

Data Science

Shisham Adhikari

04/22/2020

```
# Put all necessary libraries here  
library(tidyverse)  
library(viridis)  
library(ggmap)  
library(rvest)
```

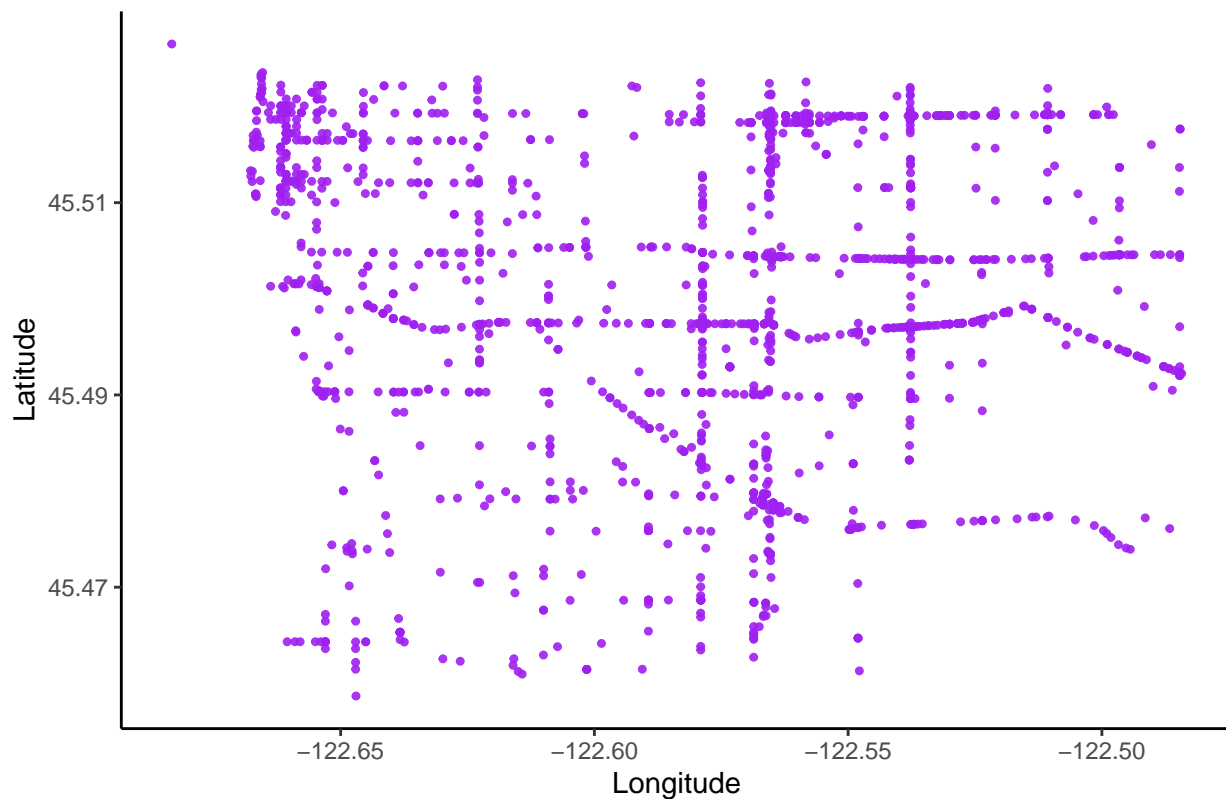
Problem 1: Mapping PDX Crashes

```
pdx_crash_2018 <- read_csv("/home/courses/math241s20/Data/pdx_crash_2018_page1.csv")
```

a. Recreating the longitude vs. latitude graph:

```
pdx_crash_2018 <- read_csv("/home/courses/math241s20/Data/pdx_crash_2018_page1.csv")  
ggplot(data = pdx_crash_2018,  
       aes(x = LONGTD_DD,  
           y = LAT_DD)) +  
  geom_point(alpha = .9,  
            size = .9,  
            color = "purple") +  
  labs(x = "Longitude",  
       y = "Latitude",  
       title = "Crashes Across Portland Based on their Latitudinal and Longitudal Data") +  
  theme_classic()
```

