

A summary of the tidyverse

Jaime Davila

2/2/2022

Transforming and visualizing with tidyverse

tidyverse is a powerful collection of functions and libraries that allows us interact and wrangle datasets in an effective manner. Before we use any of those commands we need to instruct R to the **tidyverse** library which can be done with the following command:

```
library(tidyverse)
library(ggplot2)
```

During this session we will be analyzing datasets from Defining the '90s Music Cannon so we will start by using `read_csv` to load our dataset into the tibble `song.year.tbl`.

```
file.path = "~/Mscs 341 S22/Class/Data/song.year.csv"
song.year.tbl <- read_csv(file.path)
song.year.tbl
```

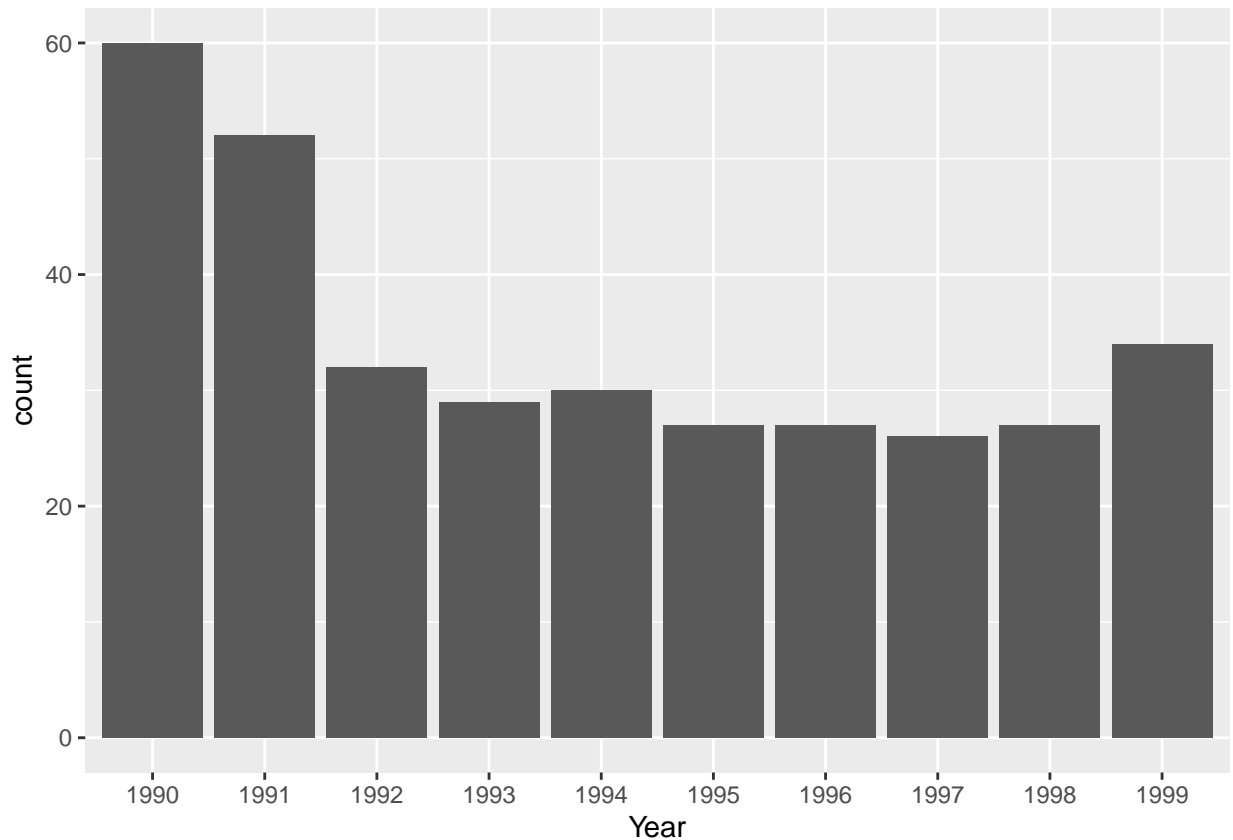
```
## # A tibble: 344 x 3
##   artist      song      year
##   <chr>      <chr>    <dbl>
## 1 2 Pac      California Love  1996
## 2 2Pac      How Do U Want It  1996
## 3 702      Where My Girls At?  1999
## 4 Ace Of Base All That She Wants  1993
## 5 Ace Of Base Don't Turn Around  1994
## 6 Ace Of Base The Sign  1994
## 7 Adina Howard Freak Like Me  1995
## 8 Aerosmith I Don't Want To Miss A Thing  1998
## 9 Aerosmith Janie's Got A Gun  1990
## 10 Alanis Morissette Ironic  1996
## # ... with 334 more rows
```

A first plot

We will start by doing a simple plot summarizing the number of songs per year of our dataset by using `ggplot`. A good presentation of `ggplot` can be found in Chapter 3: Data Visualization. A couple of quick things to note about this code:

- We will be using `geom_histogram` as our geometric object
- Notice that since `year` is a number we need to convert it to a categorical variable using the command `as.factor`

```
ggplot(song.year.tbl, aes(x=as.factor(year)))+
  geom_histogram(stat="count")+
  labs(x="Year")
```



A first summary

We can obtain the number of songs by combining the commands `group_by` and `summarize`. Again notice that:

- We are using the pipe (`%>%`) to combine the two commands together.
- `group_by` can take as an argument any combination of the columns from the tibble.
- `summarize` is a very flexible command and allows the calculation of any number of statistics like average, minimum or maximum. In this particular case we are just counting the number of elements by using the function `n()`.

```
song.year.tbl %>%
  group_by(year) %>%
  summarize(n=n())
```

```
## # A tibble: 10 x 2
##   year      n
##   <dbl> <int>
## 1 1990     60
## 2 1991     52
## 3 1992     32
## 4 1993     29
## 5 1994     30
```

```
## 6 1995 27
## 7 1996 27
## 8 1997 26
## 9 1998 27
## 10 1999 34
```

Transforming datasets

Before attempting the following exercises we recommend that you read [Chapter 5: Data Transformation] (<https://r4ds.had.co.nz/transform.html>).

In particular we will be using the following 7 commands (also called verbs) from **tidyverse**

- `group_by`
- `summarize`
- `slice`
- `arrange`
- `select`
- `filter`
- `mutate`

1. In the following exercises we will modify our original dataset to answer specific questions:
 - a. Generate a table with top-5 artists from the 90s according to their number of songs and call it `top5.tbl`.

```
top5.tbl <- song.year.tbl%>%
  group_by(artist)%>%
  summarize(count = n())%>%
  slice_max(count, n = 5)
```

- b. Let's explore Mariah Carey's career in depth. Start by summarizing the number of hits of Mariah Carey

```
song.year.tbl%>%
  filter(artist == "Mariah Carey")%>%
  group_by(year)%>%
  summarize(count = n())%>%
  arrange(desc(count))
```

```
## # A tibble: 10 x 2
##   year count
##   <dbl> <int>
## 1 1991     3
## 2 1990     2
## 3 1992     2
## 4 1993     2
## 5 1994     2
## 6 1995     2
## 7 1999     2
## 8 1996     1
## 9 1997     1
## 10 1998     1
```

```
#1991 was her best year
```

1991 was her best year with 3 hits

- c. What are the songs from her best year? Do you recognize any of them?

```
song.year.tbl%>%
  filter(artist == "Mariah Carey", year == 1991)
```

```
## # A tibble: 3 x 3
##   artist      song      year
##   <chr>      <chr>    <dbl>
## 1 Mariah Carey Emotions    1991
## 2 Mariah Carey I Don't Wanna Cry 1991
## 3 Mariah Carey Someday    1991
```

Here 3 songs are Emotions, 'I Don't Wanna Cry', and 'Someday'. I recognize Someday.

2. It seems like some artists like Mariah Carey had a song in the billboard in every year of the 90s decade, while others only had one hit in the entire decade.

