number of genes in our genome. Since we know that the column `feature_type` contains rows with `gene`, which contains obviously annotations for genes. We might want to subset those rows from the data frame.

```
f1 = filter(f, feature_type == 'gene')
#d1 = d[d$feature_type == 'gene', ]
```

## 2.4. Ensembl, Havana and CCDS.

Gene annotation for the human genome is provided by multiple organizations with different gene annotation methods and strategy. This means that information can be varying by resources, and users need to understand heterogeneity inherent in annotation databases. The GENCODE project utilizes two sources of gene annotation.

1. Havana: Manual gene annotation (detailed strategy in here)
2. Ensembl: Automatic gene annotation (detailed strategy in here)

It provides the combination of Ensembl/HAVANA gene set as the default gene annotation for the human genome. In addition, they also guarantee that all transcripts from the Consensus Coding Sequence (CCDS) set are present in the GENCODE gene set. The CCDS project is a collaborative effort to identify a core set of protein coding regions that are consistently annotated and of high quality. Initial results from the Consensus CDS (CCDS) project are now available through the appropriate Ensembl gene pages and from the CCDS project page at NCBI. The CCDS set is built by consensus among Ensembl, the National Center for Biotechnology Information (NCBI), and the HUGO Gene Nomenclature Committee (HGNC) for human (link).

Right. Then now we count how many genes annotated with HAVANA and ENSEMBL.

```
f %>% group_by(source) %>% count(source)
```

```
## # A tibble: 2 x 2
## # Groups:   source [2]
##   source     n
##   <chr>    <int>
## 1 ENSEMBL  245185
## 2 HAVANA  1511317
```

## 2.5. do.call

Since the last column `info` contains a long string for multiple annotations, we will need to split it to extract each annotation. For example, the first line for transcript annotation looks like this:

```
head(f)
```

```
## # A tibble: 6 x 9
##   chrom source feature_type start   end score strand phase info
##   <chr> <chr>   <chr>    <dbl> <dbl> <chr> <chr> <chr> <chr>
## 1 chr1  HAVANA gene      11869 14409 .      +      .      "gene_id \"ENSG00000~
## 2 chr1  HAVANA transcript 11869 14409 .      +      .      "gene_id \"ENSG00000~
## 3 chr1  HAVANA exon      11869 12227 .      +      .      "gene_id \"ENSG00000~
## 4 chr1  HAVANA exon      12613 12721 .      +      .      "gene_id \"ENSG00000~
## 5 chr1  HAVANA exon      13221 14409 .      +      .      "gene_id \"ENSG00000~
## 6 chr1  HAVANA transcript 12010 13670 .      +      .      "gene_id \"ENSG00000~
```

If you would like to split `transcript_support_level` and create a new column, you can use `strsplit` function.

```

a = 'chr1      HAVANA      transcript      11869      14409      .      +      .      gene_id "ENSG00000223972.5"; transcript_support_level\\s+')

## [[1]]
## [1] "chr1      HAVANA      transcript      11869      14409      .      +      .      gene_id \"ENSG00000223972.5\";
## [2] "1\"; hgnc_id \"HGNC:37102\"; tag \"basic\"; havana_gene \"OTTHUMG00000000961.2\"; havana_transc

```

After split the string, you can select the second item in the list ([[1]][2]).

```
strsplit(a, 'transcript_support_level\\s+')[[1]][2]
```

```
## [1] "1\"; hgnc_id \"HGNC:37102\"; tag \"basic\"; havana_gene \"OTTHUMG00000000961.2\"; havana_transc
```

You can find the 1 in the first position, which you will need to split again.

```
b = strsplit(a, 'transcript_support_level\\s+')[[1]][2]
strsplit(b, '\\\"')
```

```

## [[1]]
## [1] "1"                      ; hgnc_id "HGNC:37102"
## [4] ";" tag "basic"           ; havana_gene "
## [7] "OTTHUMG00000000961.2"   ; havana_transcript "OTTHUMT00000362751.1"
## [10] ","

```

Now you would like to apply `strsplit` function across vectors. For this, `do.call` function can be easily implemented to `strsplit` over the vectors from one column. Let's try this.

```
head(do.call(rbind.data.frame, strsplit(a, 'transcript_support_level\\s+'))[[2]])
```

```
## [1] "1\"; hgnc_id \"HGNC:37102\"; tag \"basic\"; havana_gene \"OTTHUMG00000000961.2\"; havana_transc
```

