

Week 1 - CCLE Annotations

Sidharth Jain

6/12/2020

Introduction

The **Cancer Cell Line Encyclopedia** is a *massive* collection of data that covers more than 1400 cancer cell lines (and still growing!). It contains various types of data, of which only some are listed below:

- RNAseq (gene expression profile)
- RPPA (reverse-phase protein arrays - protein production profile)
- Fusion/translocation events (chromosomal/gene level rearrangements)
- miRNA (micro-RNA expression)
- Mutations (germline - present in all cells, heritable; somatic - not heritable)
- Copy number (chromosome/gene level amplifications and deletions)
- RRBS (reduced-representation bisulfite sequencing - methylation profiles)

The human genome has somewhere between 17,000-30,000 genes (depending on who you ask), so some of these datasets can have around 1,400 x 17,000 different points of data!

How many points of data is that exactly? Let's calculate it! Remember that we can treat R like a calculator:

```
number_of_celllines <- 1400 # assign the value 1400 to the variable "number_of_celllines"
number_of_genes <- 17000 # do the same thing for 17000 to "number_of_genes"

number_of_datapoints <- number_of_genes * number_of_celllines # create a new variable called number_of_
number_of_datapoints # show our result!

## [1] 23800000
```

So a single chunk of data can have approximately 23800000 different values.

That's a LOT of data! But before we even begin to dive into the actual data itself, we need to learn information about the data.

Getting data about the data

To start off, we're going to take a look at information about the cell lines in CCLE. Each of the cell lines listed here are used in the wet lab. For example, you may have heard the story of HeLa cells, which were obtained from a woman named Henrietta Lacks. If you have some free time, then I highly recommend reading *The Immortal Life of Henrietta Lacks*!

Like HeLa cells, each of the cell lines in CCLE was obtained from real patients, isolated from tumor biopsies, and immortalized for research use. Each of has a story to tell about how cancer arose, how it grew, and hopefully, how we can stop it from growing.

The first thing we want to do is find the file that contains the meta-data we want for the Cancer Cell Line Encyclopedia. I have placed that file for you in your working directory (`Cell_lines_annotations_20181226.txt`) that contains the data we want to look at. To start, let's load this file and see what kind of data we're working with.

```
CCLE_metadata <- read.delim("Cell_lines_annotations_20181226.txt") # use the read.delim function to read

nrow(CCLE_metadata) # how many rows does our data have?
ncol(CCLE_metadata) # how many columns does our data have?

CCLE_metadata[1:5,1:10] # what does our data actually look like?
# Here, I'm getting the first 5 rows (1:5), and the first 10 columns (1:10), just to see - try increasing
```

CCLE_ID	depMapID	Name	Pathology	Site_Primary
DMS53_LUNG	ACH-000698	DMS 53	primary	lung
SW1116_LARGE_INTESTINE	ACH-000489	SW1116	primary	large_intestine
NCIH1694_LUNG	ACH-000431	NCI-H1694	metastasis	lung
P3HR1_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	ACH-000707	P3HR-1	metastasis	haematopoietic
HUT78_HAEMATOPOIETIC_AND_LYMPHOID_TISSUE	ACH-000509	HuT 78	primary	haematopoietic

So according to this file, there are 1,461 rows! Each row should correspond to a cell line, and we can see that the first 3 columns show 3 different naming or ID conventions for each cell line. The first is the CCLE ID, which just takes the name of the cell line and stitches it to the type of tissue it came from. For example, A549_LUNG is a lung cancer cell line that is conventionally called A549.

For convenience, the **type** is also present in it's own column, and may have additional details in the **type_refined** column.

(P.S., you can also find HELA_CERVIX, and the data associated with the HeLa cancer cell line)

There's a lot of other very interesting data here. Let's take a look at some of the column names to get a sense of what else we can find out about these cell lines.

```
colnames(CCLE_metadata)
```

Column names of CCLE metadata

```
CCLE_ID
depMapID
Name
Pathology
Site_Primary
Site_Subtype1
Site_Subtype2
Site_Subtype3
Histology
Hist_Subtype1
Hist_Subtype2
Hist_Subtype3
Gender
Life_Stage
Age
Race
Geo_Loc
inferred_ethnicity
Site_Of_Finding
Disease
Annotation_Source
Original.Source.of.Cell.Line
```

Column names of CCLE metadata

Characteristics
Growth.Medium
Supplements
Freezing.Medium
Doubling.Time.from.Vendor
Doubling.Time.Calculated.hrs
type
type_refined
PATHOLOGIST_ANNOTATION
mutRate
tcga_code

As background, the `Pathology` (and `site_subtypes`) and `Histology` (and `Hist_subtypes`) are annotations provided by the pathologist after they took the patient's tumor out. We also have information about each patient, including:

- `Gender`
- `Age`
- `Race`
- `Geo_Loc` (where was the patient sample obtained?)
- `inferred_ethnicity` (based on genetic analysis of the cell line)

There's also information about how the cells were grown in the lab. This includes:

- `Original.Source.of.Cell.Line` (where did the CCLE get their stock of cell line from?)
- `Characteristics` (is it an adherent cell line that sticks to the plate, or does it grow suspended in fluid like a blood cell?)
- `Growth.Medium` (what media does the cell grow in?)
- `Supplements` (what additional nutrients does the cell need to grow?)
- `Freezing.Medium` (what media was the cell stored in?)
- `Doubling.Time.From.Vendor` (according to the source, how long does it take for the cells to replicate?)
- `Doubling.Time.Calculated.Hours` (how many hours did the cells actually take to double?)