- a. We are interested in calculating the `year.span` of an artist, which is basically defined as the difference in years between their latest and earliest song in the 90s. Create a table `artist.span.tbl` with such information and include `earliest.year` and `latest.year` as columns:

```
artist.span.tbl <- song.year.tbl%>%
  group_by(artist)%>%
  summarize(earliest.year = min(year),
            latest.year = max(year))%>%
  mutate(year.span = latest.year - earliest.year+1)
```

- b. What are the top-5 artists with biggest span?

```
artist.span.tbl%>%
  arrange(desc(year.span))%>%
  slice(1:5)
```

```
## # A tibble: 5 x 4
##   artist      earliest.year latest.year year.span
##   <chr>          <dbl>      <dbl>    <dbl>
## 1 Mariah Carey    1990        1999      10
## 2 Whitney Houston 1990        1999      10
## 3 Aerosmith      1990        1998       9
## 4 Madonna        1990        1998       9
## 5 Celine Dion     1991        1998       8
```

```
artist.span.tbl%>%
  slice_max(year.span, n = 5)
```

```
## # A tibble: 6 x 4
##   artist      earliest.year latest.year year.span
##   <chr>          <dbl>      <dbl>    <dbl>
## 1 Mariah Carey    1990        1999      10
## 2 Whitney Houston 1990        1999      10
## 3 Aerosmith      1990        1998       9
## 4 Madonna        1990        1998       9
## 5 Celine Dion     1991        1998       8
## 6 TLC            1992        1999       8
```

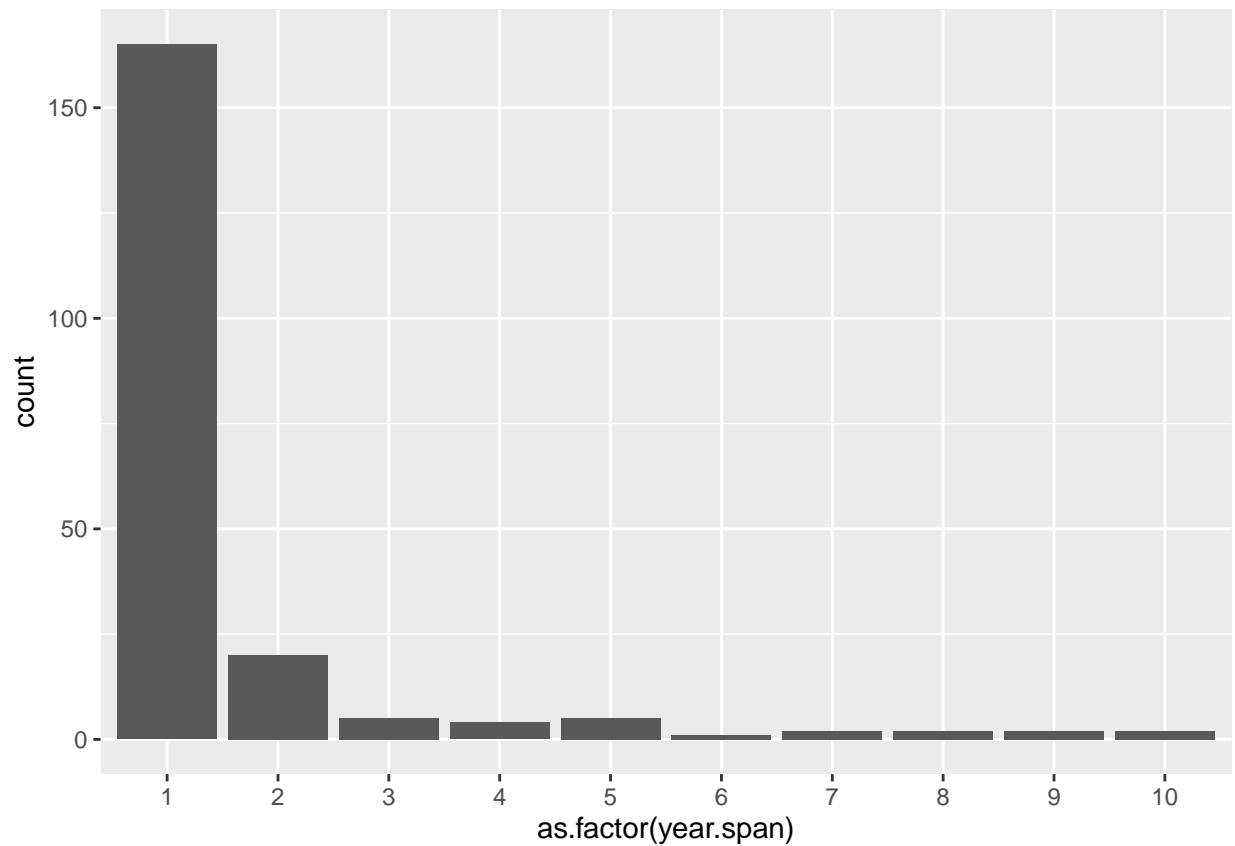
Mariah Carey, Whitney Houston, Aerosmith, Madonna, Celine Dion, TLC (tied with Celine Dion)

- c. Generate a histogram describing the span across all the artists in our table:

```
artist.span.tbl%>%
  ggplot(aes(x = as.factor(year.span)))+
```

```
geom_histogram(stat = "count")
```

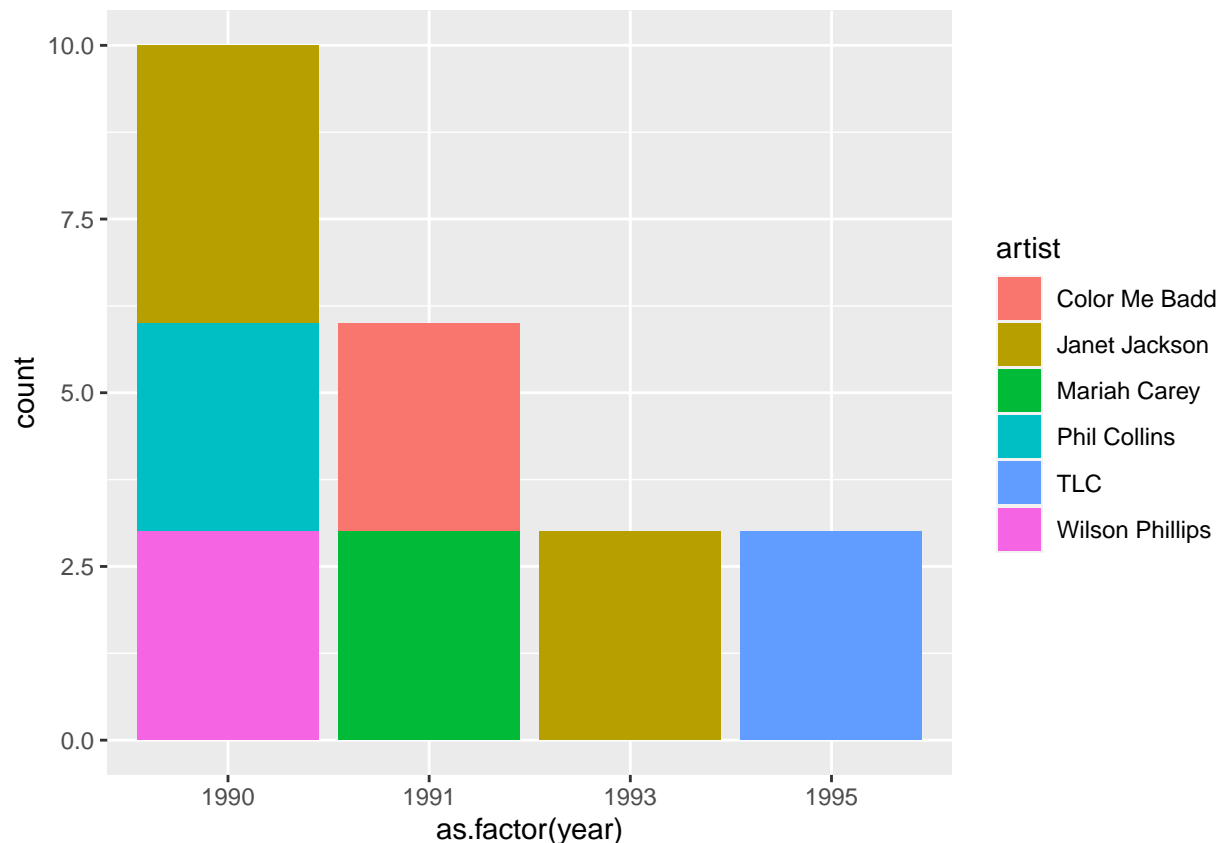
```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```



3. Create a table with year and number of songs by artist and year. Only select artists that had at least 3 hits in a year. Generate a graph depicting this info.