Now you can write two lines of codes to process two steps we discussed above.

```

# First filter transcripts and create a data frame.
f2 <- f %>% filter(feature_type == 'transcript')

# Now apply the functions.
f2$transcript_support_level <- as.character(do.call(rbind.data.frame,
                                                 strsplit(f2$info, 'transcript_support_level\\s+')))

f2$transcript_support_level <- as.character(do.call(rbind.data.frame,
                                                 strsplit( f2$transcript_support_level, '\\\"'))[[1]])

```

Now you can check the `strsplit` works.

```
head(f2$transcript_support_level)
```

```
## [1] "1"    "NA"   "NA"   "NA"   "5"    "5"
```

### 3. Exercises

Here I list the questions for your activity. Please note that it is an exercise for tidyverse functions, which you will need to use in your code. In addition, you will need to write an one-line code for each question using pipe `%>%`.

For questions, you should read some information thoroughly, including:  
\* Gene biotype.  
\* 0 or 1 based annotation in GTF, BED format  
\* Why some features have 1 bp length?  
\* What is the meaning of zero-length exons in GENCODE?  
Also fun to have a review for microexons  
\* Transcript support level (TSL)

#### 3.1. Annotation of transcripts in our genome

1. Computes the number of transcripts per gene. What is the mean number of transcripts per gene? What is the quantile (25%, 50%, 75%) for these numbers? Which gene has the greatest number of transcript?

```
library(tidyverse)
d_1 = d %>% filter(type == "transcript") %>% count(gene_id)
head(d_1)
```

```
##           gene_id n
## 1 ENSG00000000003.14 3
## 2 ENSG00000000005.6 1
## 3 ENSG00000000419.12 2
## 4 ENSG00000000457.14 3
## 5 ENSG00000000460.17 5
## 6 ENSG00000000938.13 3
```

```
mean(d_1$n)
```

```
## [1] 1.7861
```

```
quantile(d_1$n,)
```

```
##    0%   25%   50%   75% 100%
##     1     1     1     2    87
```

```
d_1[which.max(d_1$n),]
```

```
##           gene_id n
## 3662 ENSG00000109339.22 87
```

2. Compute the number of transcripts per gene among gene biotypes. For example, compare the number of transcript per gene between protein-coding genes, long noncoding genes, pseudogenes.

```
d %>% filter(type == "transcript" ) %>%
  filter(gene_type %in% c("protein_coding", "lncRNA") |
    (str_detect(gene_type, "pseudogene")))) %>%
  count(gene_type,gene_id) %>% group_by(gene_type) %>% summarise(Mean = mean(n),
  Quantile25 = quantile(n, 0.25),
  Quantile75 = quantile(n, 0.75),
  Max = which.max(n))
```

```

## # A tibble: 19 x 5
##   gene_type      Mean Quantile25 Quantile75   Max
##   <chr>        <dbl>     <dbl>     <dbl> <int>
## 1 IG_C_pseudogene    1         1         1     1
## 2 IG_J_pseudogene    1         1         1     1
## 3 IG_pseudogene      1         1         1     1
## 4 IG_V_pseudogene    1         1         1     1
## 5 lncRNA            1.48      1         1     843
## 6 polymorphic_pseudogene 2.17      1         3    10
## 7 processed_pseudogene 1.00      1         1   2913
## 8 protein_coding     2.90      1         4   3630
## 9 pseudogene          1         1         1     1
## 10 rRNA_pseudogene    1         1         1     1
## 11 TR_J_pseudogene    1         1         1     1
## 12 TR_V_pseudogene    1         1         1     1
## 13 transcribed_processed_pseudogene 1.78      1         2     89
## 14 transcribed_unitary_pseudogene   2.45      2         3    27
## 15 transcribed_unprocessed_pseudogene 2.08      2         3     6
## 16 translated_processed_pseudogene   1         1         1     1
## 17 translated_unprocessed_pseudogene 1         1         1     1
## 18 unitary_pseudogene      1         1         1     1
## 19 unprocessed_pseudogene     1.00      1         1   611