Diving into the metadata

Now that we have a pretty good understanding of what kinds of data we have, let's take a look at the data itself.

The first thing we should do is understand how some of the values look. R has some useful functions for summarizing data, some of which I'm going to demonstrate here:

The `table` function gives us a frequency table for a categorical variable (how many times does each value occur in the data?)

```
# First, let's look at the different tissue/tumor types we're working with.
table(CCLE_metadata$type)
```

type	freq
AML	37
B-cell_ALL	12
bile_duct	8
breast	60
chondrosarcoma	4

type	freq
CML	15
colorectal	63
endometrium	28
esophagus	26
Ewings_Sarcoma	12
giant_cell_tumour	3
glioma	65
kidney	37
leukemia_other	5
liver	28
lung_NSC	135
lung_small_cell	53
lymphoma_Burkitt	11
lymphoma_DLBCL	18
lymphoma_Hodgkin	13
lymphoma_other	28
medulloblastoma	4
melanoma	63
meningioma	3
mesothelioma	11
multiple_myeloma	29
neuroblastoma	17
osteosarcoma	10
other	4
ovary	55
pancreas	46
prostate	8
soft_tissue	20
stomach	39
T-cell_ALL	16
thyroid	12
upper_aerodigestive	33
urinary_tract	28

HOT TIP: Remember that we loaded the CCLE metadata as a data frame. This means we can access a column of the data frame using the `$` operator - so for accessing the `type` column, I wrote `CCLE_metadata$type`. We can also access columns using the `[]`, so I could have also written `CCLE_metadata[, "type"]`. I could also access a specific column number, by writing `CCLE_metadata[, 5]`.

We can also use the `table` function to compare two variables - so if I wanted to look at both gender and ethnicity, I could do the following:

```
table(CCLE_metadata$Gender, CCLE_metadata$inferred_ethnicity)
```

```
##
##      African_american Asian Caucasian
##           7      56      73
## female      29     102     268
## male       25     148     329
## null        0       0       0
```

We'll definitely talk more about this on Friday, but for now, I hope you see how useful the `table` function can be!

Another useful function is `summary`, which can give us important summary statistics about numeric variables. One interesting numerical variable here is mutation rate (or `mutRate`). See below for the output of `summary`

```
summary(CCLE_metadata$mutRate)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
##    53.25  102.92  134.96  192.92  178.86 3119.62     497
```

This means that, on average, the mutation rate for each cell line is 192.92. Don't worry about exactly what that means just yet, we'll cover it on Friday!

One important thing to see is the number of NA values. These are values that are missing from the data, and are handled specifically by R as missing values. When we calculate the mean and median by themselves in R, we need to make sure we omit NA values - it's like multiplying by 0, all the other numbers end up becoming NA also. Let's try just taking the mean without omitting NA, using the `mean` function.

```
mean(CCLE_metadata$mutRate) # welp, that didn't work
```

```
## [1] NA
```

```
mean(CCLE_metadata$mutRate, na.rm = T) # na.rm means "remove NAs for this calculation" - much better!
```

```
## [1] 192.9194
```

We can also use other summary statistics - for example, `sd` gives standard deviation, `var` gives variance, and there are others that we'll learn later on.

HOT TIP: If you don't already know, you can pull up the help page for any function by putting `?` before the function and running it. For example, `?mean` pulls up the help page for the mean function. You can read here what other arguments to functions are (like `na.rm`) and what their default values are (`na.rm` by default is `FALSE`, so NA values are NOT removed).

Narrowing down to one type of cancer

Just as an example, let's say we wanted to learn more about breast cancer. We would need to take a subset of the data. So let's do that here:

```
is_breast_cell <- CCLE_metadata$type == "breast" # if the type is breast, give me TRUE. otherwise, give
which_is_breast_cell <- which(is_breast_cell) # "which" is a function that gives the index of each TRUE
```

```
which_is_breast_cell
```

```
## [1] 12 55 65 66 71 81 91 92 101 173 225 259 260 278 292
## [16] 303 304 310 311 320 327 330 331 336 356 548 551 595 599 613
## [31] 630 645 680 684 728 756 757 761 764 766 812 860 880 892 897
## [46] 898 899 911 941 954 955 956 959 964 966 984 1008 1028 1035 1044
```

If that wasn't clear, what we want to do is only identify which cell lines have a `type` that is exactly "breast", so we use the `==` comparison to check if each value in the `type` column is breast or not. If it is, we get a `TRUE`, and if it's not, we get a `FALSE`. That gives us a single logical vector containing `TRUE`s and `FALSE`s for each value in `type`.

