

# Human Genome Annotation

## [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(<http://www.biolyse.ca/what-is-gene-annotation-in-bioinformatics/>).

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/gencode/Gencode_human/release_31/gencode.v31.basic.annotation.gtf.gz

# 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 is a part of tidyverse, so you can just load tidyverse to use readr functions.

```
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      v forcats 0.5.1
```

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

```
library(readr)
```

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)
d = read_delim('gencode.v31.basic.annotation.gtf.gz', delim = '\t', skip = 5, progress = F, col_names = 1)
```

```
## 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.
```

```
d
```

```
## # 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 \\"ENSG00000022~
## 2 chr1  HAVANA  transcript 11869 14409 .      +      .      "gene_id \\"ENSG00000022~
## 3 chr1  HAVANA  exon      11869 12227 .      +      .      "gene_id \\"ENSG00000022~
## 4 chr1  HAVANA  exon      12613 12721 .      +      .      "gene_id \\"ENSG00000022~
## 5 chr1  HAVANA  exon      13221 14409 .      +      .      "gene_id \\"ENSG00000022~
## 6 chr1  HAVANA  transcript 12010 13670 .      +      .      "gene_id \\"ENSG00000022~
## 7 chr1  HAVANA  exon      12010 12057 .      +      .      "gene_id \\"ENSG00000022~
## 8 chr1  HAVANA  exon      12179 12227 .      +      .      "gene_id \\"ENSG00000022~
## 9 chr1  HAVANA  exon      12613 12697 .      +      .      "gene_id \\"ENSG00000022~
## 10 chr1  HAVANA  exon      12975 13052 .      +      .      "gene_id \\"ENSG00000022~
## # ... 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")
```

```
## 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'
```

```
## Installation paths not writeable, unable to update packages
## path: C:/Program Files/R/R-4.1.1/library
## packages:
## lattice, mgcv, nlme, survival
```

```
## Old packages: 'digest', 'lubridate', 'readr', 'stringi', 'tibble', 'tidyr',
## 'xfun'
```

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
## 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 ENSE000002234944.1      <NA>
## 4 basic OTTHUMT00000362751.1      2 ENSE000003582793.1      <NA>
## 5 basic OTTHUMT00000362751.1      3 ENSE000002312635.1      <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 ([https://www.gencodegenes.org/pages/data\\_format.html](https://www.gencodegenes.org/pages/data_format.html)).

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 \_\_\_\_\_ 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.

```
# There are 8 feature types in the GENCODE.
# CDS 567862 / exon 744835 / gene 60603 / Selenocysteine 96 / start_codon 57886 / stop_codon 57775 / tr
```

### 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', ]
d1
```

```
## # A tibble: 60,603 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 \"ENSG00~
## 2 chr1  HAVANA   gene      14404  29570 .    -    .    "gene_id \"ENSG00~
## 3 chr1  ENSEMBL  gene      17369  17436 .    -    .    "gene_id \"ENSG00~
## 4 chr1  HAVANA   gene      29554  31109 .    +    .    "gene_id \"ENSG00~
## 5 chr1  ENSEMBL  gene      30366  30503 .    +    .    "gene_id \"ENSG00~
## 6 chr1  HAVANA   gene      34554  36081 .    -    .    "gene_id \"ENSG00~
## 7 chr1  HAVANA   gene      52473  53312 .    +    .    "gene_id \"ENSG00~
## 8 chr1  HAVANA   gene      57598  64116 .    +    .    "gene_id \"ENSG00~
## 9 chr1  HAVANA   gene      65419  71585 .    +    .    "gene_id \"ENSG00~
## 10 chr1 HAVANA   gene      89295  133723 .   -    .    "gene_id \"ENSG00~
## # ... with 60,593 more rows
```

### 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 ([https://asia.ensembl.org/info/genome/genebuild/manual\\_havana.html](https://asia.ensembl.org/info/genome/genebuild/manual_havana.html))
2. Ensembl: Automatic gene annotation ([https://asia.ensembl.org/info/genome/genebuild/automatic\\_coding.html](https://asia.ensembl.org/info/genome/genebuild/automatic_coding.html))

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 (<https://asia.ensembl.org/info/genome/genebuild/ccds.html>).