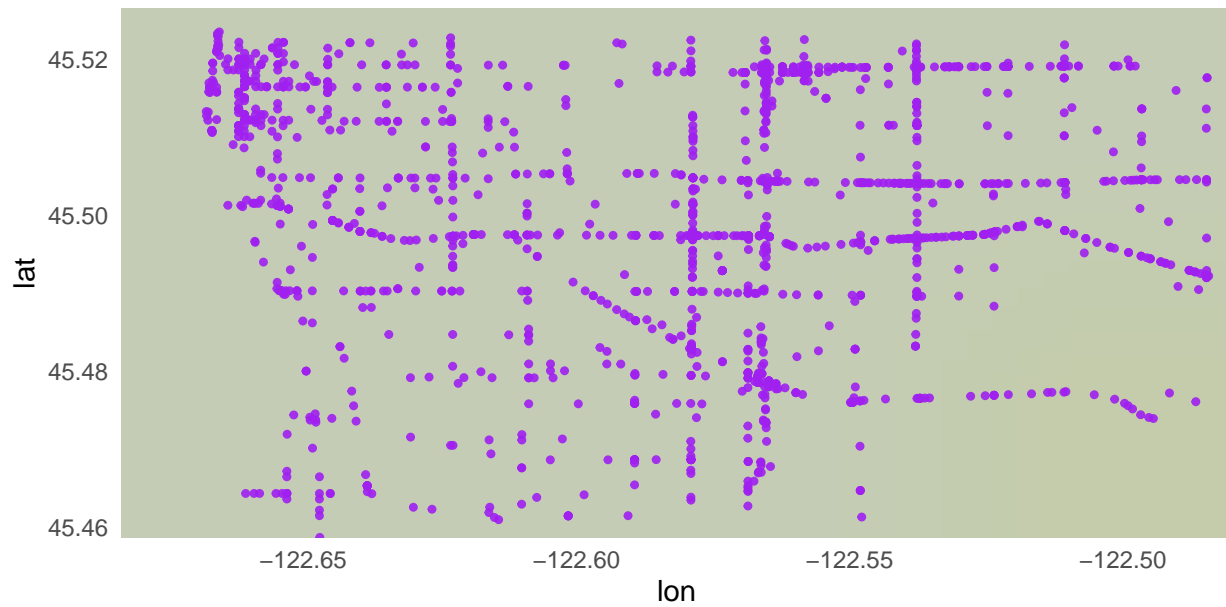
Crashes Across Portland Based on their Latitudinal and Longitudal Data



b.A (static) raster map with the crashes mapped as points on top:

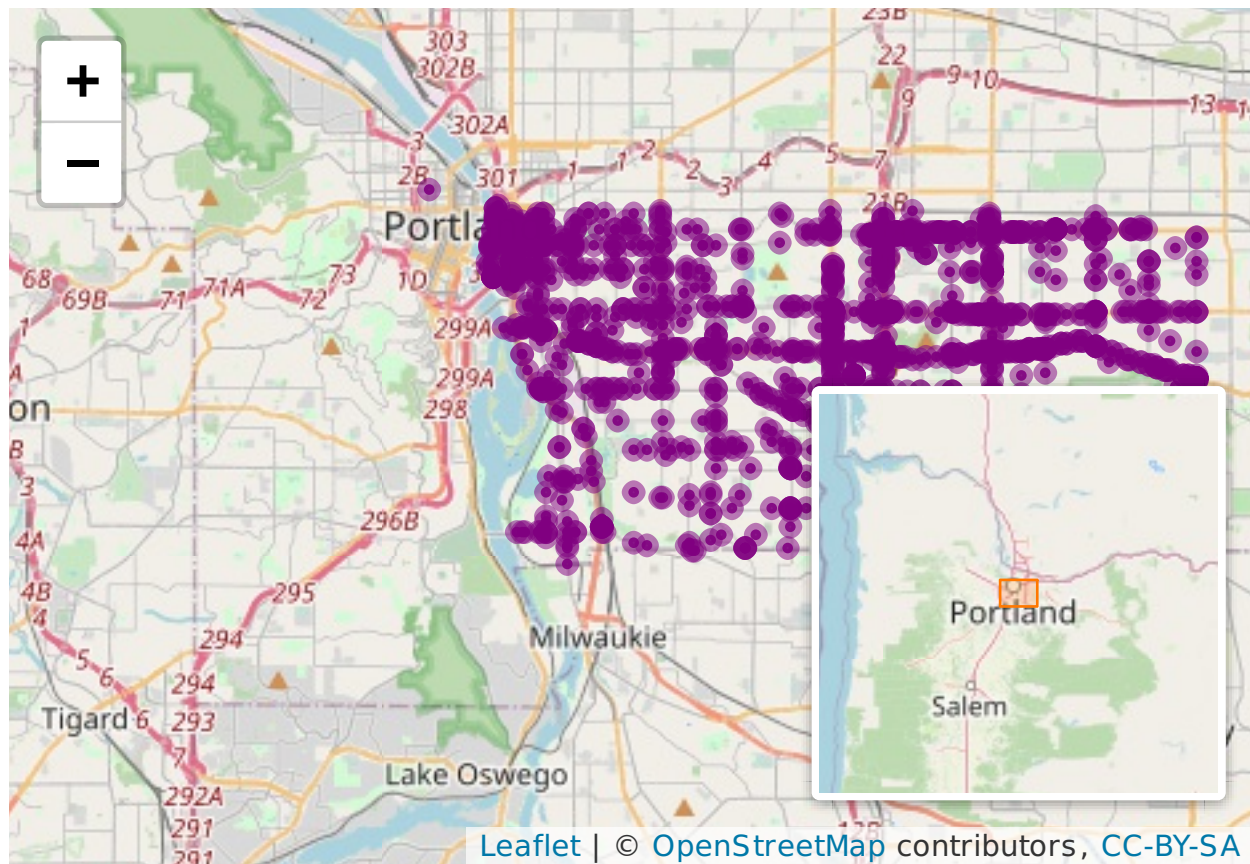
```
aleutian_box <- c(bottom = 45.45868, left = -122.6833,
                  top = 45.5265, right = -122.4800)
aleutian <- get_stamenmap(aleutian_box,
                          maptype = "terrain-background",
                          zoom = 5)

aleutian %>%
  save(aleutian, file = "aleutian.RData")
load("aleutian.RData")
aleutian %>%
  ggmap() +
  geom_point(data = pdx_crash_2018,
            aes(x = LONGTD_DD,
                y = LAT_DD),
            alpha = .9,
            size = .9,
            color = "purple") +
  theme_minimal()
```