Then, we use `which` to give us the index (or position) of each `TRUE`. From this, we know that the 12th value is a `TRUE`, which means that row 12 in `CCLE_metadata` data frame is a breast cancer cell line. We can check that here:

```
CCLE_metadata[12,] # give us the 12th row
```

	CCLE_ID	depMapID	Name	Pathology	Site_Primary	Site_Subtype1	Site_Subtype2	Site
12	HCC2157_BREAST	ACH-000691	HCC2157	primary	breast	NS	NS	NS

Yep! The 12th cell line is HCC2157_BREAST, which is indeed a breast cancer cell line.

Okay, so now that we have all of the indices (plural of index, not indexes!) of the breast cell lines, I'm curious - which breast cancer cell line has the highest mutation rate?

```
CCLE_metadata_breast <- CCLE_metadata[which_is_breast_cell,] # make a new subsetting dataframe that only
max(CCLE_metadata_breast$mutRate, na.rm = T) # The highest mutation rate - don't forget to remove NAs!

## [1] 486.8665

which_breast_highest_mutRate <- which.max(CCLE_metadata_breast$mutRate) # Get the index using the which
CCLE_metadata_breast$CCLE_ID[which_breast_highest_mutRate]

## [1] "HCC1569_BREAST"
```

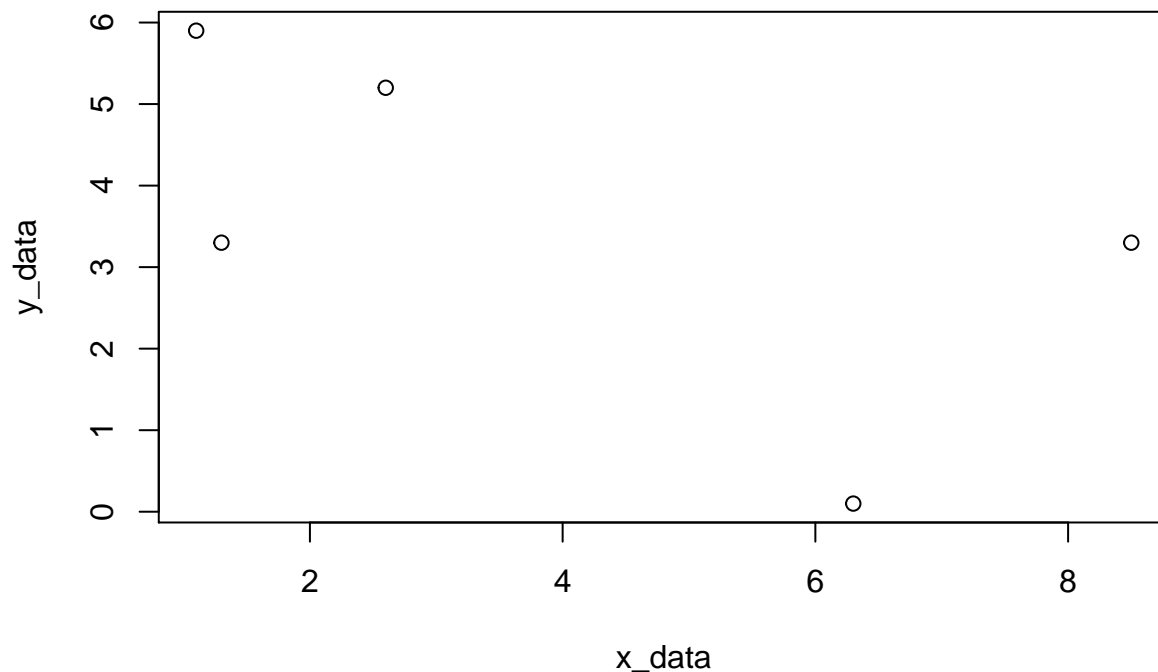
QUESTIONS

Create a new R chunk or click the “Insert -> R” near the top right of the script-writing panel, and write code to solve the following questions.

1. Now that we know how many cell lines and columns there are, can you calculate the number of data points in the metadata file? Hint: Remember that R is a calculator!
2. A software was used to infer the ethnicity of the cell line from genomic data. Which ethnicities are most/least represented in CCLE? Compare that to the reported race of the patient. Can you think about what the consequences of this might be? What about the software that is used to infer ethnicity? Hint: Use the `table` function to answer this question.
3. Cancer is often thought of as a disease of aging, but in truth, it affects people of all ages. What are some of the summary statistics about the ages of the patients from whom the CCLE was derived? BONUS: Can you identify which cancers predominantly affect people under age 18? Hint: Use the `summary` function to answer this question.
4. Growing cells in a dish requires that you understand how frequently the cell population doubles. Cells that grow very quickly need to be “passaged”, or transferred from one plate to another to allow the cells more room to grow. The “doubling time” has been recorded for all CCLE cell lines, and we can use it as a proxy for how fast the tumor might have been growing in the patient. What is the average doubling time (use calculated hours) for colorectal cancer cell lines? Hint: You'll need to subset for colorectal cancer cell lines.
5. Create a plot showing the relationship between doubling time and mutation rate. BONUS: Are these variables correlated? Extra bonus for anyone who can use the `ggplot2` package (load using `library(ggplot2)`) to create their plot, and/or who can color their plot by cancer type. Hint: The base plot function works by taking in data for the x-axis and y-axis as follows:

```
x_data <- c(1.3, 2.6, 6.3, 8.5, 1.1)
y_data <- c(3.3, 5.2, 0.1, 3.3, 5.9)

plot(x_data, y_data)
```



Extra

hint: The `qplot` function of `ggplot2` also works similarly! BONUS hint: To learn more about the correlation function in R, try typing `?cor` to read the help page.

