

# [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.

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.

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/genocode/Gencode_human/release_31/genocode.v31.basic.annotation.
# Once you downloaded the file, you won't need to download it again. So please comment out the command
```

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 in 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  3.1.5      v dplyr  1.0.7  
## v tidyr   1.1.4      v stringr 1.4.0  
## v readr   2.0.2      v forcats 0.5.1  
  
## -- Conflicts ----- tidyverse_conflicts() --  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()     masks stats::lag()  
  
d = 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.
```

```
#delim: what is it seperated by. in this case, tab.  
#skip 5: header.
```

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")
```

```
## Bioconductor version 3.13 (BiocManager 1.30.16), R 4.1.1 (2021-08-10)
```

```
## Warning: package(s) not installed when version(s) same as current; use `force = TRUE` to  
## re-install: 'rtracklayer'
```

```
## Old packages: 'htmlTable', 'lubridate', 'S4Vectors', 'xfun'
```

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

```
library(tidyverse)  
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 OTTHUMT00000002844.2      <NA>      <NA> PGO:0000019
##          protein_id ccidsid
## 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.

```

d = 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(d)
```

```
## # 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 \"ENSG000000~
## 2 chr1  HAVANA  transcript    11869 14409 .    +    .    "gene_id \"ENSG000000~
## 3 chr1  HAVANA  exon         11869 12227 .    +    .    "gene_id \"ENSG000000~
## 4 chr1  HAVANA  exon         12613 12721 .    +    .    "gene_id \"ENSG000000~
## 5 chr1  HAVANA  exon         13221 14409 .    +    .    "gene_id \"ENSG000000~
## 6 chr1  HAVANA  transcript    12010 13670 .    +    .    "gene_id \"ENSG000000~
```

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

```
summary(d)
```

```
##      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 `group_by` function.

```
library(tidyverse)
d %>% 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).

```
#There are 8 feature types.
# selenocysteine: analogue of the more common cysteine with selenium in place of the sulfur.
```

## 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.

```
d1 = filter(d, 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.

```
d %>% 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:

```
#chr1      HAVANA      transcript      11869      14409      .      +      .      gene_id "ENSG00000223972.5"; transcrip
```

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"; transcrip
strsplit(a, 'transcript_support_level\\s+')
## [[1]]
## [1] "chr1      HAVANA      transcript      11869      14409      .      +      .      gene_id \"ENSG00000223972.5\"; transcrip
## [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, '\\s+')
## [[1]]
## [1] "1"
## [4] "; tag "
## [7] "OTTHUMG00000000961.2" "; havana_transcript " "OTTHUMT00000362751.1"
## [10] ";"
```

From this, you will get the first item in the list (`[[1]][1]`). 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
```

```
#rbind -> binding two data frames together
```

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

```
# First filter transcripts and create a data frame.
d2 <- d %>% filter(feature_type == 'transcript')
# Now apply the functions.
d2$transcript_support_level <- as.character(do.call(rbind.data.frame,
                                                    strsplit(d2$info, 'transcript_support_level\\s+')))

d2$transcript_support_level <- as.character(do.call(rbind.data.frame,
                                                    strsplit( d2$transcript_support_level, '\\s'))[[1]]
```

Now you can check the `strsplit` works.

```
head(d2$transcript_support_level)
```

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

You can use the same method to extract other annotations, like `gene_id`, `gene_name` etc.

```
#Extracting gene_name
# Now apply the functions.
d2$gene_name <- as.character(do.call(rbind.data.frame,
                                     strsplit(d2$info, 'gene_name\\s+'))[[2]])

d2$gene_name <- as.character(do.call(rbind.data.frame,
                                     strsplit( d2$gene_name, '\\s'))[[1]])
head(d2$gene_name)
```

```
## [1] "DDX11L1" "DDX11L1" "WASH7P" "MIR6859-1" "MIR1302-2HG"
## [6] "MIR1302-2HG"
```

