# Tutorial_2_Human Genome Annotation_2020250137

## [Tutorial] Human Genome Annotation

Tutorial for Tidyverse(Chapter 4)

## 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.

```
#wget ftp://ftp.ebi.ac.uk/pub/databases/gencode/Gencode_human/release_31/gencode.v31.basic.annotation.g
```

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

```
#zcat 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  3.1.3     v dplyr   1.0.7
## v tidyr   1.1.3     v stringr 1.4.0
## v readr   2.0.1     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.
```

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 OTTHUMT00000002844.2        <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.

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.

```r
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.

```r
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 \"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(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 _____ function.

```
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.

#8 feature-types. 567862 items in CDS(CoDing Sequence: The basic unit of gene sequence that can be translated from gene to protein.), 74485 items in exon, 60603 items in gene, 96 items in Selenocysteine(21st proteinogenic amino acid. It is an analogue of the more common cystein with selenium in place of the sulfur.), 57886 items in start_codon, 57775 items in stop_codon, 108243 items in transcript and 159202 items in UTR.

## 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 row 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')
head(d1)
```

```
## # 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 \"ENSG0000~
## 2 chr1  HAVANA  gene         14404 29570 .     -      .     "gene_id \"ENSG0000~
## 3 chr1  ENSEMBL gene         17369 17436 .     -      .     "gene_id \"ENSG0000~
## 4 chr1  HAVANA  gene         29554 31109 .     +      .     "gene_id \"ENSG0000~
## 5 chr1  ENSEMBL gene         30366 30503 .     +      .     "gene_id \"ENSG0000~
## 6 chr1  HAVANA  gene         34554 36081 .     -      .     "gene_id \"ENSG0000~
```

```
#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 heterogeniety inherent in annotation databases.

The GENCODE project utilizes two sources of gene annotation.

1. Havana: Manual gene annotation
2. Ensembl: Automatic gene annotation

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.

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"; transcri
```

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"; transc

strsplit(a,'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] ";"
```

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 imple-
mented 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.

```
d2 <- d %>% filter(feature_type=="transcript")
```

```
d2$transcript_support_level <-as.character(do.call(rbind.data.frame, strsplit(d2$info, 'transcript_supp
d2$transcript_support_level <- as.character(do.call(rbind.data.frame,
```

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.

```
gene_id<-as.character(do.call(rbind.data.frame, strsplit(d2$info, 'gene_id\\s+"'))[[2]])
gene_id<-as.character(do.call(rbind.data.frame, strsplit( gene_id, '\\"'))[[1]])
```

```
head(gene_id)
```

```
## [1] "ENSG00000223972.5" "ENSG00000223972.5" "ENSG00000227232.5"
## [4] "ENSG00000278267.1" "ENSG00000243485.5" "ENSG00000243485.5"
```

```
gene_name<-as.character(do.call(rbind.data.frame, strsplit(d2$info, 'gene_name\\s+"'))[[2]])
gene_name<-as.character(do.call(rbind.data.frame, strsplit( gene_name, '\\"'))[[1]])
```

```
head(gene_name)
```

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

```
transcript_id<-as.character(do.call(rbind.data.frame, strsplit(d2$info, 'transcript_id\\s+"'))[[2]])
transcript_id<-as.character(do.call(rbind.data.frame, strsplit( transcript_id, '\\"'))[[1]])
```

```
head(transcript_id)
```

```
## [1] "ENST00000456328.2" "ENST00000450305.2" "ENST00000488147.1"
## [4] "ENST00000619216.1" "ENST00000473358.1" "ENST00000469289.1"
```

# 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)

```
d_10 = rtracklayer::import('gencode.v31.basic.annotation.gtf.gz')
d_10=d_10 %>% as.data.frame()
```

## 3.1 Annotation of transcripts in our genome

1. Compute 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?