Answer key

1. See below code:

```
number_of_metadatapoints <- nrow(CCLE_metadata) * ncol(CCLE_metadata)
```

2. This is one possible solution:

```
table(CCLE_metadata$inferred_ethnicity)
```

```
##
## African_american      Asian      Caucasian
##           61           306           670
```

```
table(CCLE_metadata$Race)
```

```
##
##
##           african african_american american_indian
##           469           2           33           1
##           asian      caucasian      east_indian      north_african
##           181           359           1           1
##           turkish
##           1
```

Let's make a table to compare stated race with inferred ethnicity

```
table(CCLE_metadata$inferred_ethnicity, CCLE_metadata$Race)
```

```
##
##           african african_american american_indian asian caucasian
## African_american  20           2           31           0           3           5
## Asian            130           0           0           0          173           2
## Caucasian         302           0           2           1           4          350
```

```
##
##          east_indian north_african turkish
## African_american      0           0      0
## Asian                  0           0      0
## Caucasian              1           1      1
```

This teaches us about the importance of building the right reference, and creating the right representation - why does the algorithm only consider African American, Asian, and Caucasian ethnicities? And why are there not other races represented in CCLE?

3. This is one possible solution:

```
summary(CCLE_metadata$Age)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.     NA's
##      0.25  39.00   54.00   49.18  64.00   92.00     662
```

```
# BONUS: This was a hard one, I'll be honest. Don't sweat if you didn't get this!
```

```
is_under_18 <- CCLE_metadata$Age < 18
```

```
table(CCLE_metadata$type, is_under_18) # we can get our answer visually, by looking at this table
```

```
##          is_under_18
##          FALSE TRUE
## AML                26   6
## B-cell_ALL         4   7
## bile_duct          0   0
## breast             58   0
## chondrosarcoma      3   0
## CML                 13   2
## colorectal         47   0
## endometrium        20   0
## esophagus          17   0
## Ewings_Sarcoma      4   5
## giant_cell_tumour    1   2
## glioma             35   2
## kidney             11   0
## leukemia_other      2   3
## liver              23   3
## lung_NSC            94   1
## lung_small_cell     48   0
## lymphoma_Burkitt     1   8
## lymphoma_DLBCL      12   1
## lymphoma_Hodgkin     10   1
## lymphoma_other      21   5
## medulloblastoma      0   4
## melanoma           44   1
## meningioma          0   0
## mesothelioma         9   0
## multiple_myeloma     24   0
## neuroblastoma        0  13
## osteosarcoma         2   6
## other                4   0
## ovary               35   0
## pancreas            37   0
## prostate            6   0
## soft_tissue          7  12
```



```
##      stomach          29      0
##      T-cell_ALL       5      11
##      thyroid          11      0
##      upper_aerodigestive 24      1
##      urinary_tract    18      0
```

```
# or programmatically by saving the table as a variable, and finding out which rows have more TRUE than
under_18_table <- table(CCLE_metadata$type, is_under_18)
which(under_18_table[,1] < under_18_table[,2]) # These are cancer cell lines where more samples were de
```

```
##      B-cell_ALL      Ewings_Sarcoma giant_cell_tumour      leukemia_other
##              2              10              11              14
## lymphoma_Burkitt      medulloblastoma      neuroblastoma      osteosarcoma
##              18              22              27              28
##      soft_tissue      T-cell_ALL
##              33              35
```

4. This is one possible solution:

```
is_colorectal_cell <- CCLE_metadata$type == "colorectal" # if the type is colorectal, give me TRUE. oth
which_is_colorectal_cell <- which(is_colorectal_cell) # "which" is a function that gives the index of e

# Now let's subset for only the colorectal cancer cells (which we identified above)
CCLE_metadata_colorectal <- CCLE_metadata[which_is_colorectal_cell,]

# Get the average doubling time of CRC cells
mean(CCLE_metadata_colorectal$Doubling.Time.Calculated.hrs, na.rm = T)
```

```
## [1] 83.4
```

```
# Identify the fastest growing CRC cell
which_fastest_growing_cell <- which.min(CCLE_metadata_colorectal$Doubling.Time.Calculated.hrs)
CCLE_metadata_colorectal$CCLE_ID[which_fastest_growing_cell]
```

```
## [1] "SW620_LARGE_INTESTINE"
```

```
# EDIT: Identify the fastest growing overall cell (because my question was vague and bad).
which_fastest_growing_colo_cell <- which.min(CCLE_metadata_colorectal$Doubling.Time.Calculated.hrs)
CCLE_metadata_colorectal$CCLE_ID[which_fastest_growing_cell]
```

```
## [1] "SW620_LARGE_INTESTINE"
```