### 3. Exercises

Here I list the questions for group activity. You will need to pick up one session for three questions for your group. It will be your quiz on the next class. If you have done your session, you can of course go ahead and take other sessions for your practice. 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)



```

# Before starting, I would like to create the columns for all the informations that are needed

#gene_id
d$gene_id <- as.character(do.call(rbind.data.frame,
                                strsplit(d$info, 'gene_id\\s+'))[[2]])

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

#gene name
d$gene_name <- as.character(do.call(rbind.data.frame,
                                    strsplit(d$info, 'gene_name\\s+'))[[2]])

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

#hgnc_id
d$hgnc_id <- as.character(do.call(rbind.data.frame,
                                  strsplit(d$info, 'hgnc_id\\s+'))[[2]])

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

#gene_type
d$gene_type <- as.character(do.call(rbind.data.frame,
                                    strsplit(d$info, 'gene_type\\s+'))[[2]])

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

#CCDS id
d$ccdsid <- as.character(do.call(rbind.data.frame,
                                 strsplit(d$info, 'ccdsid\\s+'))[[2]])

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

#TSL
d$TSL <- as.character(do.call(rbind.data.frame,
                              strsplit(d$info, 'transcript_support_level\\s+'))[[2]])

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

View(d)

```

##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?

```

d_trans <- d %>% filter(feature_type == 'transcript')
transcripts_per_gene <- d_trans %>% group_by(gene_id) %>% count()
mean(transcripts_per_gene$n)

```

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

```
quantile(transcripts_per_gene$n)
```

```
##    0%   25%   50%   75%  100%  
##     1     1     1     2    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. Final task is to compute the number of transcripts per gene per chromosome.

```
transcripts_per_biotype <- d_trans %>% group_by(gene_type) %>% count()  
transcripts_per_biotype
```

```
## # A tibble: 40 x 2  
## # Groups:   gene_type [40]  
##   gene_type      n  
##   <chr>      <int>  
## 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_pseudogene     1  
## 7 IG_V_gene     144  
## 8 IG_V_pseudogene  188  
## 9 lncRNA      24993  
## 10 miRNA      1881  
## # ... with 30 more rows
```

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

```
transcripts_per_chromosome <- d_trans %>% group_by(chrom) %>% count()  
transcripts_per_chromosome
```

```
## # A tibble: 25 x 2  
## # Groups:   chrom [25]  
##   chrom      n  
##   <chr> <int>  
## 1 chr1   9827  
## 2 chr10  4157  
## 3 chr11  6265  
## 4 chr12  5612  
## 5 chr13  2209  
## 6 chr14  3958  
## 7 chr15  3863  
## 8 chr16  4639  
## 9 chr17  5730  
## 10 chr18  2133  
## # ... with 15 more rows
```

## 3.2. Gene length in the GENCODE

1. What is the average length of human genes?

```
#get gene length
d <- d %>% mutate(length = abs(end - start))
#filter out only genes
d_gene <- d %>% filter(feature_type == 'gene')
d_gene %>% summarise(mean = mean(length))
```

```
## # A tibble: 1 x 1
##   mean
##   <dbl>
## 1 32628.
```

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. 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.

```
#Autosomal vs. Sex
d <- d %>%
  mutate(chr_type = ifelse(chrom %in% c('chrX','chrY'), 'sex', 'aut'))
View(d)
d_gene <- d %>% filter(feature_type == 'gene')
d_gene %>%
  filter(chr_type %in% c('aut', 'sex')) %>%
  group_by(chr_type) %>%
  summarise(length_dist = quantile(length))
```

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

```
## # A tibble: 10 x 2
## # Groups:   chr_type [2]
##   chr_type length_dist
##   <chr>      <dbl>
## 1 aut         7
## 2 aut        563
## 3 aut       3768.
## 4 aut      25787.
## 5 aut     2473536
## 6 sex         47
## 7 sex        472
## 8 sex       1911
## 9 sex      13501
## 10 sex     2241764
```

```
#length distribution in biotype
d_gene %>%
  group_by(gene_type) %>%
  summarise(distribution = quantile(length))
```

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