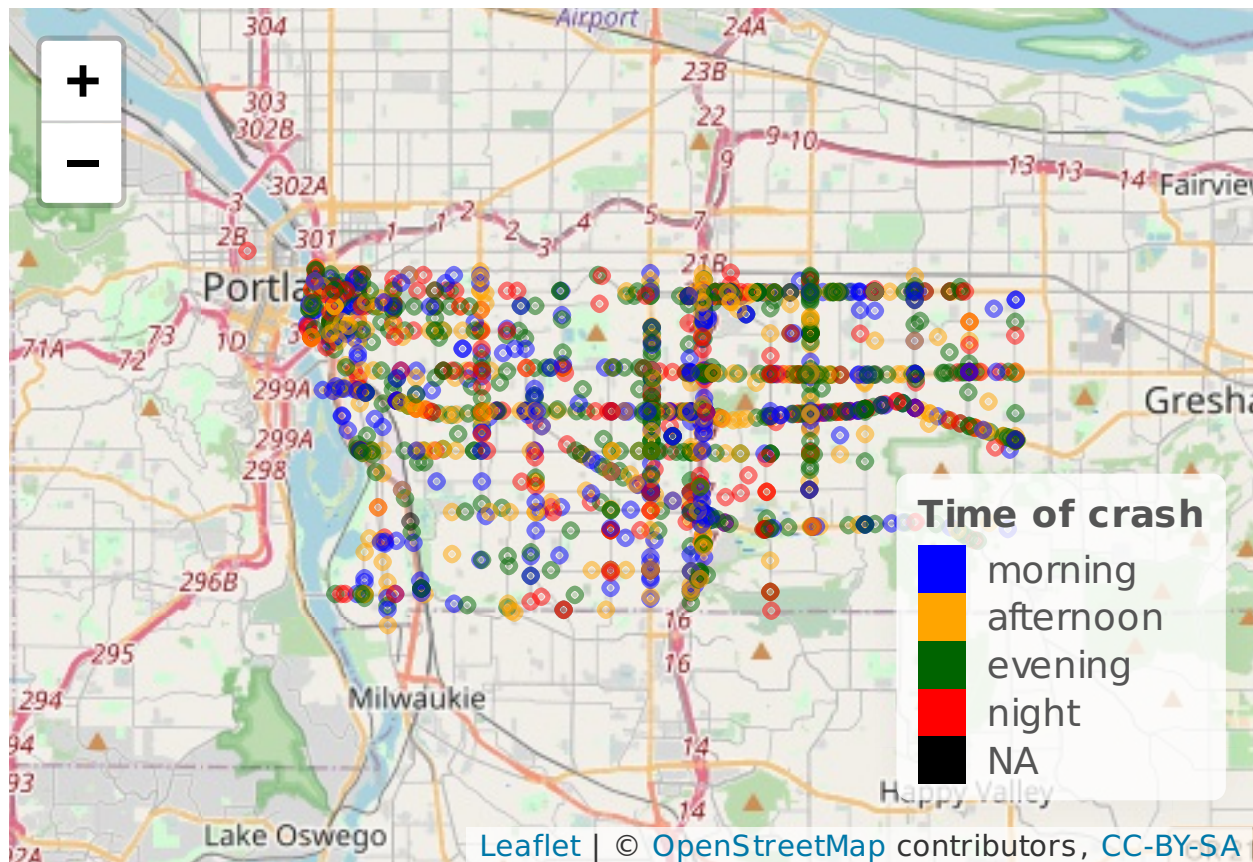
c. An interactive map of the crashes:

```
library(leaflet)
leaflet() %>%
  setView(lng = -122.6308, lat = 45.4811,
    zoom = 11) %>%
  addTiles() %>%
  addCircleMarkers(lng = ~LONGTD_DD , lat = ~LAT_DD,
    data = pdx_crash_2018, color="purple", fillOpacity = 0.9, radius=2) %>%
  addMiniMap()
```



d.