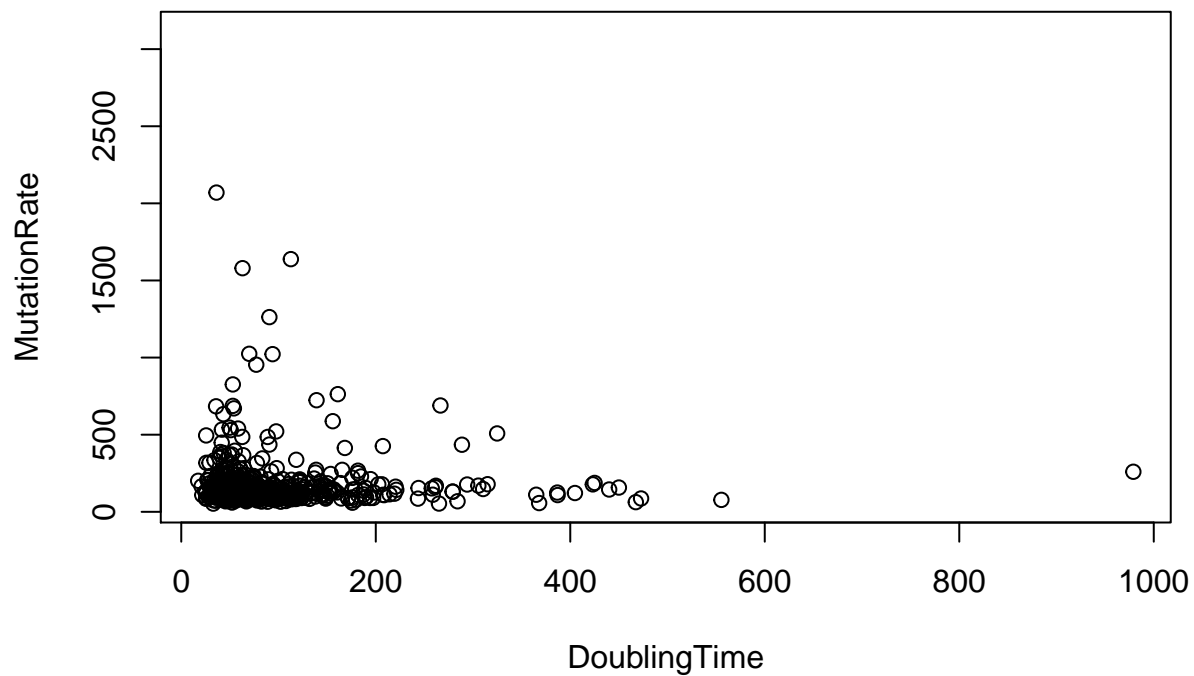
```
which_fastest_growing_cell <- which.min(CCLE_metadata$Doubling.Time.Calculated.hrs)
CCLE_metadata$CCLE_ID[which_fastest_growing_cell]
```

```
## [1] "KYSE410_OESOPHAGUS"
```

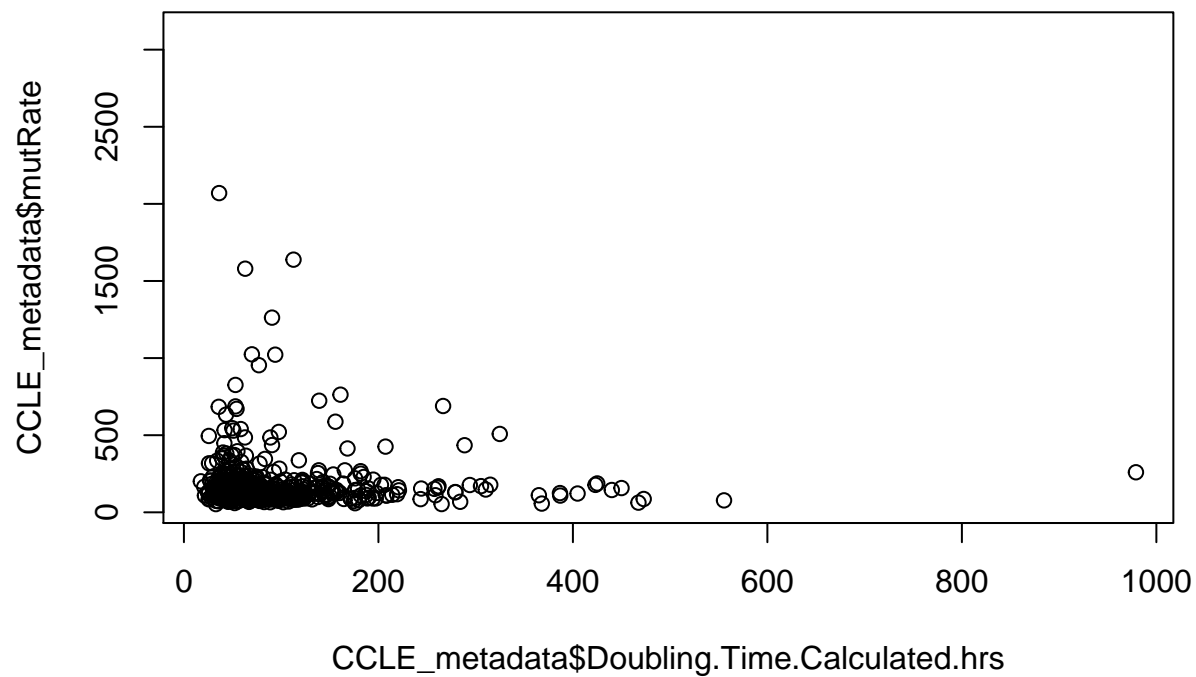
5. This is one possible solution:

```
DoublingTime <- CCLE_metadata$Doubling.Time.Calculated.hrs
MutationRate <- CCLE_metadata$mutRate

plot(DoublingTime, MutationRate)
```