```
gene_id2<-as.data.frame(gene_id)
```

```
gene_id_n<-gene_id2 %>% group_by(gene_id) %>% count()
gene_id_n
```

```
## # A tibble: 60,603 x 2
## # Groups:   gene_id [60,603]
##    gene_id              n
##    <chr>            <int>
##  1 ENSG00000000003.14   3
##  2 ENSG00000000005.6    1
##  3 ENSG00000000419.12   2
##  4 ENSG00000000457.14   3
##  5 ENSG00000000460.17   5
##  6 ENSG00000000938.13   3
##  7 ENSG00000000971.15   3
##  8 ENSG00000001036.13   1
##  9 ENSG00000001084.13   5
## 10 ENSG00000001167.14   2
## # ... with 60,593 more rows
```

```
m<-mean(gene_id_n$n)
m
```

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

```
quantile(gene_id_n$n,c(0.25,0.5,0.75))
```

```
## 25% 50% 75%
##   1   1   2
```

```
max(gene_id_n$n)
```

```
## [1] 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.

```
gene_type<-as.character(do.call(rbind.data.frame, strsplit(d2$info, 'gene_type\\s+"')))[[2]])
gene_type<-as.character(do.call(rbind.data.frame, strsplit( gene_type, '\\"')))[[1]])
```

```
d2$gene_type<-gene_type
d2$gene_id<-gene_id
d2 %>% filter(gene_type%in%c("pseudogene","protein_coding","lncRNA")) %>% group_by(gene_type,gene_id) %>%
```

```
## # A tibble: 36,833 x 3
## # Groups:   gene_type, gene_id [36,833]
##    gene_type gene_id              n
##    <chr>     <chr>            <int>
##  1 lncRNA    ENSG00000082929.8    2
##  2 lncRNA    ENSG00000083622.8    1
##  3 lncRNA    ENSG00000093100.13   1
##  4 lncRNA    ENSG00000099869.7    1
##  5 lncRNA    ENSG00000103472.10   1
##  6 lncRNA    ENSG00000115934.11   1
##  7 lncRNA    ENSG00000116652.6    1
##  8 lncRNA    ENSG00000116883.8    1
##  9 lncRNA    ENSG00000117242.7    1
## 10 lncRNA    ENSG00000120664.11   2
## # ... with 36,823 more rows
```

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

```
d2 %>% group_by(chrom,gene_id) %>% count()
```

```
## # A tibble: 60,603 x 3
## # Groups:   chrom, gene_id [60,603]
##    chrom gene_id              n
##    <chr> <chr>            <int>
##  1 chr1  ENSG00000000457.14   3
```

```
##  2 chr1  ENSG00000000460.17      5
##  3 chr1  ENSG00000000938.13      3
##  4 chr1  ENSG00000000971.15      3
##  5 chr1  ENSG00000001460.18      3
##  6 chr1  ENSG00000001461.17      4
##  7 chr1  ENSG00000004455.16      8
##  8 chr1  ENSG00000004487.16      3
##  9 chr1  ENSG00000006555.11      2
## 10 chr1  ENSG00000007341.19      6
## # ... with 60,593 more rows
```

## 3.2 Gene length in the GENCODE

1. What is the average length of human genes?

```
mean(d_10$width)
```

```
## [1] 4069.518
```

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.

```
dfilter<-d_10 %>% filter(type=="transcript")
d_sex_chromosomes<-dfilter %>% filter(seqnames %in% c('chrX','chrY'))
d_autosomal<-dfilter %>% filter(seqnames!=c('chrX','chrY','chrM'))

sex_chromosome<-quantile(d_sex_chromosomes$width,c(0,0.25,0.5,0.75,1))
autosomal<-quantile(d_autosomal$width,c(0,0.25,0.5,0.75,1))
data.frame(sex_chromosome,autosomal)
```

```
##       sex_chromosome  autosomal
## 0%             48.00       8.00
## 25%          1005.00    1710.00
## 50%          6309.50   11238.50
## 75%         36558.75   41822.25
## 100%      2092292.00 2471657.00
```

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.

```
dfilter %>% group_by(gene_type) %>% summarize(distribution = quantile(width,c(0,0.25,0.5,0.75,1)))
```

```
## `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               441
##  2 IG_C_gene               477.
```

```
##  3 IG_C_gene              4590.
##  4 IG_C_gene              5479.
##  5 IG_C_gene              7278
##  6 IG_C_pseudogene         248
##  7 IG_C_pseudogene         313
##  8 IG_C_pseudogene         317
##  9 IG_C_pseudogene         734
## 10 IG_C_pseudogene        5211
## # ... 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?