```
song.year.tbl%>%  
  group_by(artist, year)%>%  
  summarize(count = n())%>%  
  filter(count >= 3)%>%  
  ggplot(aes(x = as.factor(year)))+  
    geom_bar(aes(y = count, fill = artist), stat = "identity")
```

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



Joinining datasets

A common situation in analysis is that all of the information is not contained in just a single dataset. Let's start by loading a different dataset which has recognition metrics for all of our previous songs

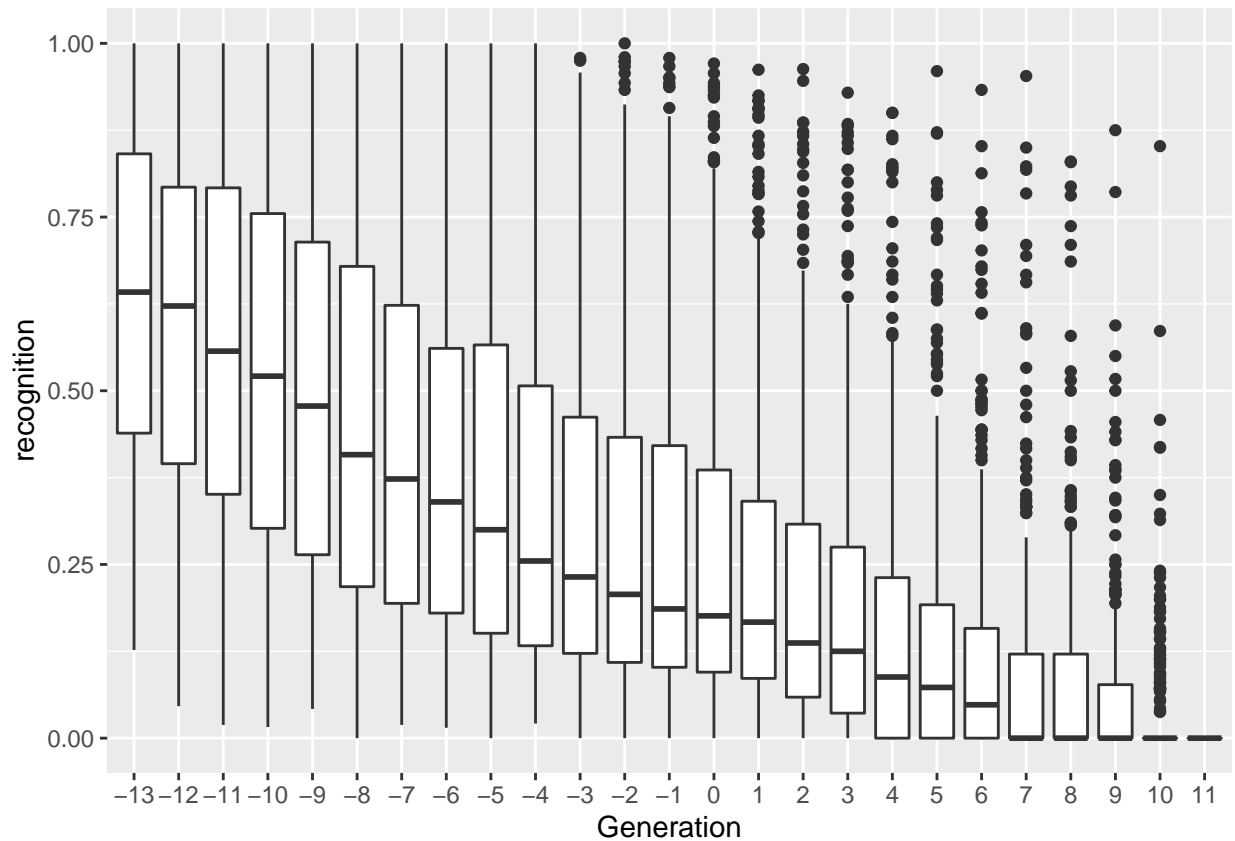
```
file.path = "~/Mscs 341 S22/Class/Data/song.recognition.csv"
```

```
song.recognition.tbl <- read_csv(file.path)
```

Notice how `song.recognition.tbl` has `artist`, `song`, `generation` and `recognition` as variables. `generation` is measured as the age of the respondents when the song was released. For example, Macarena was released in 1996, so -10 represents the people who *were 10 years old* when it was released (and were born in 1986), while 5 represents the people who *were born 5 years after the song* (and were born in 2001).

Let's take a look at the entire dataset using a boxplot

```
ggplot(song.recognition.tbl, aes(as.factor(generation), recognition)) +
  geom_boxplot() +
  labs(x = "Generation")
```



In the following exercises we will explore this dataset in more detail. Make sure to consult Chapter 13: Relational data and pay close attention to the function `inner_join()`

4. Let's explore the decaying trend in recognition is the same according to the year that the song came out. Let's do that using the following steps:

1. Combine `song.recognition.tbl` and `song.year.tbl` in a single table called `song.decay.tbl`

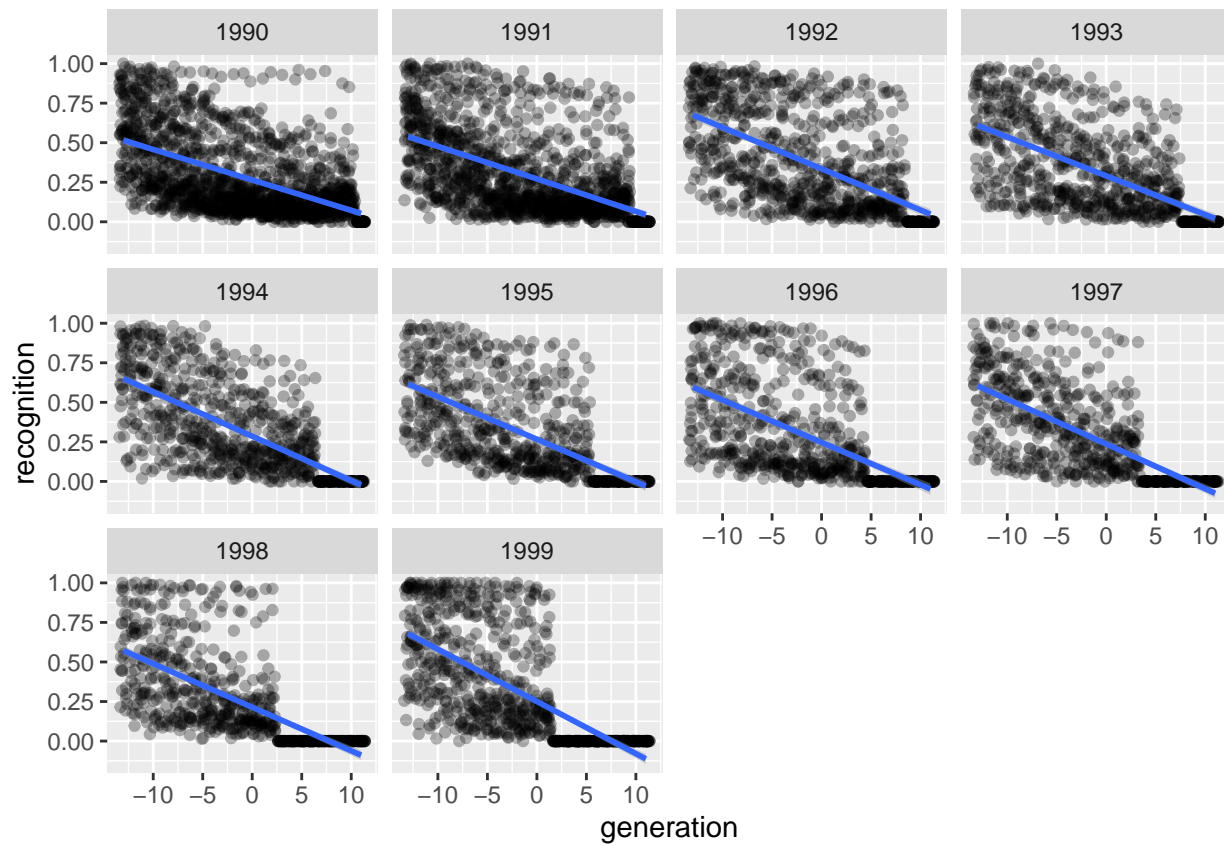
```
song.decay.tbl <- song.recognition.tbl%>%
  inner_join(song.year.tbl)
```