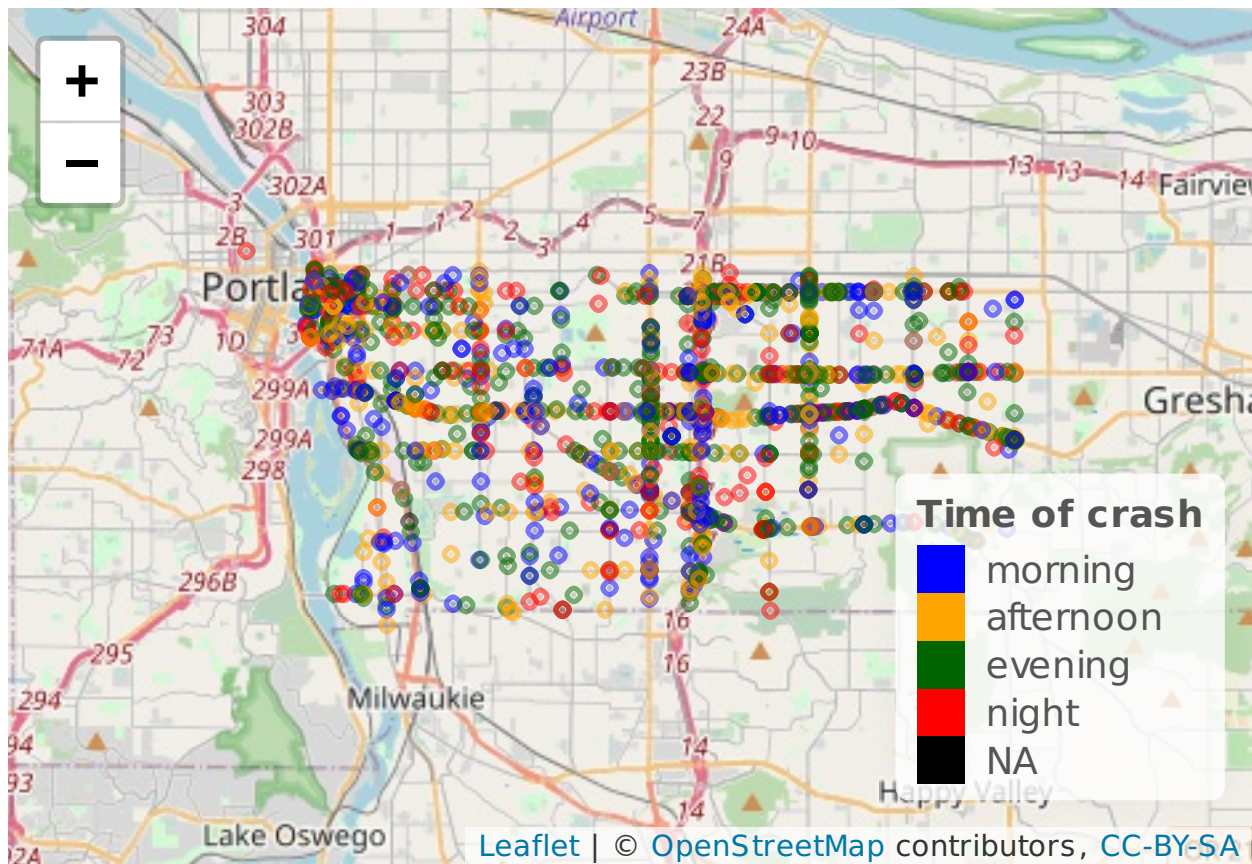
```
pdx_crash_2018 <- pdx_crash_2018 %>%
  mutate(CRASH_HR_NO= as.integer(CRASH_HR_NO)) %>%
  mutate(day_time=factor(CRASH_HR_NO, levels=c(5:11, 12:15, 16:19, 20:23, 00:05, 99),
    labels = c(rep("morning",7), rep("afternoon",4), rep("evening",4), rep("night",4))))
pal <- colorFactor(palette = c("blue", "orange", "darkgreen", "red", "black"), domain = pdx_crash_2018$day_time)
pdx_crash_2018 %>%
  leaflet(options = leafletOptions(minZoom = 10, maxZoom = 15)) %>%
  addTiles() %>%
  addCircles(lng = ~LONGTD_DD, lat = ~LAT_DD,
    color=~pal(day_time)) %>%
  addLegend("bottomright", pal=pal, values=~day_time, title = "Time of crash", opacity=1)
```



Comments: From the plot, we see that crash locations vary by parts of the day but there is no strong trend. Most crashes are happening along highways or big streets (yellow, red lines in the plot) and also in the intersections.

e. Adding a pop-up to the interactive map that provides the exact address of the crash:

```
content <- paste("<b>", pdx_crash_2018$CNTY_NM, pdx_crash_2018$CITY_SECT_NM,
               "</b></br>", "Street of the crash:",
               pdx_crash_2018$ST_FULL_NM, "(", pdx_crash_2018$LAT_DD, ",", pdx_crash_2018$LONGTD_DD,
pdx_crash_2018 %>%
  leaflet(options = leafletOptions(minZoom = 10, maxZoom = 15)) %>%
  addTiles() %>%
  addCircles(lng = ~LONGTD_DD, lat = ~LAT_DD,
             color=~pal(day_time), popup = content) %>%
  addLegend("bottomright", pal=pal, values=~day_time, title = "Time of crash", opacity=1)
```