```

- Final task is to compute the number of transcripts per gene per chromosome.

```

d %>% filter(type == "transcript") %>%
  group_by(seqnames) %>%
  count(gene_id) %>%
  dplyr::summarise(`Mean transcript per gene per chromosome` = mean(n),
                  Quantile25 = quantile(n, 0.25),
                  Quantile75 = quantile(n, 0.75),
                  Max = which.max(n))

## # A tibble: 25 x 5
##   seqnames `Mean transcript per gene per chromosome` Quantile25 Quantile75   Max
##   <fct>        <dbl>             <dbl>     <dbl> <int>
## 1 chr1           1.80              1         2     519
## 2 chr2           1.77              1         2   1046
## 3 chr3           1.93              1         2     84
## 4 chr4           1.76              1         2     88
## 5 chr5           1.74              1         2     79
## 6 chr6           1.78              1         2   2849
## 7 chr7           1.76              1         2     349
## 8 chr8           1.75              1         2   1753
## 9 chr9           1.70              1         2     25
## 10 chr10          1.78              1         2     392
## # ... with 15 more rows

```

### 3.2. Gene length in the GENCODE

- What is the average length of human genes?

```
d %>% filter(type == "gene") %>%
  summarise(average_length_of_human_gene = mean(width))
```

```
##   average_length_of_human_gene
## 1                      32629.02
```

2. Is the distribution of gene length differed by autosomal and sex chromosomes? Please calculate the quantiles (0%, 25%, 50%, 75%, 100%) of the gene length for each group.

```
p <- seq(0, 1, 0.25)
d <- d %>%
  mutate(group = case_when(seqnames %in% c("chrX", "chrY") ~ "Sex",
                           TRUE ~ "Autosomal"))
d%>%
  group_by(group) %>%
  summarise(quantile0 = quantile(width, p[1]),
            quantile25 = quantile(width, p[2]),
            quantile50 = quantile(width, p[3]),
            quantile75 = quantile(width, p[4]),
            quantile100 = quantile(width, p[5]))
```

```
## # A tibble: 2 x 6
##   group      quantile0 quantile25 quantile50 quantile75 quantile100
##   <chr>        <dbl>     <dbl>     <dbl>     <dbl>      <dbl>
## 1 Autosomal    1         80        129       222      2473537
## 2 Sex           1         78        127       230      2241765
```

3. Is the distribution of gene length differed by gene biotype? Please calculate the quantiles (0%, 25%, 50%, 75%, 100%) of the gene length for each group.

```
p <- seq(0, 1, 0.25)
d%>%
  group_by(gene_type) %>%
  summarise(quantile0 = quantile(width, p[1]),
            quantile25 = quantile(width, p[2]),
            quantile50 = quantile(width, p[3]),
            quantile75 = quantile(width, p[4]),
            quantile100 = quantile(width, p[5])) %>%
  head(10)
```

```
## # A tibble: 10 x 6
##   gene_type      quantile0 quantile25 quantile50 quantile75 quantile100
##   <chr>        <dbl>     <dbl>     <dbl>     <dbl>      <dbl>
## 1 IG_C_gene      3         92       312.      336       8914
## 2 IG_C_pseudogene 34        293       316       424       5211
## 3 IG_D_gene       3         17        20        31        37
## 4 IG_J_gene       3         38        49        70        176
## 5 IG_J_pseudogene 50        50        55        60        60
## 6 IG_pseudogene   306       306       306       306       306
## 7 IG_V_gene       3         46       296       436      176628
## 8 IG_V_pseudogene 28        235       291       437       792
## 9 lncRNA          1         126       339       2658      1375317
## 10 miRNA         41         70        80        91        180
```

### 3.3. Transcript support levels (TSL)

The GENCODE TSL provides a consistent method of evaluating the level of support that a GENCODE transcript annotation is actually expressed in humans.