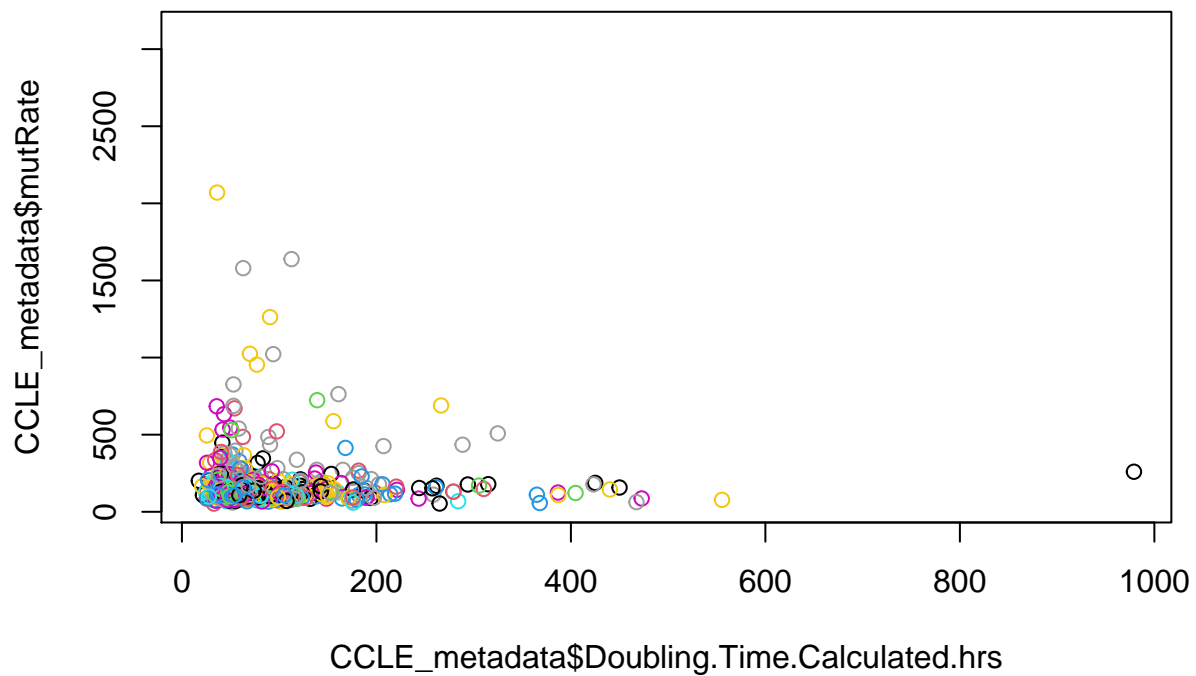


```
plot(CCLE_metadata$Doubling.Time.Calculated.hrs, CCLE_metadata$mutRate)
```