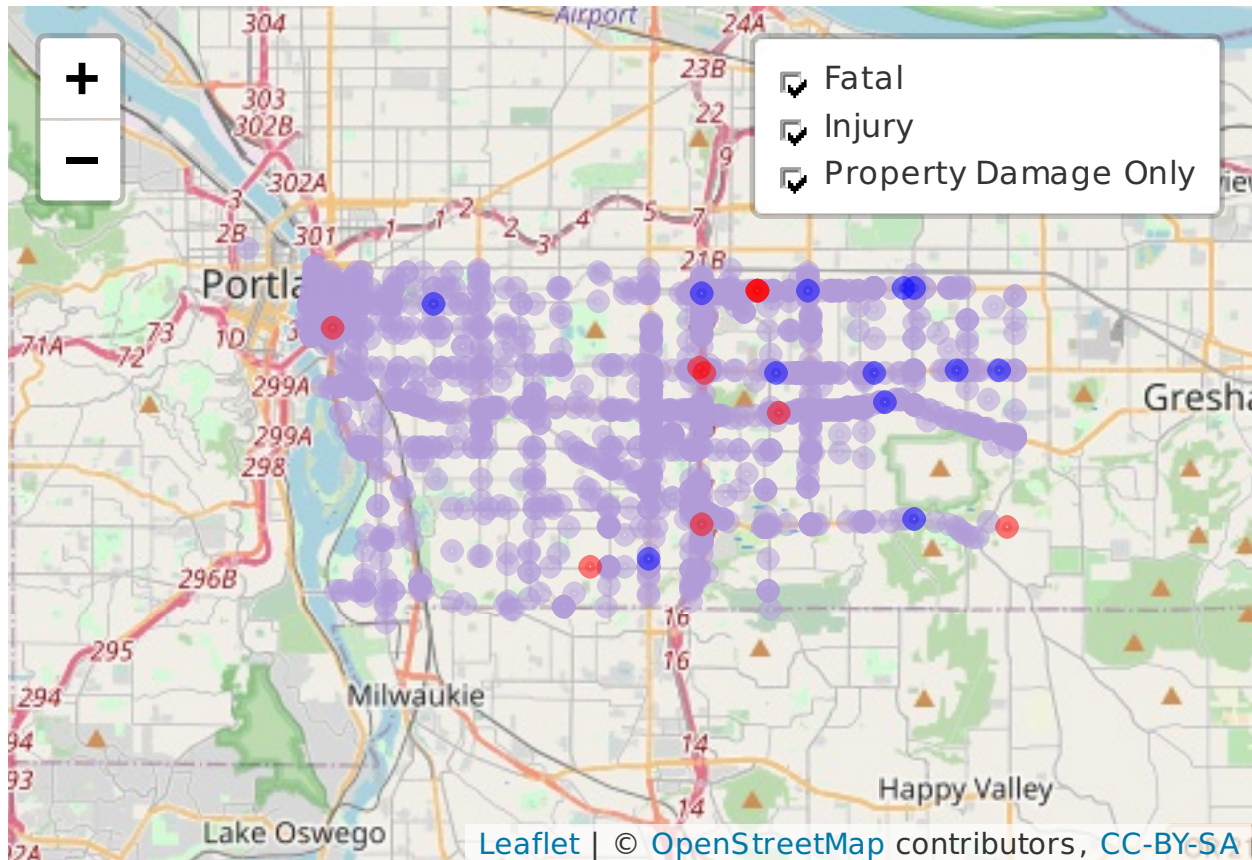



f. A leaflet graph displaying the different crash severities:

```
fatal <- pdx_crash_2018 %>%
  filter(CRASH_SVRTY_SHORT_DESC == "FAT")
injury <- pdx_crash_2018 %>%
  filter(CRASH_SVRTY_SHORT_DESC == "INJ")
property_damage <- pdx_crash_2018 %>%
  filter(CRASH_SVRTY_SHORT_DESC == "PDO")

pdx_crash_2018 %>%
  leaflet() %>%
  addTiles() %>%
    addCircleMarkers(lng = ~LONGTD_DD, lat = ~LAT_DD,
                     radius = 2,
                     data = injury, color = "#b19cd9",
                     group = "Injury") %>%
    addCircleMarkers(lng = ~LONGTD_DD, lat = ~LAT_DD,
                     radius = 2,
                     data = property_damage, color = "#ff0000",
                     group = "Property Damage Only") %>%
    addCircleMarkers(lng = ~LONGTD_DD, lat = ~LAT_DD,
                     radius = 2,
                     data = fatal, color = "Blue",
                     group = "Fatal") %>%
  # Layers control
  addLayersControl(
    overlayGroups = c("Fatal", "Injury", "Property Damage Only"),
```

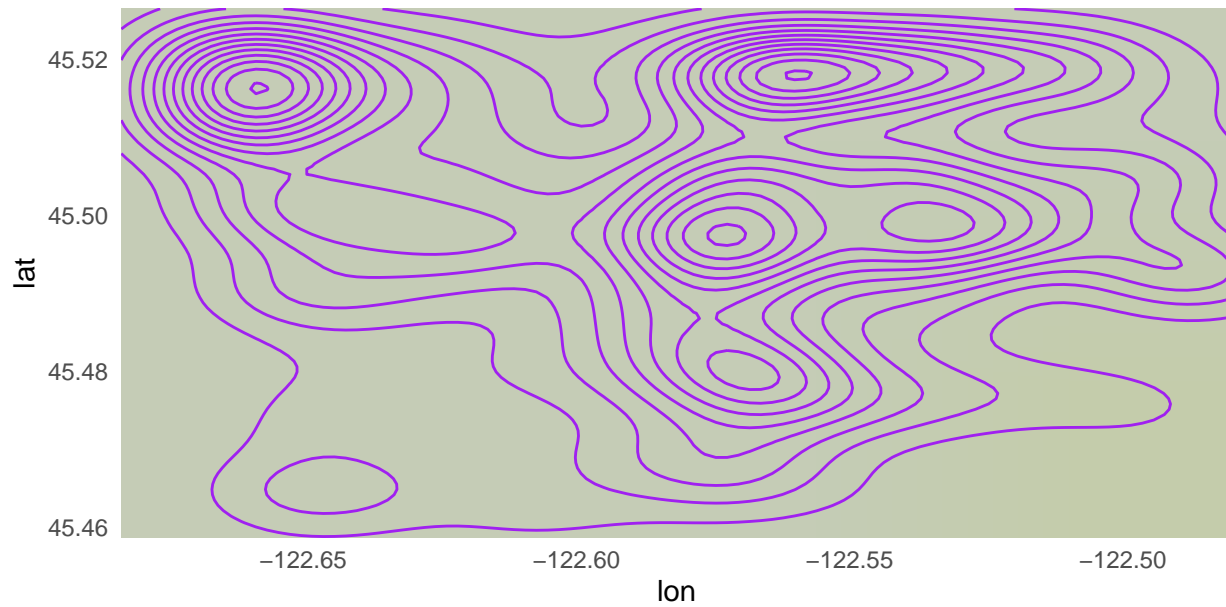
```
options = layersControlOptions(collapsed = FALSE))
```