```
## Joining, by = c("artist", "song")
```

2. Create the following which shows the trends for each different year. As the years get closer to the

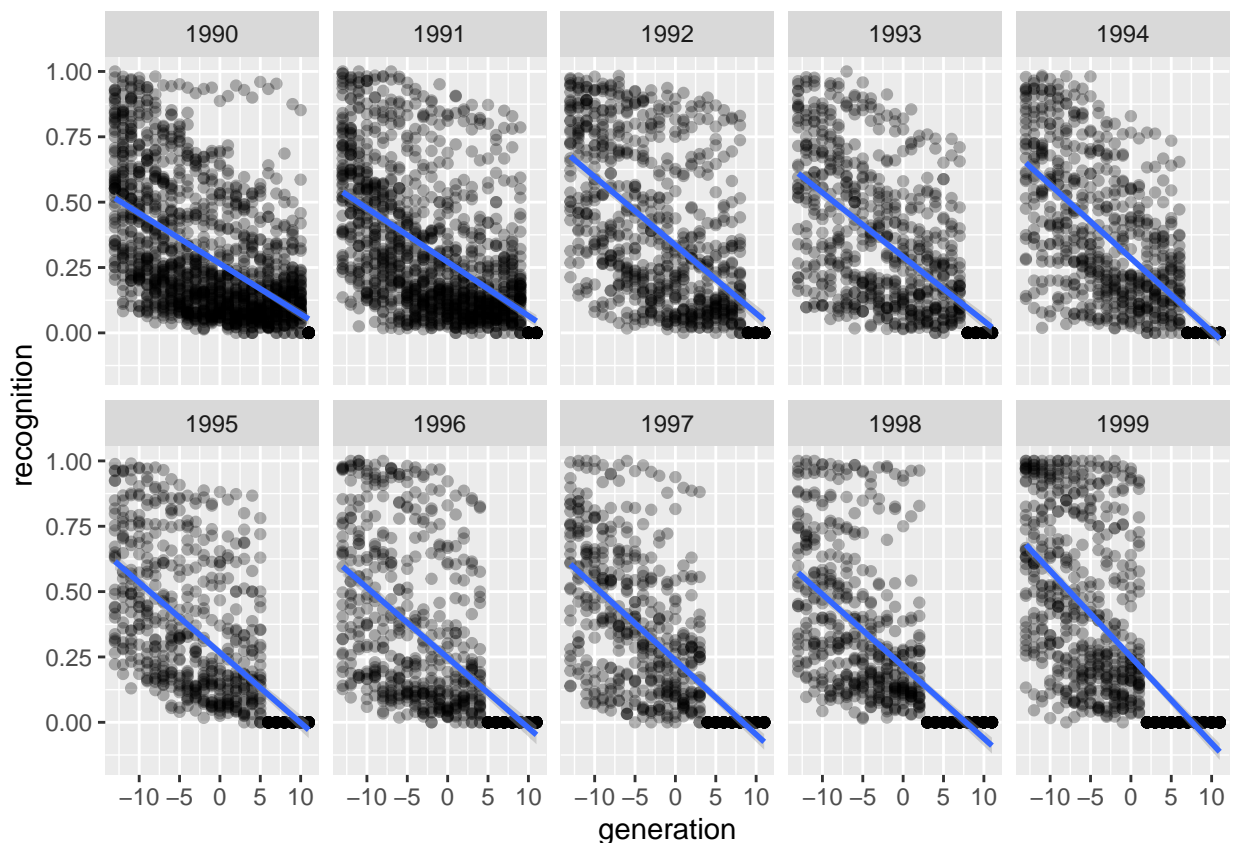
```
song.decay.tbl%>%
  mutate(year = as.factor(year))%>%
  ggplot(aes(x = generation, y = recognition))+
    geom_jitter(alpha = 0.3)+
    geom_smooth(method = "lm")+
    facet_wrap(~year)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```



```
ggplot(song.decay.tbl, aes(x = generation, y = recognition)) +
  geom_point(alpha=0.3) +
  geom_smooth(method="lm")+
  facet_wrap(vars(year), nrow=2)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

```
song.decay.tbl%>%
  filter(year == 1999, generation > 4)
```

```
## # A tibble: 238 x 5
##   artist      song      generation recognition  year
##   <chr>      <chr>      <dbl>      <dbl>  <dbl>
## 1 702        Where My Girls At?      5          0  1999
## 2 Brandy     Have You Ever?          5          0  1999
## 3 Brian McKnight Back At One             5          0  1999
## 4 Britney Spears ..Baby One More Time     5          0  1999
## 5 Busta Rhymes What's It Gonna Be?!     5          0  1999
## 6 Cher       Believe                5          0  1999
## 7 Christina Aguilera Genie In A Bottle       5          0  1999
## 8 Destiny's Child Bills, Bills, Bills      5          0  1999
## 9 Enrique Iglesias Bailamos                 5          0  1999
## 10 Jennifer Lopez  If You Had My Love      5          0  1999
## # ... with 228 more rows
```

As the years get closer to 1999, we get more y values that are 0 because they don't have data on people who are younger (maybe laws prohibiting them from collecting data from minors?) It does appear that the cut off is whenever the generation is 18 - for example, generation 11 in 1990 represent people who would have been ~17-18 in 2018 so it makes sense that there wouldn't be data on them. As the year increases, this becomes true for more and more generations.

- Let's look at the trends for the top 5 most popular artist by doing a join with `top5.tbl` and subset to only negative generations.

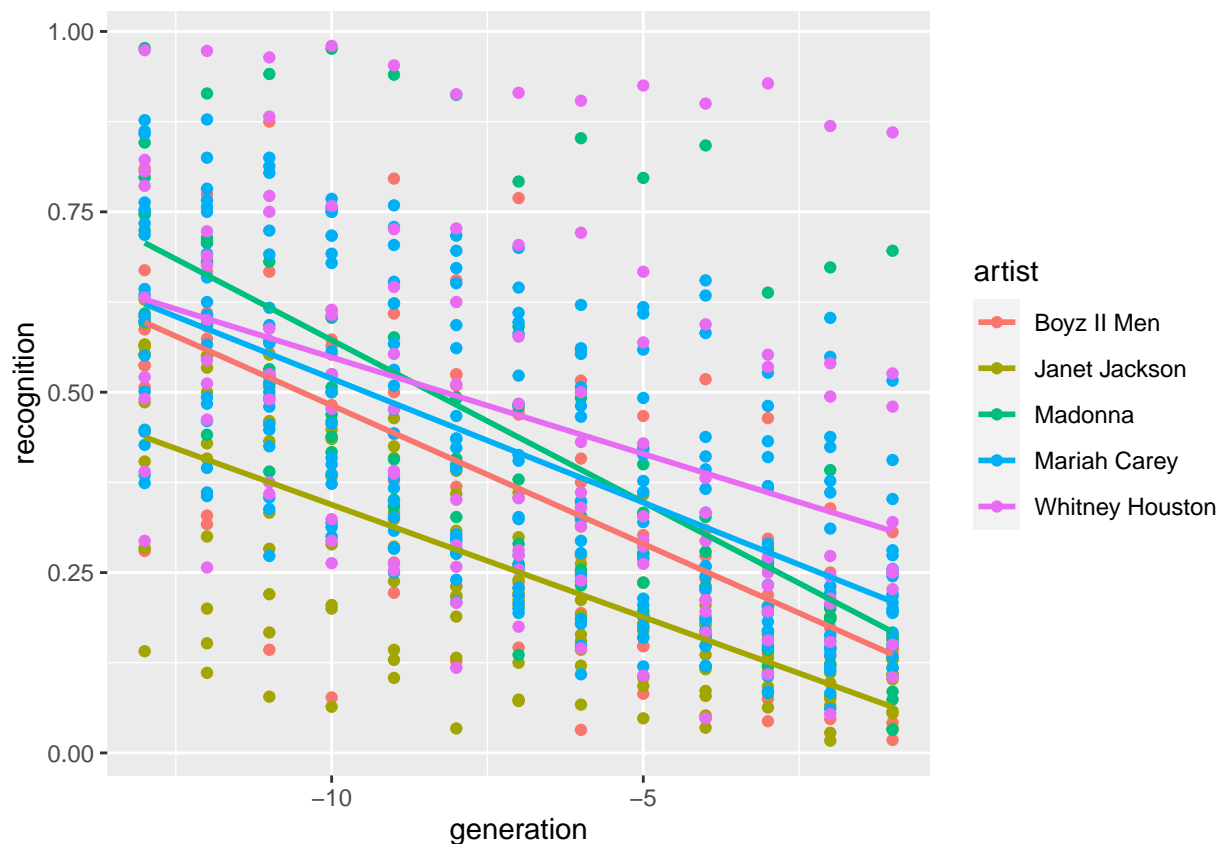
```
top5.decay <- song.decay.tbl%>%
  inner_join(top5.tbl)%>%
  filter(generation<0)
```