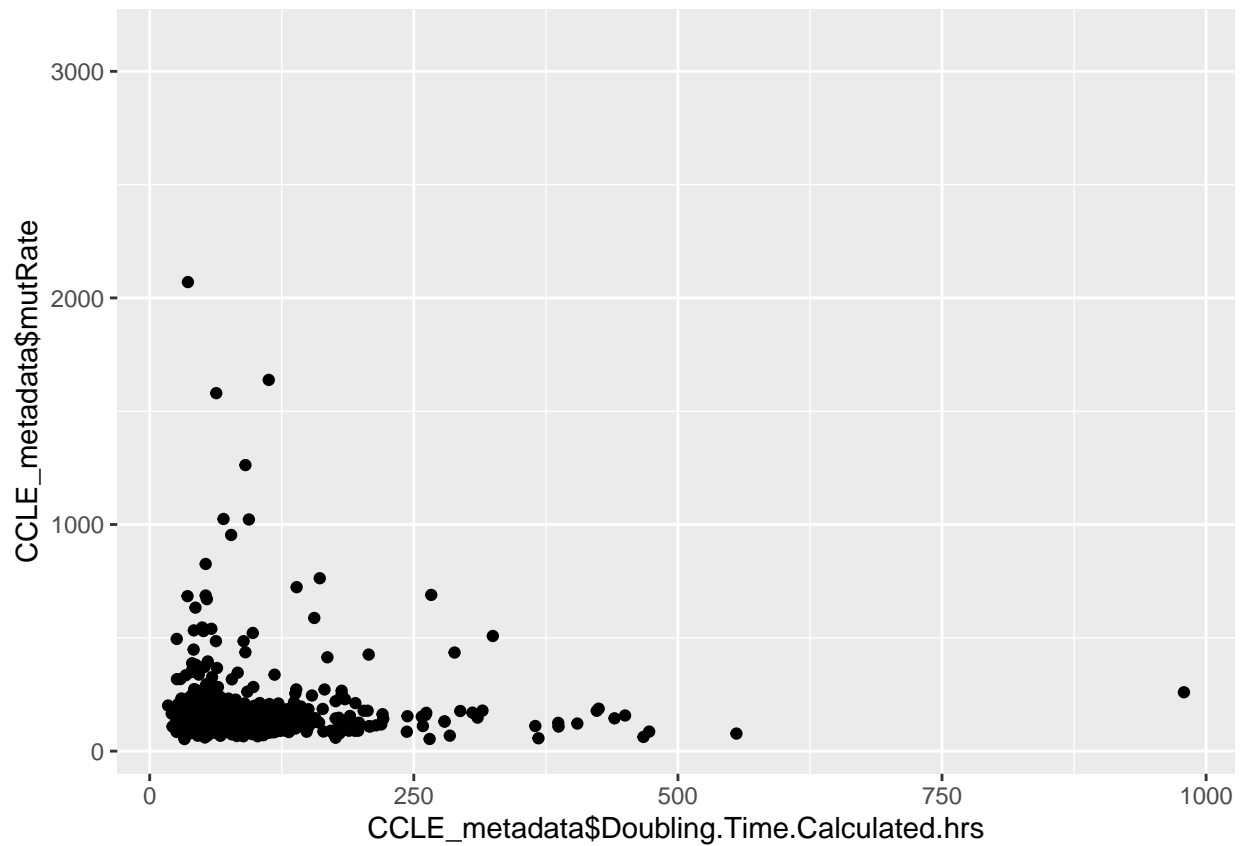
```
# We can also add colors, but this requires us transforming one of the variable types to a factor...
plot(CCLE_metadata$Doubling.Time.Calculated.hrs, CCLE_metadata$mutRate, col = as.factor(CCLE_metadata$...))

# Using qplot from ggplot2 - qplot is like the base plot function -- easy to use, but doesn't give us a
library(ggplot2)
```



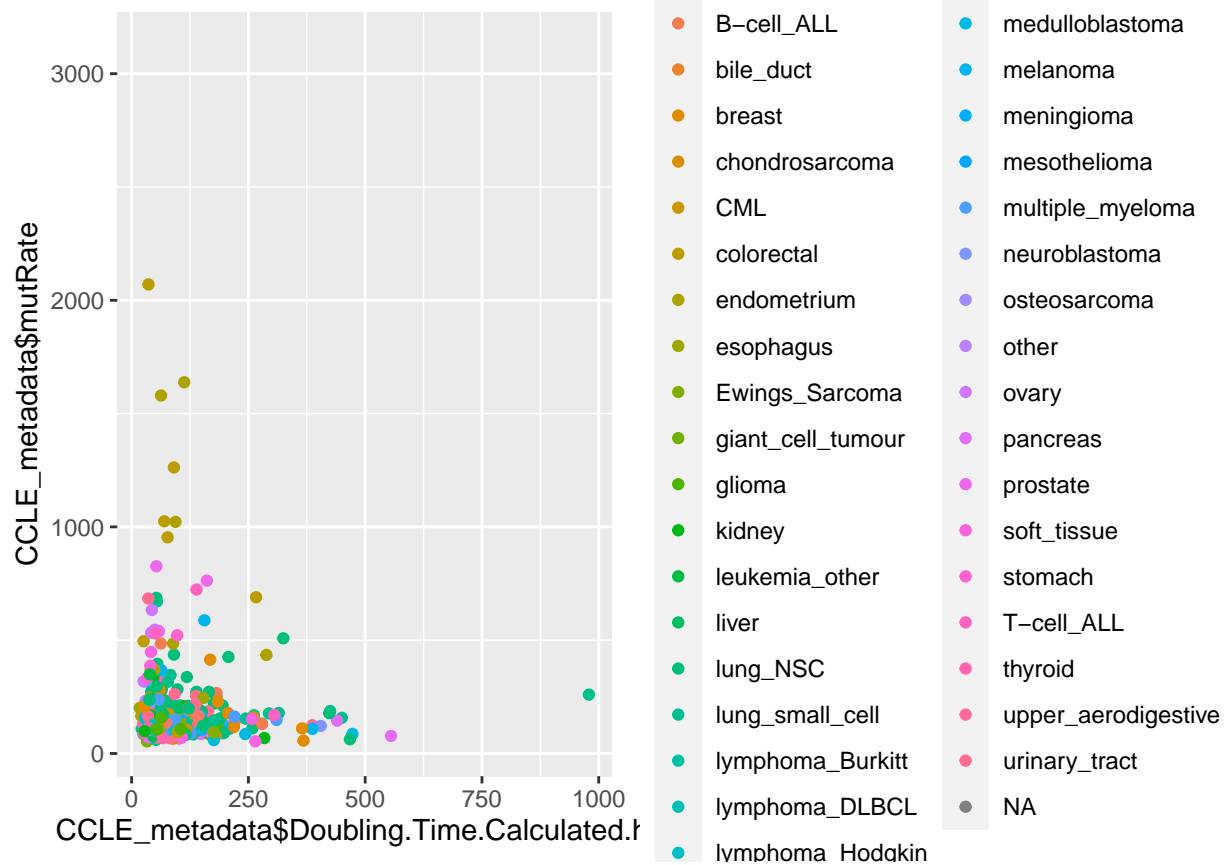
```
qplot(CCLE_metadata$Doubling.Time.Calculated.hrs, CCLE_metadata$mutRate)
```

```
## Warning: Removed 941 rows containing missing values (geom_point).
```