We see that most of the crashes are with injuries and is common across all locations. We see the fatal crashes are mostly on eastern regions and on big highways and streets (red and orange lines in the map). The crashes with only property destroyed are fewer in number but are spread out. In the map, they are along the highway near downtown or other big streets and highway (205) which are red and orange lines in the plot.

g.

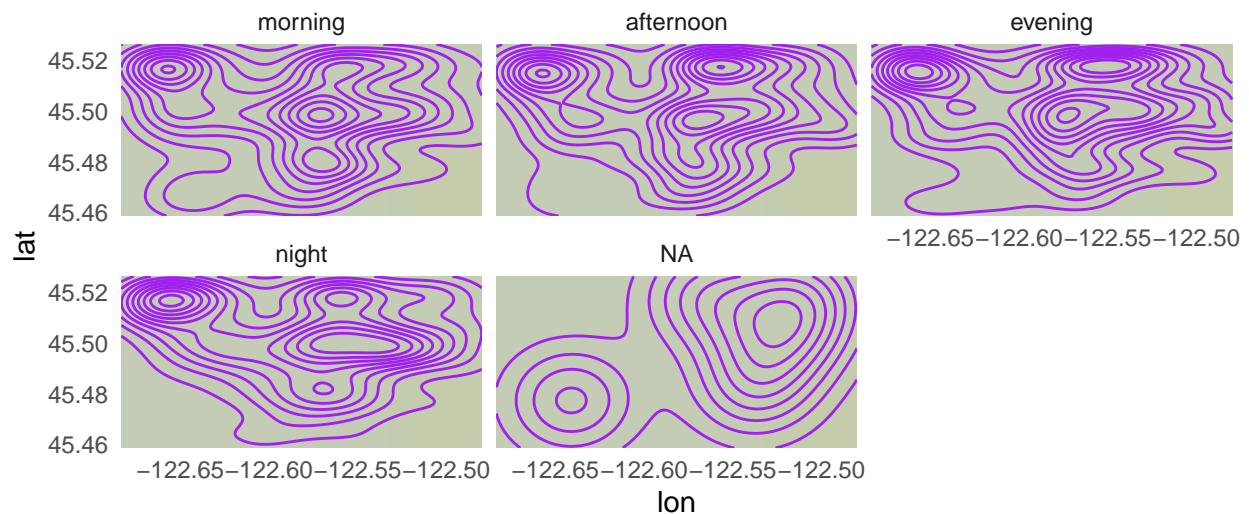
```
aleutian %>%
  ggmap() +
  geom_density2d(data = pdx_crash_2018,
    aes(x = LONGTD_DD,
      y = LAT_DD,
      color = "purple") +
  theme_minimal()
```



This map tells us that the car crashes in the SE are mostly centered around the highway on the way to downtown (upper, left) and the highway, probably 205 (top, right). There are also more crashes around the central SE region. This story is similar to the map using `geom_point()` but the crashes clustered around the center and top-right is not as visible in the map with `geom_point()`. The crashes centered near downtown (top-left) is pretty visible in both the plots.

h.

```
aleutian %>%
  ggmap() +
  geom_density2d(data = pdx_crash_2018,
    aes(x = LONGTD_DD,
      y = LAT_DD,
      color = "purple") +
  theme_minimal()+facet_wrap(~day_time)
```



For all parts of day, the distribution on accidents does not vary as much. The crashes are still centered around the highway near downtown (top, left) and still around the highway, probably 205 (top right) for afternoon and evening. Just like before, they all seem to have accidents around center too.