```
## Joining, by = "artist"
```

1. Plot the trends for every artist. Who is the most recognized artist (dare I say *diva*) from the 90s

```
top5.decay %>%
  ggplot(aes(x = generation, y = recognition, color = artist))+
  geom_point()+
  geom_smooth(method = "lm", se = FALSE)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

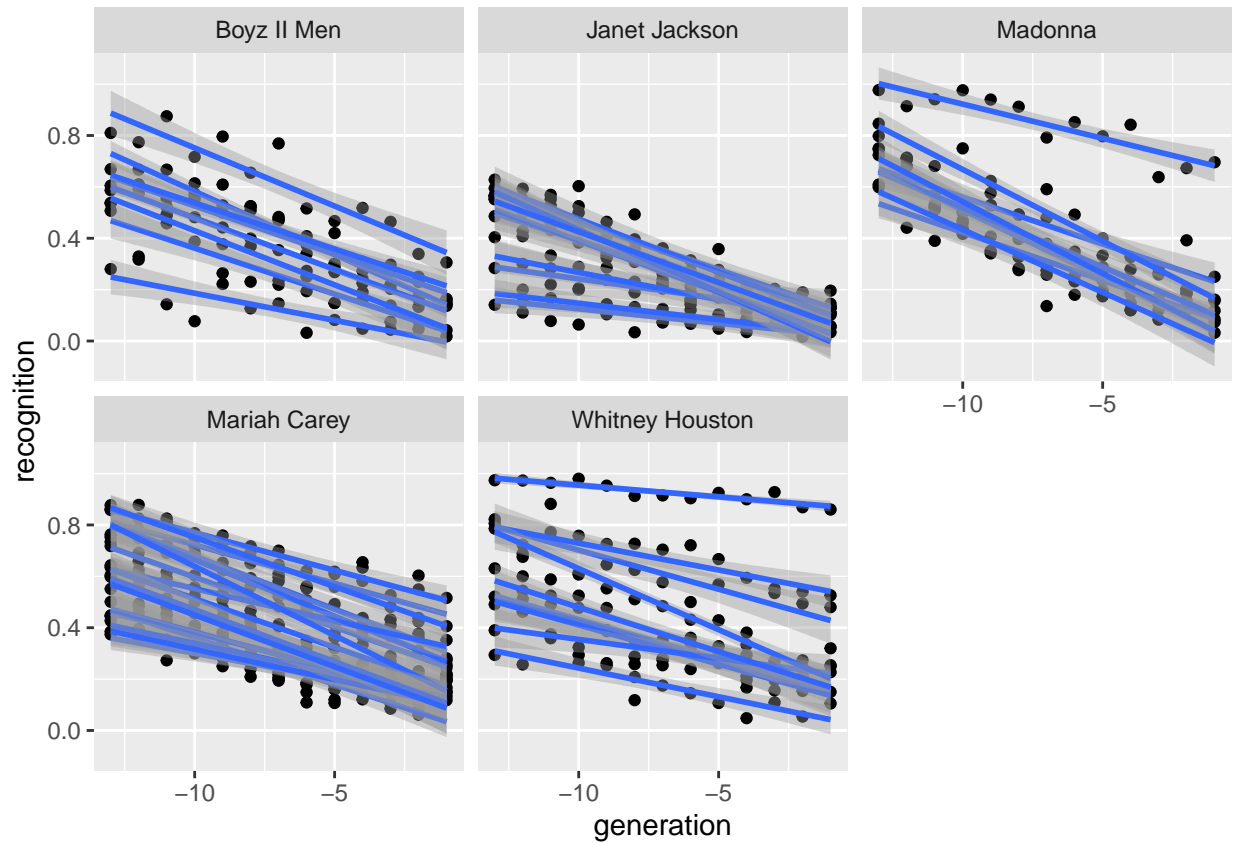


The most recognized artist appears to be Whitney Houston, although Madonna is more recognized among the older generations (-13 to -9) and has the highest single recognition value of close to 0.75 among generation -13. After that generation -9 though, Whitney Houston becomes the most recognized.

2. Look at the individual trends for each artist by generating the following plot.

```
top5.decay %>%
  group_by(artist, song)%>%
  ggplot(aes(x = generation, y = recognition))+
  geom_point()+
  geom_smooth(aes(group = song), method = "lm")+
  facet_wrap(~artist)
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

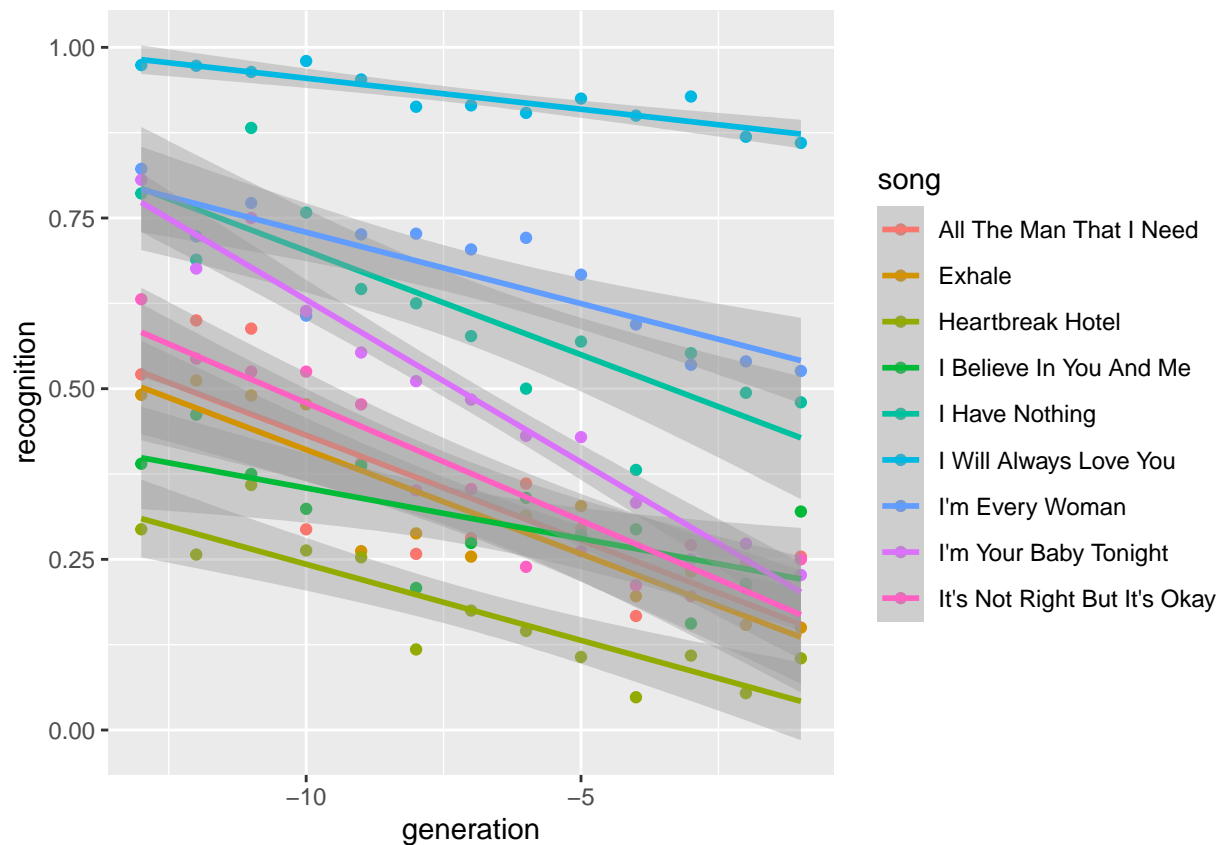


6. In the following exercise let's focus on songs by Whitney Houston.

1. Identify the most and least recognized songs by Whitney Houston by generating the following plot