Figure 2. Comparison of CCDS and Gencode (<https://twitter.com/ensembl/status/441959722376499200>).

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"; transcr
```

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"; transcr
strsplit(a, 'transcript_support_level\\s+')'
```

```
## [[1]]
## [1] "chr1    HAVANA    transcript    11869    14409    .    +    .    gene_id \"ENSG00000223972.5\""; transcr
## [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"                                "; 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 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.
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_sup
d2$transcript_support_level <- as.character(do.call(rbind.data.frame, strsplit(d2$transcript_support_le
```

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.

### 3. Exercises

Here I list the questions for group 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)

```
library(tidyverse)
library(readr)

d2 = rtracklayer::import('gencode.v31.basic.annotation.gtf.gz') %>% as.data.frame()
```

#### 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?

```
d2 %>% filter(type == 'transcript') %>% count(gene_id) %>% summarise(num = n) %>% pull(num) %>% mean()
```

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

```
d2 %>% filter(type == 'transcript') %>% count(gene_id) %>% summarise(num = n) %>% pull(num) %>% quantile()
```

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

```
d2 %>% filter(type == 'transcript') %>% group_by(gene_id) %>% count() %>% ungroup() %>% top_n(1, n)
```

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

```
d2 %>% filter(type == 'transcript') %>% count(gene_type) %>% summarise(num = n) %>% pull(num) %>% mean()
```

```
## [1] 2706.075
```

```
d2 %>% filter(type == 'transcript') %>% count(gene_type) %>% summarise(num = n) %>% pull(num) %>% quantile()
```

```
##    25%    50%    75%
##    5.75  50.50 891.00
```

```
d2 %>% filter(type == 'transcript') %>% group_by(gene_type) %>% count() %>% ungroup() %>% top_n(1, n)
```

```
## # A tibble: 1 x 2
##   gene_type          n
##   <chr>         <int>
## 1 protein_coding 57846
```

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

```
d2 %>% filter(type == 'transcript') %>% count(seqnames) %>% summarise(num = n) %>% pull(num) %>% mean()
```

```
## [1] 4329.72
```

```
d2 %>% filter(type == 'transcript') %>% count(seqnames) %>% summarise(num = n) %>% pull(num) %>% quantile()
```

```
## 25% 50% 75%
## 2618 4350 5612
```

```
d2 %>% filter(type == 'transcript') %>% group_by(seqnames) %>% count() %>% ungroup() %>% top_n(1, n)
```

```
## # A tibble: 1 x 2
##   seqnames          n
##   <fct>         <int>
## 1 chr1          9827
```

### 3.2. Gene length in the GENCODE

1. What is the average length of human genes?

```
d2 %>% filter(type == 'gene') %>% summarise(mean(width))
```

```
##   mean(width)
## 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.

```
d2 %>% mutate(chrom = ifelse(seqnames %in% c('chrX', 'chrY'), 'sex chromosome', ifelse(seqnames == 'chr1', 'chr1', 'chr2')))
```

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

```
## # A tibble: 10 x 2
## # Groups:   chrom [2]
##   chrom      qnt
##   <chr>    <dbl>
## 1 autosome      8
## 2 autosome    565
## 3 autosome   3779
## 4 autosome  25813
## 5 autosome 2473537
## 6 sex chromosome    48
## 7 sex chromosome   473
## 8 sex chromosome  1912
## 9 sex chromosome 13502
## 10 sex chromosome 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.

```
d2 %>% group_by(gene_type) %>% summarise(qnt = quantile(width, ))
```

## '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      qnt
##   <chr>    <dbl>
## 1 IG_C_gene      3
## 2 IG_C_gene     92
## 3 IG_C_gene   312.
## 4 IG_C_gene   336
## 5 IG_C_gene  8914
## 6 IG_C_pseudogene  34
## 7 IG_C_pseudogene 293
## 8 IG_C_pseudogene 316
## 9 IG_C_pseudogene 424
## 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 %>% filter(type == 'transcript') %>% count(transcript_support_level)
```