```
d2 %>% group_by(transcript_support_level) %>% count()
```

```
## # A tibble: 7 x 2
## # Groups:   transcript_support_level [7]
##   transcript_support_level      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.

```
d2 %>% group_by(gene_type,transcript_support_level) %>% count()
```

```
## # A tibble: 91 x 3
## # Groups:   gene_type, transcript_support_level [91]
##    gene_type        transcript_support_level      n
##    <chr>            <chr>                      <int>
##  1 IG_C_gene        "1"                            1
##  2 IG_C_gene        "5"                            1
##  3 IG_C_gene        "gene_id "                     5
##  4 IG_C_gene        "NA"                           7
##  5 IG_C_pseudogene  "NA"                           9
##  6 IG_D_gene        "NA"                          37
##  7 IG_J_gene        "NA"                          18
##  8 IG_J_pseudogene  "NA"                           3
##  9 IG_pseudogene    "NA"                           1
## 10 IG_V_gene        "5"                            3
## # ... with 81 more rows
```

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

```
d2 %>% group_by(source,transcript_support_level) %>% count()
```

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

```
gene<-d_10 %>% select(gene_id,gene_name, gene_type,hgnc_id, seqnames,start, end, strand,type,ccdsid)
gene1<-gene %>% filter(type=="transcript"& gene_type=="protein_coding")
gene1$isHgnc<-ifelse(is.na(gene1$hgnc_id)!="TRUE",1,0)
gene1$isCCDS<-ifelse(is.na(gene1$ccdsid)!="TRUE",1,0)
head(gene1)
```

```
##              gene_id gene_name     gene_type    hgnc_id seqnames  start    end
## 1   ENSG00000186092.6     OR4F5 protein_coding HGNC:14825     chr1  65419  71585
## 2   ENSG00000186092.6     OR4F5 protein_coding HGNC:14825     chr1  69055  70108
## 3   ENSG00000284733.1    OR4F29 protein_coding HGNC:31275     chr1 450703 451697
## 4   ENSG00000284662.1    OR4F16 protein_coding HGNC:15079     chr1 685679 686673
## 5 ENSG00000187634.12    SAMD11 protein_coding HGNC:28706     chr1 925731 944574
## 6 ENSG00000187634.12    SAMD11 protein_coding HGNC:28706     chr1 925741 944581
##   strand       type     ccdsid isHgnc isCCDS
## 1      + transcript       <NA>      1      0
## 2      + transcript CCDS30547.1      1      1
## 3      - transcript CCDS72675.1      1      1
## 4      - transcript CCDS41221.1      1      1
## 5      + transcript    CCDS2.2      1      1
## 6      + transcript       <NA>      1      0
```

2. Please count the number of hgnc by gene biotypes.

```
gene %>% group_by(gene_type) %>% count(hgnc_id)
```

```
## # A tibble: 37,565 x 3
## # Groups:   gene_type [40]
##    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 37,555 more rows
```

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).

```
d_level<-ifelse(is.na(d_10$hgnc_id)!="TRUE",1,0)
d_10 %>% filter(d_level=="1") %>% group_by(level) %>% count()
```

```
## # A tibble: 3 x 2
## # Groups:   level [3]
##   level       n
##   <chr>   <int>
## 1 1      107054
## 2 2     1279964
## 3 3      237265
```

## 3.5 Transcripts in the GENCODE

1. Which gene has the largest number of transcripts?

```
max(gene_id_n$n)
```

```
## [1] 87
```

```
gene_id_n %>% filter(n=="87")
```

```
## # A tibble: 1 x 2
## # Groups:   gene_id [1]
##   gene_id                n
##   <chr>              <int>
## 1 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.