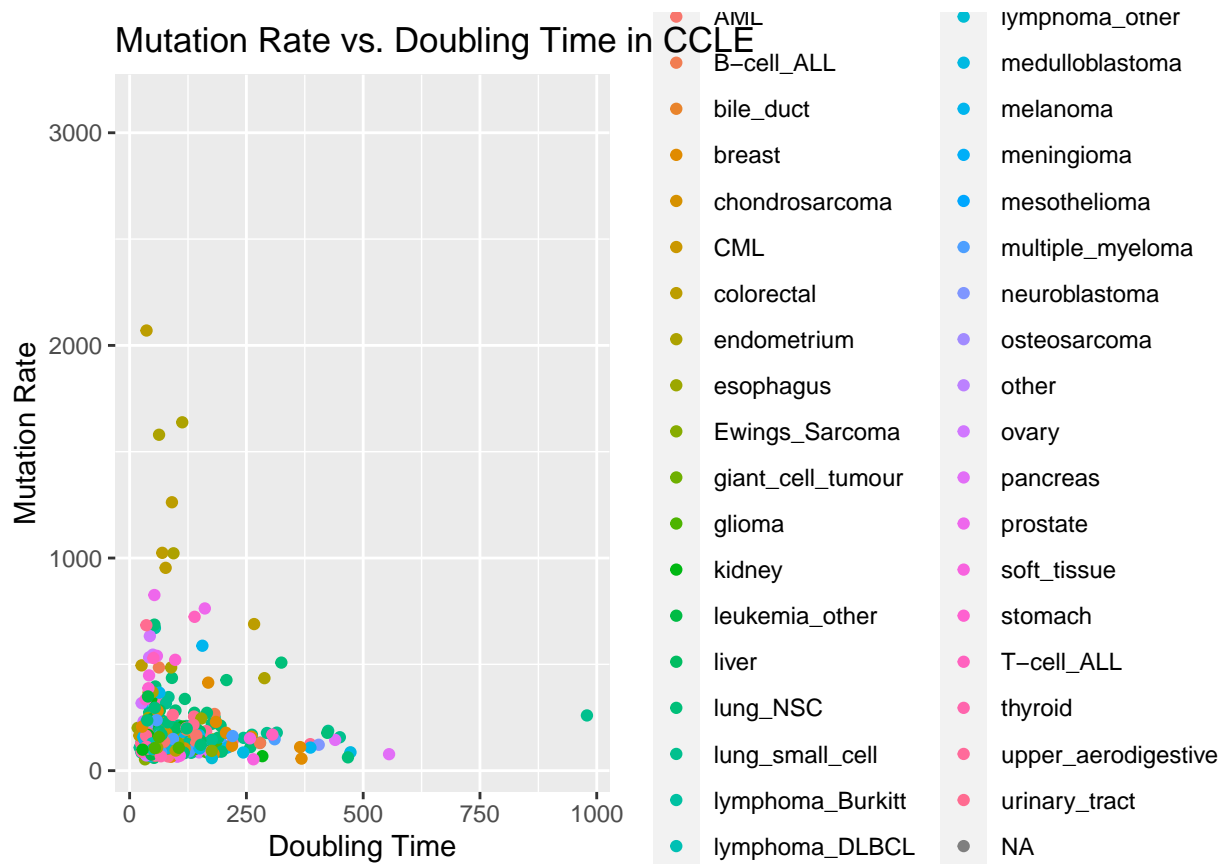
```
qplot(CCLE_metadata$Doubling.Time.Calculated.hrs, CCLE_metadata$mutRate, col = CCLE_metadata$type, geom
```

```
## Warning: Removed 941 rows containing missing values (geom_point).
```



```
# If we wanted to use full-blown ggplot2 (we'll learn this next week, don't worry)
ggplot(data = CCLE_metadata,
  aes(x = Doubling.Time.Calculated.hrs,
    y = mutRate,
    color = type)) +
  geom_point() +
  labs(x = "Doubling Time", y = "Mutation Rate", title = "Mutation Rate vs. Doubling Time in CCLE")
```

```
## Warning: Removed 941 rows containing missing values (geom_point).
```



Is there a correlation?

```
cor(CCLE_metadata$Doubling.Time.Calculated.hrs, CCLE_metadata$mutRate, use = 'complete.obs') # Note: we
```

```
## [1] -0.03234549
```