```
## transcript_support_level      n
## 1                      1 31801
## 2                      2 13372
## 3                      3  7228
## 4                      4  2245
## 5                      5 13674
## 6                     NA 27843
## 7                     <NA> 12080
```

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

```
d2 %>% filter(type == 'transcript') %>% group_by(gene_type, transcript_support_level) %>% count(transcript)
```

```
## # 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      NA                      7
## 4 IG_C_gene      <NA>                   5
## 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 %>% filter(type == 'transcript') %>% group_by(source, transcript_support_level) %>% count(transcript)
```

```
## # A tibble: 14 x 3
## # Groups:   source, transcript_support_level [14]
##   source transcript_support_level      n
##   <fct>   <chr>                  <int>
## 1 HAVANA  1                29434
## 2 HAVANA  2                12052
## 3 HAVANA  3                 6964
## 4 HAVANA  4                 2116
## 5 HAVANA  5                10157
```

```
## 6 HAVANA NA 19962
## 7 HAVANA <NA> 11901
## 8 ENSEMBL 1 2367
## 9 ENSEMBL 2 1320
## 10 ENSEMBL 3 264
## 11 ENSEMBL 4 129
## 12 ENSEMBL 5 3517
## 13 ENSEMBL NA 7881
## 14 ENSEMBL <NA> 179
```

### 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.

```
df <- d2 %>% filter(type == 'gene') %>% summarise(gene_id = gene_id, gene_name = gene_name, hgnc_id = hgnc_id, ccds_id = ccds_id)
```

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

```
df %>% group_by(gene_type) %>% filter(isHgnc == 1) %>% count()
```

```
## # A tibble: 36 x 2
## # Groups:   gene_type [36]
##   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_V_gene     142
## 7 IG_V_pseudogene 185
## 8 lncRNA       3970
## 9 miRNA       1856
## 10 misc_RNA    1033
## # ... with 26 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).

```
d2 %>% mutate(isHgnc = ifelse(is.na(hgnc_id) == T, 0, 1)) %>% filter(isHgnc == 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?

```
d2 %>% filter(type == 'transcript') %>% group_by(gene_id) %>% count() %>% ungroup() %>% top_n(1, n)
```

```
## # A tibble: 1 x 2
##   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.

```
d2 %>% filter(type == 'gene', gene_type %in% c('lncRNA', 'protein_coding')) %>% group_by(gene_type) %>%
```

```
## '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      qnt
##   <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.

```
d2 %>% 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.

```
d2 %>% filter(type == 'gene') %>% group_by(seqnames) %>% count()
```

```
## # A tibble: 25 x 2
## # Groups:   seqnames [25]
##   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
## # ... with 15 more rows
```

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

```
d2 %>% mutate(chrom = ifelse(seqnames %in% c('chrX', 'chrY'), 'sex chromosome', ifelse(seqnames == 'chr1', 'chr1', 'autosome')))
```

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

```
## # A tibble: 2 x 3
##   chrom          mean median
##   <chr>         <dbl>  <dbl>
## 1 autosome      2617.   2604.
## 2 sex chromosome 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.

```
d2 %>% filter(type == 'gene', gene_type %in% c('lncRNA', 'protein_coding')) %>% group_by(seqnames, gene_type)
```

```
## # A tibble: 49 x 3
## # Groups:   seqnames, gene_type [49]
##   seqnames gene_type     n
##   <fct>     <chr>     <int>
## 1 chr1      lncRNA       1416
## 2 chr1      protein_coding  2048
## 3 chr2      lncRNA       1241
## 4 chr2      protein_coding  1247
## 5 chr3      lncRNA        861
## 6 chr3      protein_coding  1075
## 7 chr4      lncRNA        790
## 8 chr4      protein_coding   751
## 9 chr5      lncRNA        950
## 10 chr5      protein_coding   886
## # ... with 39 more rows
```