```
top5.decay %>%
  filter(artist == "Whitney Houston")%>%
  ggplot(aes(x = generation, y = recognition, color = song))+
    geom_point()+
    geom_smooth(method = "lm")

## 'geom_smooth()' using formula 'y ~ x'
```



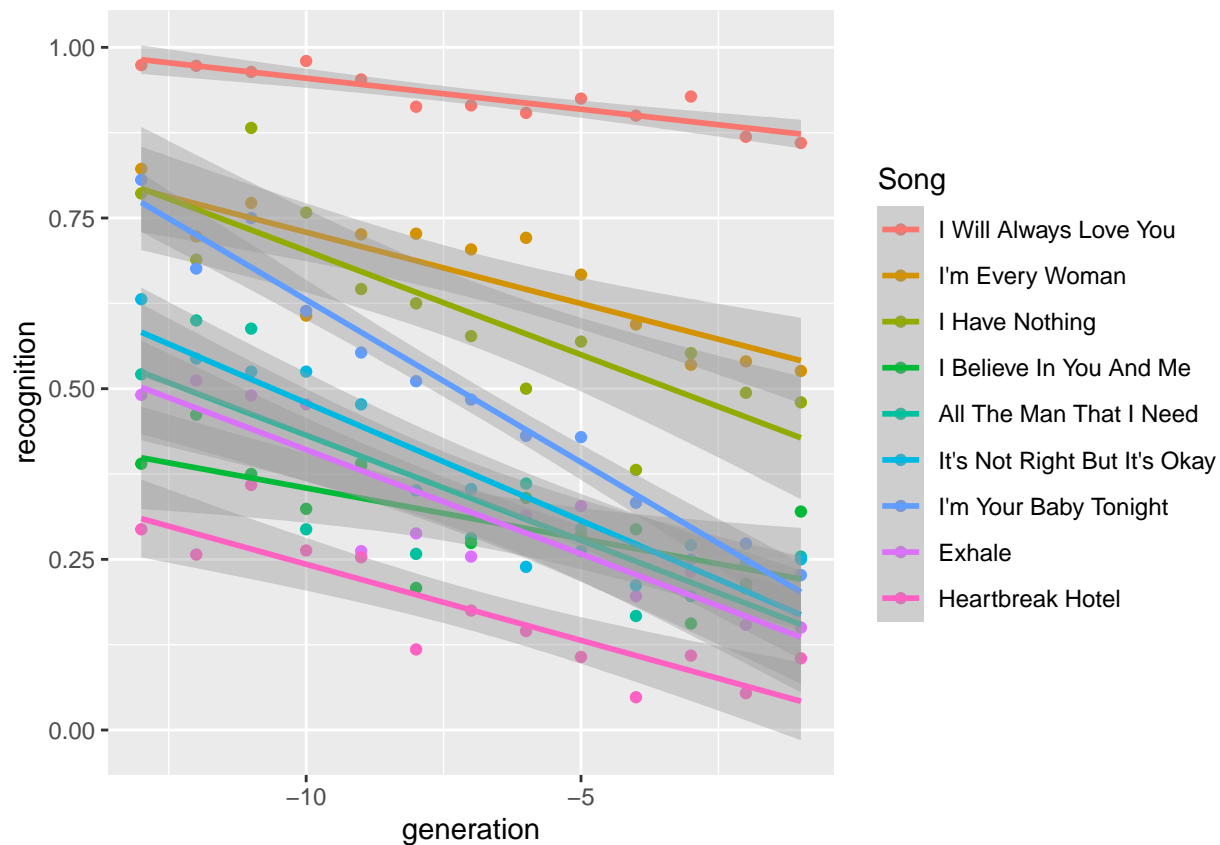
#My color scheme different than the plot given in pdf but same data representation / results

2. Modify the previous plot so that the names are in order of recognition (*Hint*: Create a new category)

#I used fct_reorder2 instead of creating a categorical variable

```
top5.decay %>%
  filter(artist == "Whitney Houston")%>%
  ggplot(aes(x = generation, y = recognition, color = fct_reorder2(song, generation, recognition)))+
    geom_point()+
    geom_smooth(method = "lm")+
    labs(color = "Song")
```

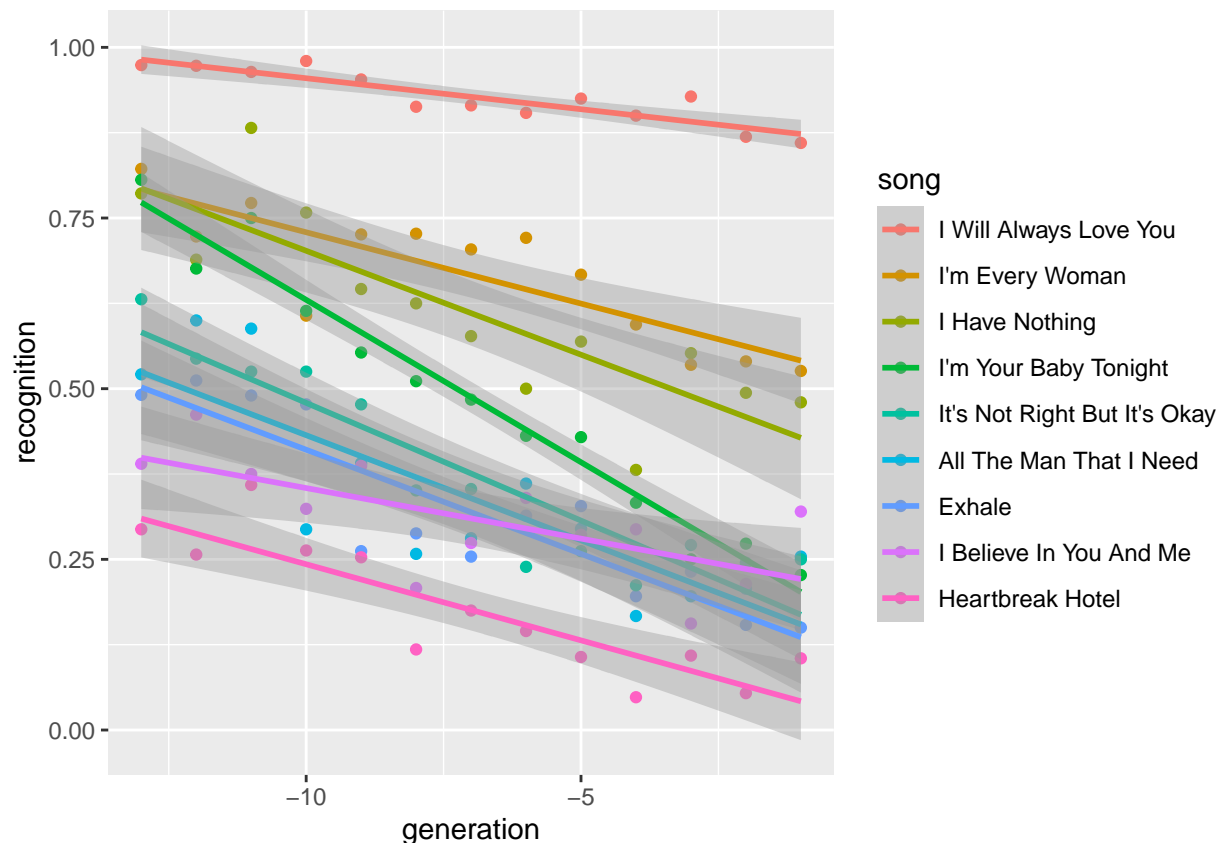
'geom_smooth()' using formula 'y ~ x'