```
length1<-d_10 %>%  select("gene_type","width") %>% filter(gene_type %in% c('protein_coding','lncRNA'))
length1 %>% summarize(distribution = quantile(width,c(0,0.25,0.5,0.75,1)))
```

```
## `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       distribution
##     <chr>                  <dbl>
##  1 lncRNA                     1
##  2 lncRNA                   126
##  3 lncRNA                   339
##  4 lncRNA                  2658
##  5 lncRNA               1375317
##  6 protein_coding             1
##  7 protein_coding            76
##  8 protein_coding           123
##  9 protein_coding           193
## 10 protein_coding       2473537
```

3. Please count the number of transcripts by chromosomes.

```
d2 %>% group_by(chrom,gene_id) %>% count()
```

```
## # A tibble: 60,603 x 3
## # Groups:   chrom, gene_id [60,603]
##     chrom gene_id                 n
##     <chr> <chr>               <int>
##  1 chr1  ENSG00000000457.14      3
##  2 chr1  ENSG00000000460.17      5
##  3 chr1  ENSG00000000938.13      3
##  4 chr1  ENSG00000000971.15      3
##  5 chr1  ENSG00000001460.18      3
##  6 chr1  ENSG00000001461.17      4
##  7 chr1  ENSG00000004455.16      8
##  8 chr1  ENSG00000004487.16      3
##  9 chr1  ENSG00000006555.11      2
## 10 chr1  ENSG00000007341.19      6
## # ... with 60,593 more rows
```

## 3.6 Autosomal vs. Sex chromosomes

1. Please calculate the number of genes per chromosome.

```
d_10 %>% group_by(seqnames) %>% count(gene_name)
```

```
## # A tibble: 59,479 x 3
## # Groups:   seqnames [25]
##    seqnames gene_name       n
##    <fct>    <chr>       <int>
##  1 chr1     A3GALT2        15
##  2 chr1     AADACL3        14
##  3 chr1     AADACL4        13
##  4 chr1     ABCA4         210
##  5 chr1     ABCB10         32
##  6 chr1     ABCD3          74
##  7 chr1     ABHD17AP1       4
##  8 chr1     ABHD17AP3       4
##  9 chr1     ABL2          227
## 10 chr1     AC000032.1      4
## # ... with 59,469 more rows
```

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

```
gene_sex<-d_10 %>%select(seqnames,gene_name) %>%  filter(seqnames %in% c('chrX','chrY'))
gene_autosomal<-d_10 %>% filter(seqnames!=c('chrX','chrY','chrM'))
```

```
## Warning in `!=.default`(seqnames, c("chrX", "chrY", "chrM")):
##
```

```
## Warning in is.na(e1) | is.na(e2):
```

```
sex_chromosome2<-gene_sex %>% count(gene_name) %>% summarise(Mean_sex=mean(n), Median_sex=median(n))
autosomal2<-gene_autosomal %>% count(gene_name) %>% summarise(Mean_autosomal=mean(n), Median_autosomal=
data.frame(sex_chromosome2,autosomal2)
```

```
##   Mean_sex Median_sex Mean_autosomal Median_autosomal
## 1 25.09771          4       29.33634                5
```

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_10 %>%  select("gene_type","gene_name","seqnames") %>% filter(gene_type %in% c('protein_coding','lncR
```

```
## # A tibble: 36,797 x 4
## # Groups:   gene_type, seqnames [49]
##    gene_type seqnames gene_name       n
##    <chr>     <fct>    <chr>       <int>
##  1 lncRNA    chr1     AC000032.1      4
##  2 lncRNA    chr1     AC004865.2      5
##  3 lncRNA    chr1     AC011700.1      3
##  4 lncRNA    chr1     AC091614.1      4
##  5 lncRNA    chr1     AC091614.2      6
##  6 lncRNA    chr1     AC092017.4      7
##  7 lncRNA    chr1     AC092265.1      4
##  8 lncRNA    chr1     AC092783.1      6
##  9 lncRNA    chr1     AC092800.1      4
## 10 lncRNA    chr1     AC092801.1      7
## # ... with 36,787 more rows
```