Problem 2: Choropleth Maps

```
api_key <- "e20f7e545ee9474d9d353d401cc0e1d00d31deeb"
```

a. Grabbing data on the median gross rent for Multnomah county:

```
library(tidycensus)
library(tigris)
county <- get_acs(geography = "county subdivision",
                  variables = "B25064_001",
                  county = "multnomah",
                  state = "Oregon",
                  geometry = TRUE,
                  key = api_key,
                  cache_table = TRUE)
```

##

		0%
	==	2%
	==	3%
	====	5%
	====	6%
	=====	8%
	=====	9%
	=====	11%
	=====	12%
	=====	14%
	=====	15%
	=====	17%
	=====	18%
	=====	20%
	=====	21%
	=====	23%
	=====	24%
	=====	26%
	=====	27%

=====	29%
=====	30%
=====	32%
=====	33%
=====	35%
=====	36%
=====	38%
=====	39%
=====	41%
=====	42%
=====	43%
=====	44%
=====	46%
=====	47%
=====	49%
=====	50%
=====	52%
=====	53%
=====	55%
=====	56%
=====	58%
=====	59%
=====	61%
=====	62%
=====	64%
=====	65%
=====	67%

			68%
	=====		
	=====		70%
	=====		
	=====		71%
	=====		
	=====		73%
	=====		
	=====		74%
	=====		
	=====		76%
	=====		
	=====		77%
	=====		
	=====		79%
	=====		
	=====		80%
	=====		
	=====		82%
	=====		
	=====		83%
	=====		
	=====		85%
	=====		
	=====		86%
	=====		
	=====		88%
	=====		
	=====		89%
	=====		
	=====		91%
	=====		
	=====		92%
	=====		
	=====		94%
	=====		
	=====		95%
	=====		
	=====		97%
	=====		
	=====		98%
	=====		
	=====		100%

```
tract <- get_acs(
  geography = "tract",
  variables = "B25064_001",
  state = "OR",
  county = "Multnomah",
  geometry = TRUE,
  key = api_key
)
```

```
##
```

	0%
=	2%
==	2%
===	4%
====	5%
=====	6%
=====	7%
=====	9%
=====	9%
=====	11%
=====	12%
=====	13%
=====	14%
=====	16%
=====	17%
=====	18%
=====	19%
=====	20%
=====	21%
=====	23%
=====	24%
=====	25%
=====	26%
=====	27%
=====	28%
=====	30%
=====	31%

	=====		32%
	=====		33%
	=====		34%
	=====		35%
	=====		37%
	=====		38%
	=====		39%
	=====		40%
	=====		41%
	=====		42%
	=====		44%
	=====		45%
	=====		46%
	=====		47%
	=====		48%
	=====		49%
	=====		51%
	=====		52%
	=====		53%
	=====		54%
	=====		55%
	=====		56%
	=====		58%
	=====		59%
	=====		60%
	=====		61%
	=====		62%

	=====		63%
	=====		65%
	=====		66%
	=====		67%
	=====		68%
	=====		69%
	=====		70%
	=====		72%
	=====		73%
	=====		74%
	=====		75%
	=====		77%
	=====		79%
	=====		80%
	=====		81%
	=====		82%
	=====		84%
	=====		84%
	=====		86%
	=====		87%
	=====		88%
	=====		89%
	=====		91%
	=====		91%
	=====		93%
	=====		94%
	=====		95%

	=====	96%
	=====	98%
	=====	98%
	=====	100%

```

block <- get_acs(
  geography = "block group",
  variables = "B25064_001",
  state = "OR",
  county = "Multnomah",
  geometry = TRUE,
  key = api_key
)

```

##

		0%
	=	1%
	==	2%
	==	3%
	===	4%
	====	5%
	====	6%
	=====	6%
	=====	7%
	=====	8%
	=====	8%
	=====	9%
	=====	10%
	=====	11%
	=====	11%
	=====	12%
	=====	12%
	=====	13%

=====			14%
=====			15%
=====			15%
=====			16%
=====			16%
=====			17%
=====			18%
=====			19%
=====			20%
=====			21%
=====			22%
=====			23%
=====			24%
=====			25%
=====			26%
=====			27%
=====			28%
=====			29%
=====			30%
=====			31%
=====			32%
=====			32%
=====			33%
=====			34%
=====			35%
=====			36%
=====			37%

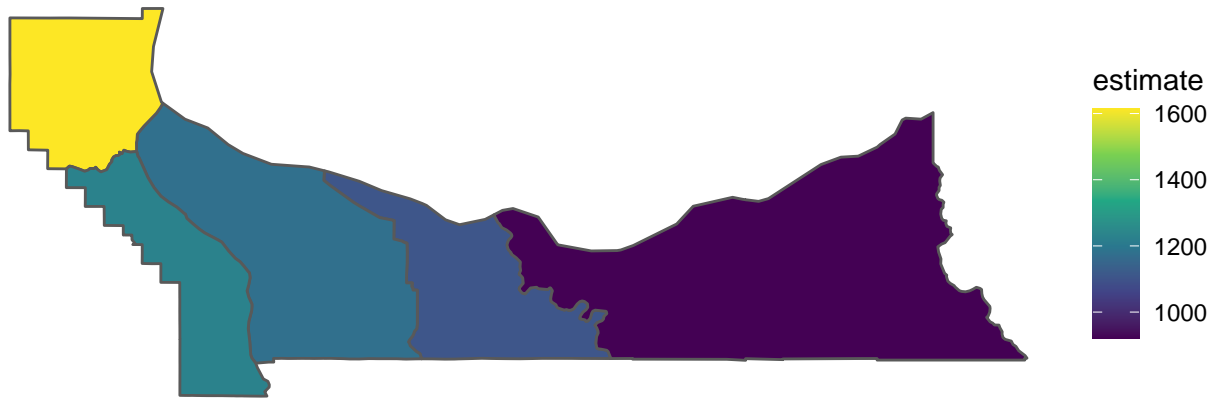
=====	38%
=====	39%
=====	40%
=====	41%
=====	41%
=====	42%
=====	42%
=====	43%
=====	44%
=====	45%
=====	45%
=====	46%
=====	47%
=====	48%
=====	49%
=====	49%
=====	50%
=====	51%
=====	52%
=====	53%
=====	54%
=====	55%
=====	56%
=====	57%
=====	58%
=====	58%
=====	59%