```
#Here's the other way by creating the levels of the song
#and then changing the factor of the song column in the
#tibble when graphing
song.levels <- (top5.decay %>%
  filter(artist == "Whitney Houston")%>%
  group_by(song)%>%
  summarize(avg_rec = mean(recognition))%>%
  arrange(-avg_rec))$song

top5.decay %>%
  filter(artist == "Whitney Houston")%>%
  mutate(song = factor(song, levels = song.levels))%>%
  ggplot(aes(x = generation, y = recognition, color = song))+
    geom_point()+
    geom_smooth(method = "lm")

## 'geom_smooth()' using formula 'y ~ x'
```



7. In this exercise we will plot the trends of particular sets of songs. For the purposes of the exercise we will subset `song.decay.tbl` to songs before generation 0

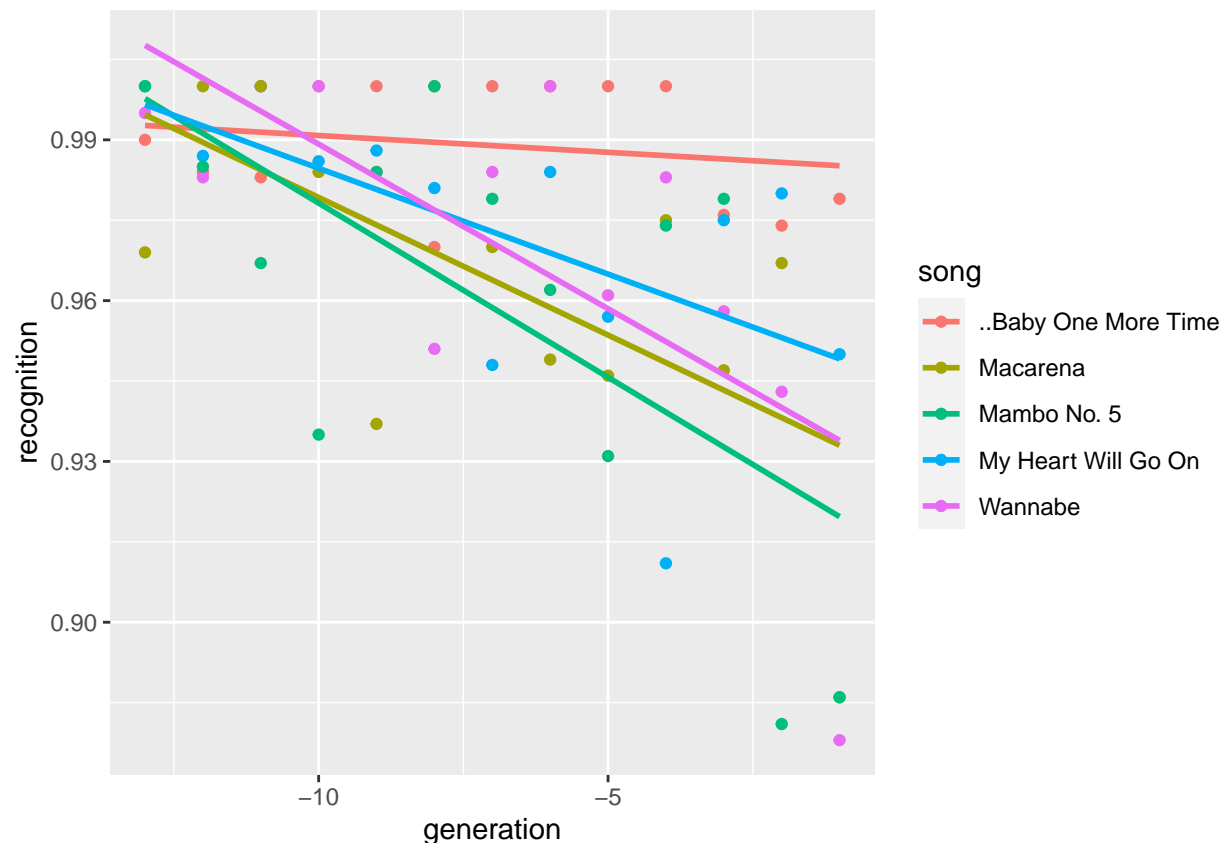
```
song.decay.filter.tbl <- song.decay.tbl%>%
  filter(generation<0)
```

1. Identify the top-5 songs with the highest recognition and plot their trends. Make sure to use 'inner.

```
top5_recognition <- song.decay.filter.tbl%>%
  group_by(song)%>%
  summarize(avg_recognition = mean(recognition))%>%
  slice_max(avg_recognition, n = 5)

top5_recognition%>%
  inner_join(song.decay.filter.tbl)%>%
  ggplot(aes(x = generation, y = recognition, color = song))+
    geom_point()+
    geom_smooth(method = "lm", se = FALSE)
```

```
## Joining, by = "song"
## 'geom_smooth()' using formula 'y ~ x'
```



The top 5 songs are “..Baby One More Time” (I’m assuming this is “Hit Me Baby One More Time”), Macarena, Mambo No. 5, My Heart Will Go On, and Wannabe

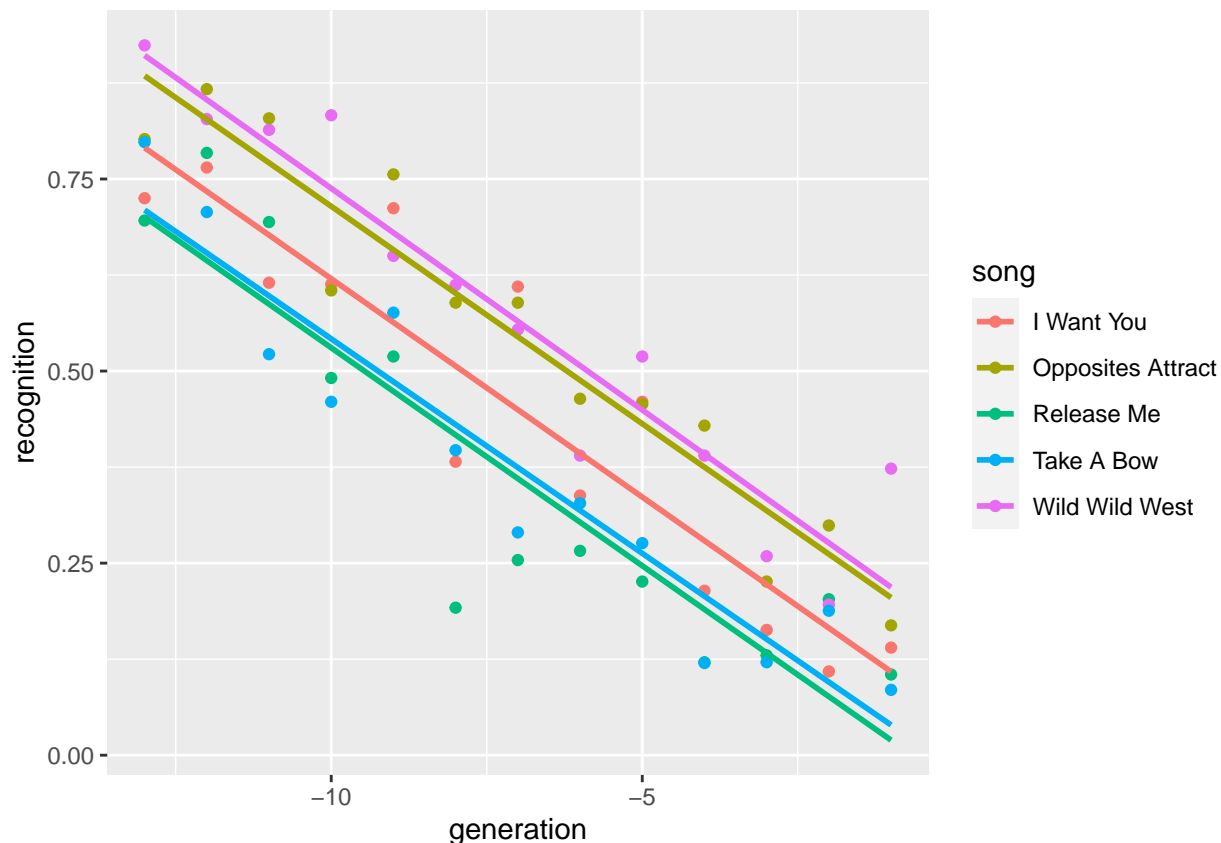
1. Identify the top-5 songs with the highest variability and plot their trends. Make sure to use ‘inner.’

#Unsure what 'variability' refers to here, but I decided to use variance of recognition among the songs