```
## # A tibble: 200 x 2
## # Groups:   gene_type [40]
##   gene_type      distribution
##   <chr>          <dbl>
## 1 IG_C_gene      440
## 2 IG_C_gene     476.
## 3 IG_C_gene    4588.
## 4 IG_C_gene    5478.
## 5 IG_C_gene    8913
## 6 IG_C_pseudogene 247
## 7 IG_C_pseudogene 312
## 8 IG_C_pseudogene 316
## 9 IG_C_pseudogene 733
## 10 IG_C_pseudogene 5210
## # ... with 190 more rows
```

### 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_trans %>% group_by(TSL) %>% count()
```

```
## # A tibble: 7 x 2
## # Groups:   TSL [7]
##   TSL      n
##   <chr>  <int>
## 1 "1"    31801
## 2 "2"    13372
## 3 "3"     7228
## 4 "4"     2245
## 5 "5"    13674
## 6 "gene_id " 12080
## 7 "NA"    27843
```

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

```
d_trans %>% group_by(TSL, gene_type) %>% count()
```

```
## # A tibble: 91 x 3
## # Groups:   TSL, gene_type [91]
##   TSL gene_type      n
##   <chr> <chr>      <int>
## 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
## # ... with 81 more rows
```

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

```
d_trans %>% group_by(TSL, source) %>% count()
```

```
## # A tibble: 14 x 3
## # Groups:   TSL, source [14]
##   TSL      source      n
##   <chr>    <chr>  <int>
## 1 "1"      ENSEMBL  2367
## 2 "1"      HAVANA  29434
## 3 "2"      ENSEMBL  1320
## 4 "2"      HAVANA  12052
## 5 "3"      ENSEMBL   264
## 6 "3"      HAVANA  6964
## 7 "4"      ENSEMBL   129
## 8 "4"      HAVANA  2116
## 9 "5"      ENSEMBL  3517
## 10 "5"     HAVANA  10157
## 11 "gene_id " ENSEMBL   179
## 12 "gene_id " HAVANA  11901
## 13 "NA"      ENSEMBL   7881
## 14 "NA"      HAVANA  19962
```

### 3.4. CCDS in the GENCODE

1. 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.

```
#isHGNC, isCCDS
d <- d %>% mutate(isCCDS = ifelse(ccdsid == 'gene_id ', 0, 1))
d <- d %>% mutate(isHGNC = ifelse(hgnc_id == 'gene_id ', 0, 1))

#creating data frame
gene <- d %>% select(gene_id, gene_name, hgnc_id, gene_type, chrom, start, end, strand, isHGNC, isCCDS)
```

Please count the number of hgnc by gene biotypes.

```
gene %>% filter(isHGNC == 1) %>% group_by(gene_type) %>% count(hgnc_id)
```

```
## # A tibble: 37,541 x 3
## # Groups:   gene_type [36]
##   gene_type hgnc_id      n
```

```
##      <chr>      <chr>      <int>
##  1 IG_C_gene HGNC:5478      14
##  2 IG_C_gene HGNC:5479      14
##  3 IG_C_gene HGNC:5480      16
##  4 IG_C_gene HGNC:5522      16
##  5 IG_C_gene HGNC:5525      16
##  6 IG_C_gene HGNC:5526      16
##  7 IG_C_gene HGNC:5527      22
##  8 IG_C_gene HGNC:5528      16
##  9 IG_C_gene HGNC:5541      16
## 10 IG_C_gene HGNC:5716       6
## # ... with 37,531 more rows
```

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).

```
#creating level
d$level <- as.character(do.call(rbind.data.frame,
                                strsplit(d$info, "level\\s+"))[[2]])

d$level <- as.character(do.call(rbind.data.frame,
                                strsplit(d$level, ';'))[[1]])

#add to gene data frame
gene <- gene %>% mutate(level = d$level)

#count
gene %>% filter(isHGNC == 1) %>% group_by(hgnc_id) %>% count(level)
```

```
## # A tibble: 49,078 x 3
## # Groups:   hgnc_id [37,535]
##   hgnc_id   level     n
##   <chr>      <chr> <int>
##  1 HGNC:100    1         26
##  2 HGNC:100    2         58
##  3 HGNC:10000  2         91
##  4 HGNC:10001  1         16
##  5 HGNC:10001  2         43
##  6 HGNC:10001  3         17
##  7 HGNC:10002  1         37
##  8 HGNC:10002  2        346
##  9 HGNC:10002  3        107
## 10 HGNC:10003  2        156
## # ... with 49,068 more rows
```