1. With transcript, how many transcripts are categorized for each TSL?

```
d %>% filter(type == "transcript" & transcript_support_level != "NA") %>%
  count(transcript_support_level)
```

```
##   transcript_support_level     n
## 1                      1 31801
## 2                      2 13372
## 3                      3  7228
## 4                      4  2245
## 5                      5 13674
```

2. From the first question, please count the number of transcript for each TSL by gene biotype.

```
d %>% filter(type == "transcript" & transcript_support_level != "NA") %>%
  count(transcript_support_level, gene_type) %>% head(10)
```

```
##   transcript_support_level      gene_type     n
## 1                      1          IG_C_gene    1
## 2                      1          lncRNA 1620
## 3                      1 polymorphic_pseudogene 30
## 4                      1      protein_coding 29783
## 5                      1 transcribed_processed_pseudogene 42
## 6                      1 transcribed_unitary_pseudogene 58
## 7                      1 transcribed_unprocessed_pseudogene 267
## 8                      2          lncRNA 2970
## 9                      2 polymorphic_pseudogene  3
## 10                     2      protein_coding 10104
```

3. From the first question, please count the number of transcript for each TSL by source.

```
d %>% filter(type == "transcript" & transcript_support_level != "NA") %>%
  count(transcript_support_level, source)
```

```
##   transcript_support_level source     n
## 1                      1 HAVANA 29434
## 2                      1 ENSEMBL 2367
## 3                      2 HAVANA 12052
## 4                      2 ENSEMBL 1320
## 5                      3 HAVANA 6964
## 6                      3 ENSEMBL 264
## 7                      4 HAVANA 2116
## 8                      4 ENSEMBL 129
## 9                      5 HAVANA 10157
## 10                     5 ENSEMBL 3517
```

### 3.4. CCDS in the GENCODE

- With gene, please create a data frame with the columns - gene\_id, gene\_name, hgnc\_id, gene\_type, chromosome, start, end, and strand. Then, please create new columns for presence of hgnc and ccds. For example, you can put 1 in the column isHgnc, if hgnc annotation is available, or 0 if not. Then, you can put 1 in the column isCCDS, if ccds annotation is available, or 0 if not.

```
ref <- d %>% filter(type == "gene") %>%
  select(gene_id, gene_name, hgnc_id, gene_type, seqnames, start, end, strand, ccdsid) %>%
  mutate(isHgnc = case_when(is.na(hgnc_id) ~ 0,
                            TRUE ~ 1),
         isCCDS = case_when(is.na(ccdsid) ~ 0,
                            TRUE ~ 1))
ref %>% head(10)

##           gene_id   gene_name      hgnc_id          gene_type
## 1 ENSG00000223972.5    DDX11L1 HGNC:37102 transcribed_unprocessed_pseudogene
## 2 ENSG00000227232.5    WASH7P HGNC:38034      unprocessed_pseudogene
## 3 ENSG00000278267.1    MIR6859-1 HGNC:50039        miRNA
## 4 ENSG00000243485.5    MIR1302-2HG HGNC:52482      lncRNA
## 5 ENSG00000284332.1    MIR1302-2 HGNC:35294        miRNA
## 6 ENSG00000237613.2    FAM138A HGNC:32334      lncRNA
## 7 ENSG00000268020.3    OR4G4P HGNC:14822      unprocessed_pseudogene
## 8 ENSG00000240361.2    OR4G11P HGNC:31276 transcribed_unprocessed_pseudogene
## 9 ENSG00000186092.6    OR4F5 HGNC:14825      protein_coding
## 10 ENSG00000238009.6   AL627309.1 <NA>          lncRNA
##   seqnames start     end strand ccdsid isHgnc isCCDS
## 1   chr1 11869 14409      +   <NA>    1     0
## 2   chr1 14404 29570      -   <NA>    1     0
## 3   chr1 17369 17436      -   <NA>    1     0
## 4   chr1 29554 31109      +   <NA>    1     0
## 5   chr1 30366 30503      +   <NA>    1     0
## 6   chr1 34554 36081      -   <NA>    1     0
## 7   chr1 52473 53312      +   <NA>    1     0
## 8   chr1 57598 64116      +   <NA>    1     0
## 9   chr1 65419 71585      +   <NA>    1     0
## 10  chr1 89295 133723     -   <NA>    0     0
```

- Please count the number of hgnc by gene biotypes.

```
ref %>% filter(isHgnc == "1") %>% count(gene_type) %>% head(10)

##       gene_type     n
## 1      IG_C_gene  14
## 2  IG_C_pseudogene   9
## 3      IG_D_gene  37
## 4      IG_J_gene  18
## 5  IG_J_pseudogene   3
## 6      IG_V_gene 142
## 7  IG_V_pseudogene 185
## 8      lncRNA 3970
## 9      miRNA 1856
## 10     misc_RNA 1033
```