```
top5_variance <- song.decay.filter.tbl%>%
  group_by(song)%>%
  summarize(variance = var(recognition))%>%
  slice_max(variance, n = 5)

top5_variance%>%
  inner_join(song.decay.filter.tbl)%>%
  ggplot(aes(x = generation, y = recognition, color = song))+
  geom_point()+
  geom_smooth(method = "lm", se = FALSE)
```

```
## Joining, by = "song"
## 'geom_smooth()' using formula 'y ~ x'
```



Using pivots

Make sure to read Chapter 12: Tidy data and keep in mind the following commands:

- `pivot_longer`
- `pivot_wider`
- `separate`

8. We are interested in calculating how people from different generation remember songs from the 90s. In particular we are interested in calculating the average song recognition across millennials (people born between 1980 to 1994) and generation Z (people born after 1994)

To do that create a new table using the following steps:

- Calculate the year a person was born based on the `year` and `generation` fields.
- Keep only entries of people who were born starting in 1980.
- People who were born between 1980 and 1994 will be millennials and the rest will be generation Z.
- Finally calculate the average recognition grouping by the generational group
- Name the resulting table `song.gen.tbl`

```
#Version where we don't get rid of years where no data exists
# song.gen.tbl2 <- song.decay.tbl%>%
#   mutate(year_born = year + generation)%>%
#   filter(year_born >= 1980)%>%
#   mutate(generation.group = ifelse(year_born <= 1994, "Millenial", "Z"))%>%
```



```
# group_by(song, generation.group)%>%
# summarize(avg_rec = mean(recognition))

#Option used in PDF where years with 0 recognition (no data present) are excluded
song.gen.tbl <- song.decay.tbl%>%
  mutate(year_born = year + generation)%>%
  filter(year_born >= 1980)%>%
  mutate(generation.group = ifelse(year_born <= 1994, "Millenial", "Z"))%>%
  filter(recognition != 0)%>%
  group_by(song, generation.group)%>%
  summarize(avg_rec = mean(recognition))
```

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

```
song.gen.tbl
```

```
## # A tibble: 684 x 3
## # Groups:   song [342]
##   song                generation.group avg_rec
##   <chr>                <chr>          <dbl>
## 1 ..Baby One More Time Millenial      0.992
## 2 ..Baby One More Time Z            0.963
## 3 4 Seasons Of Loneliness Millenial    0.139
## 4 4 Seasons Of Loneliness Z           0.066
## 5 A Whole New World    Millenial      0.868
## 6 A Whole New World    Z            0.706
## 7 Achy Breaky Heart    Millenial      0.838
## 8 Achy Breaky Heart    Z            0.562
## 9 Adia                 Millenial      0.387
## 10 Adia                 Z            0.119
## # ... with 674 more rows
```

9. As we can see `song.gen.tbl` has one different row for each `generation.group`, however we would like to have just one row per song and have columns with the average recognition for each generation. Use a pivot function to obtain the following table and name it `song.pivot.tbl`

```
song.pivot.tbl <- song.gen.tbl%>%
  pivot_wider(names_from = generation.group, values_from = avg_rec)%>%
  arrange(-Millenial)
```

```
song.pivot.tbl
```

```
## # A tibble: 342 x 3
## # Groups:   song [342]
##   song                Millenial      Z
##   <chr>                <dbl> <dbl>
## 1 ..Baby One More Time    0.992 0.963
## 2 Wannabe                 0.982 0.904
## 3 Believe                 0.978 0.899
## 4 My Heart Will Go On    0.974 0.962
## 5 Mambo No. 5             0.971 0.928
## 6 Macarena                0.970 0.862
## 7 Everybody               0.965 0.864
## 8 I Believe I Can Fly     0.958 0.843
## 9 All Star                 0.946 0.932
```

```
## 10 The Power          0.932 0.876
## # ... with 332 more rows
```

10. Consider one of the tables for the original publication `time.series.tbl`

```
url <- "https://raw.githubusercontent.com/the-pudding/song-decay-clean/master/src/assets/data/time_series"
time.series.tbl <- read_csv(url)
time.series.tbl
```

```
## # A tibble: 344 x 49
##   artist_song '-13' '-12' '-11' '-10' '-9' '-8' '-7' '-6' '-5' '-4' '-3'
##   <chr>       <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1 2 Pac|||Ca~ 0.648 0.54  0.712 0.435 0.644 0.594 0.6   0.421 0.446 0.391 0.309
## 2 2Pac|||How~ 0.265 0.226 0.382 0.296 0.230 0.284 0.139 0.169 0.215 0.1   0.109
## 3 702|||Wher~ 0.631 0.614 0.660 0.415 0.532 0.375 0.415 0.276 0.304 0.408 0.213
## 4 Ace Of Bas~ 0.963 0.94  0.976 0.877 0.955 0.954 1     0.853 0.9   0.875 0.779
## 5 Ace Of Bas~ 0.899 0.909 0.827 0.8   0.757 0.773 0.667 0.516 0.557 0.629 0.625
## 6 Ace Of Bas~ 0.979 0.955 0.984 0.933 0.930 0.908 0.943 0.893 0.905 0.841 0.806
## 7 Adina Howa~ 0.415 0.28  0.302 0.297 0.2   0.169 0.179 0.190 0.238 0.172 0.25
## 8 Aerosmith|~ 0.959 0.953 0.966 0.987 0.870 0.931 0.862 0.8   0.848 0.875 0.818
## 9 Aerosmith|~ 0.820 0.846 0.757 0.833 0.683 0.590 0.571 0.549 0.735 0.585 0.478
## 10 Alanis Mor~ 0.970 0.979 0.986 0.95  0.926 0.937 0.973 0.84  0.919 0.778 0.766
## # ... with 334 more rows, and 37 more variables: -2 <dbl>, -1 <dbl>, 0 <dbl>,
## #   1 <dbl>, 2 <dbl>, 3 <dbl>, 4 <dbl>, 5 <dbl>, 6 <dbl>, 7 <dbl>, 8 <dbl>,
## #   9 <dbl>, 10 <dbl>, 11 <dbl>, 12 <dbl>, 13 <dbl>, 14 <dbl>, 15 <dbl>,
## #   16 <dbl>, 17 <dbl>, 18 <dbl>, 19 <dbl>, 20 <dbl>, 21 <dbl>, 22 <dbl>,
## #   23 <dbl>, 24 <dbl>, 25 <dbl>, 26 <dbl>, 27 <dbl>, 28 <dbl>, 29 <dbl>,
## #   30 <dbl>, 31 <dbl>, 32 <dbl>, 33 <dbl>, 34 <dbl>
```

Convert this table into a table with the following format (Note that generation is an integer and recognition is a double) and make sure to subset generation from -13 to 10

```
time.series.tbl%>%
  separate(artist_song, into = c("artist", "song"), sep = "\\|\\|\\|\\|\\|")%>%
  # select(1:26)%>%
  pivot_longer(3:50, names_to = "generation", values_to = "recognition", names_transform = list(generat.
  filter(generation <=10)
```

```
## # A tibble: 8,256 x 4
##   artist song          generation recognition
##   <chr> <chr>          <int>      <dbl>
## 1 2 Pac California Love      -13      0.648
## 2 2 Pac California Love      -12      0.54
## 3 2 Pac California Love      -11      0.712
## 4 2 Pac California Love      -10      0.435
## 5 2 Pac California Love       -9      0.644
## 6 2 Pac California Love       -8      0.594
## 7 2 Pac California Love       -7      0.6
## 8 2 Pac California Love       -6      0.421
## 9 2 Pac California Love       -5      0.446
## 10 2 Pac California Love      -4      0.391
## # ... with 8,246 more rows
```