### 3.5. Transcripts in the GENCODE

1. Which gene has the largest number of transcripts?

```
d_trans %>% group_by(gene_id) %>% count(sort = TRUE)
```

```
## # A tibble: 60,603 x 2
```

```
## # Groups:   gene_id [60,603]
##   gene_id      n
##   <chr>      <int>
## 1 ENSG00000109339.22    87
## 2 ENSG00000075711.21    39
## 3 ENSG00000156113.23    39
## 4 ENSG00000168036.18    39
## 5 ENSG00000224078.15    38
## 6 ENSG00000169255.15    37
## 7 ENSG00000127990.18    34
## 8 ENSG00000285219.2     34
## 9 ENSG00000126091.20    33
## 10 ENSG00000196628.17   32
## # ... with 60,593 more rows
```

*ENSG00000109339.22 is the gene that has the most transcripts*

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

```
d_gene %>%
  filter(gene_type == c('protein_coding', 'lncRNA')) %>%
  group_by(gene_type) %>%
  summarize(length_quantile = quantile(length))
```

```
## Warning in gene_type == c("protein_coding", "lncRNA"): longer object length is
## not a multiple of shorter object length
```

```
## `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      length_quantile
##   <chr>      <dbl>
## 1 lncRNA          67
## 2 lncRNA        1888
## 3 lncRNA        6378
## 4 lncRNA       24881
## 5 lncRNA     1117566
## 6 protein_coding    116
## 7 protein_coding   9485.
## 8 protein_coding  26788
## 9 protein_coding  70564
## 10 protein_coding 2473536
```

3. Please count the number of transcripts by chromosomes.

```
d_trans %>%
  group_by(chrom) %>%
  count()
```

```
## # A tibble: 25 x 2
## # Groups:   chrom [25]
##   chrom      n
##   <chr> <int>
## 1 chr1    9827
## 2 chr10   4157
## 3 chr11   6265
## 4 chr12   5612
## 5 chr13   2209
## 6 chr14   3958
## 7 chr15   3863
## 8 chr16   4639
## 9 chr17   5730
## 10 chr18  2133
## # ... with 15 more rows
```

### 3.6. Autosomal vs. Sex chromosomes.

1. Please calculate the number of genes per chromosome.

```
d_gene %>% group_by(chrom) %>% count()
```

```
## # A tibble: 25 x 2
## # Groups:   chrom [25]
##   chrom      n
##   <chr> <int>
## 1 chr1    5471
## 2 chr10   2332
## 3 chr11   3360
## 4 chr12   3054
## 5 chr13   1397
## 6 chr14   2282
## 7 chr15   2221
## 8 chr16   2556
## 9 chr17   3060
## 10 chr18  1242
## # ... with 15 more rows
```

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

```
gene_num <- d_gene %>%
  filter(chr_type %in% c('aut', 'sex')) %>%
  group_by(chr_type, chrom) %>%
  count()

gene_num %>% group_by(chr_type) %>% summarise(mean = mean(n), median = median(n))
```

```
## # A tibble: 2 x 3
##   chr_type mean median
##   <chr>    <dbl> <dbl>
## 1 aut      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_gene %>%
  filter(feature_type == 'gene' & gene_type %in% c('protein_coding', 'lncRNA')) %>%
  group_by(gene_type, chrom) %>%
  count()
```

```
## # A tibble: 49 x 3
## # Groups:   gene_type, chrom [49]
##   gene_type chrom      n
##   <chr>      <chr> <int>
## 1 lncRNA     chr1    1416
## 2 lncRNA     chr10    695
## 3 lncRNA     chr11    798
## 4 lncRNA     chr12    938
## 5 lncRNA     chr13    480
## 6 lncRNA     chr14    593
## 7 lncRNA     chr15    689
## 8 lncRNA     chr16    842
## 9 lncRNA     chr17    875
## 10 lncRNA    chr18    522
## # ... with 39 more rows
```