3. Please count the number of hgnc by level. Please note that level in this question is not TSL. Please find information in this link: 1 (verified loci), 2 (manually annotated loci), 3 (automatically annotated loci).

```
ref %>%
  mutate(level = d$level[d$type == "gene"]) %>%
  filter(isHgnc == "1") %>% count(level)

##   level     n
## 1      1 11347
## 2      2 20517
## 3      3  5733
```

### 3.5. Transcripts in the GENCODE

1. Which gene has the largest number of transcripts?

```
ref2 <- d %>% filter(type == "transcript") %>%
  count(gene_id)

ref2[which.max(ref2$n),]

##           gene_id   n
## 3662 ENSG00000109339.22 87
```

2. Please calculate the quantiles (0%, 25%, 50%, 75%, 100%) of the gene length for protein coding genes and long noncoding genes.

```
p <- seq(0, 1, 0.25)

d %>%
  filter(type == "gene") %>%
  filter(gene_type %in% c("protein_coding", "lncRNA")) %>%
  group_by(gene_type) %>%
  summarise(quintile(width,p))

## `summarise()` has grouped output by 'gene_type'. You can override using the '.groups' argument.

## # A tibble: 10 x 2
## # Groups:   gene_type [2]
##   gene_type   `quintile(width, p)`
##   <chr>          <dbl>
## 1 lncRNA          68
## 2 lncRNA         1874.
## 3 lncRNA         6272.
## 4 lncRNA        24774.
## 5 lncRNA       1375317
## 6 protein_coding    117
## 7 protein_coding   9632.
## 8 protein_coding  27212
## 9 protein_coding  70809
## 10 protein_coding 2473537
```

3. Please count the number of transcripts by chromosomes.

```
d %>% filter(type == "transcript") %>%
  group_by(seqnames) %>%
  count()

## # A tibble: 25 x 2
## # Groups:   seqnames [25]
##   seqnames     n
##   <fct>    <int>
##   1 chr1      9827
##   2 chr2      7432
##   3 chr3      6157
##   4 chr4      4662
##   5 chr5      5203
##   6 chr6      5455
##   7 chr7      5292
##   8 chr8      4350
##   9 chr9      3949
##  10 chr10     4157
## # ... with 15 more rows
```

### 3.6. Autosomal vs. Sex chromosomes.

1. Please calculate the number of genes per chromosome.

```
d %>% filter(type == "gene") %>%
  group_by(seqnames) %>%
  count() %>%
  head(10)
```

```
## # A tibble: 10 x 2
## # Groups:   seqnames [10]
##   seqnames     n
##   <fct>    <int>
##   1 chr1      5471
##   2 chr2      4196
##   3 chr3      3185
##   4 chr4      2651
##   5 chr5      2983
##   6 chr6      3059
##   7 chr7      3014
##   8 chr8      2482
##   9 chr9      2327
##  10 chr10     2332
```

2. Please compare the number of genes between autosomal and sex chromosome (Mean, Median).

```
ref2 <- d %>% mutate(group = case_when(seqnames %in% c("chrX", "chrY") ~ "Sex",
                                         TRUE ~ "Autosomal")) %>%
  group_by(group, seqnames) %>%
```

```

filter(type == "gene") %>% count()

ref2 %>%
  group_by(group) %>%
  dplyr::summarise(Mean = mean(n), median(n))

## # A tibble: 2 x 3
##   group      Mean `median(n)`
##   <chr>     <dbl>     <dbl>
## 1 Autosomal 2505.     2556
## 2 Sex        1494.    1494.

```

3. Please divide the genes into groups ‘protein coding’ and ‘long noncoding’, and then compare the number of genes in each chromosomes within groups.

```

d %>% filter(type == "gene" &
               gene_type %in% c("protein_coding", "lncRNA")) %>%
  group_by(gene_type, seqnames) %>%
  count() %>%
  head()

## # A tibble: 6 x 3
## # Groups:   gene_type, seqnames [6]
##   gene_type seqnames     n
##   <chr>     <fct>     <int>
## 1 lncRNA    chr1      1416
## 2 lncRNA    chr2      1241
## 3 lncRNA    chr3       861
## 4 lncRNA    chr4       790
## 5 lncRNA    chr5      950
## 6 lncRNA    chr6      826

```