=====	60%
=====	61%
=====	62%
=====	63%
=====	64%
=====	65%
=====	66%
=====	67%
=====	68%
=====	69%
=====	70%
=====	71%
=====	72%
=====	73%
=====	74%
=====	75%
=====	75%
=====	76%
=====	77%
=====	78%
=====	78%
=====	79%
=====	79%
=====	80%
=====	81%
=====	82%
=====	83%

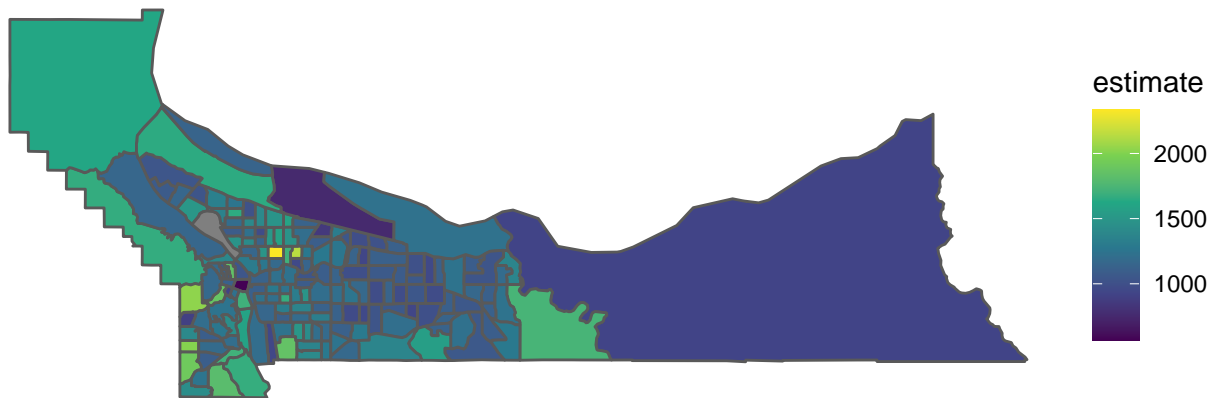
	=====		84%
	=====		85%
	=====		86%
	=====		87%
	=====		88%
	=====		88%
	=====		89%
	=====		90%
	=====		91%
	=====		92%
	=====		93%
	=====		94%
	=====		95%
	=====		96%
	=====		97%
	=====		98%
	=====		99%
	=====		99%
	=====		100%

b. Three choropleth maps of gross rent, one for each geography resolution:

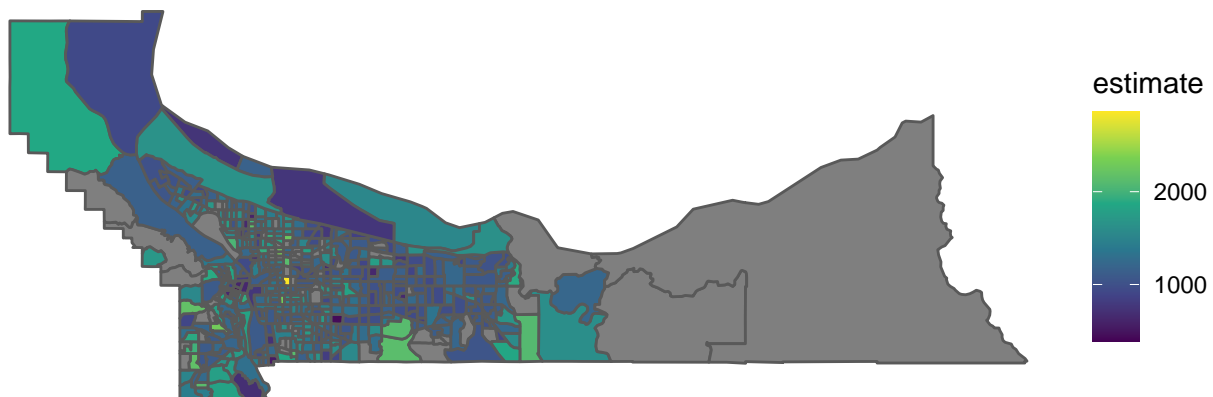
```
ggplot(data = county, mapping = aes(geometry = geometry)) +
  geom_sf(aes(fill = estimate)) +
  coord_sf() +
  scale_fill_viridis_c() + theme_void()
```



```
ggplot(data = tract, mapping = aes(geometry = geometry)) +  
  geom_sf(aes(fill = estimate)) +  
  coord_sf() +  
  scale_fill_viridis_c() + theme_void()
```



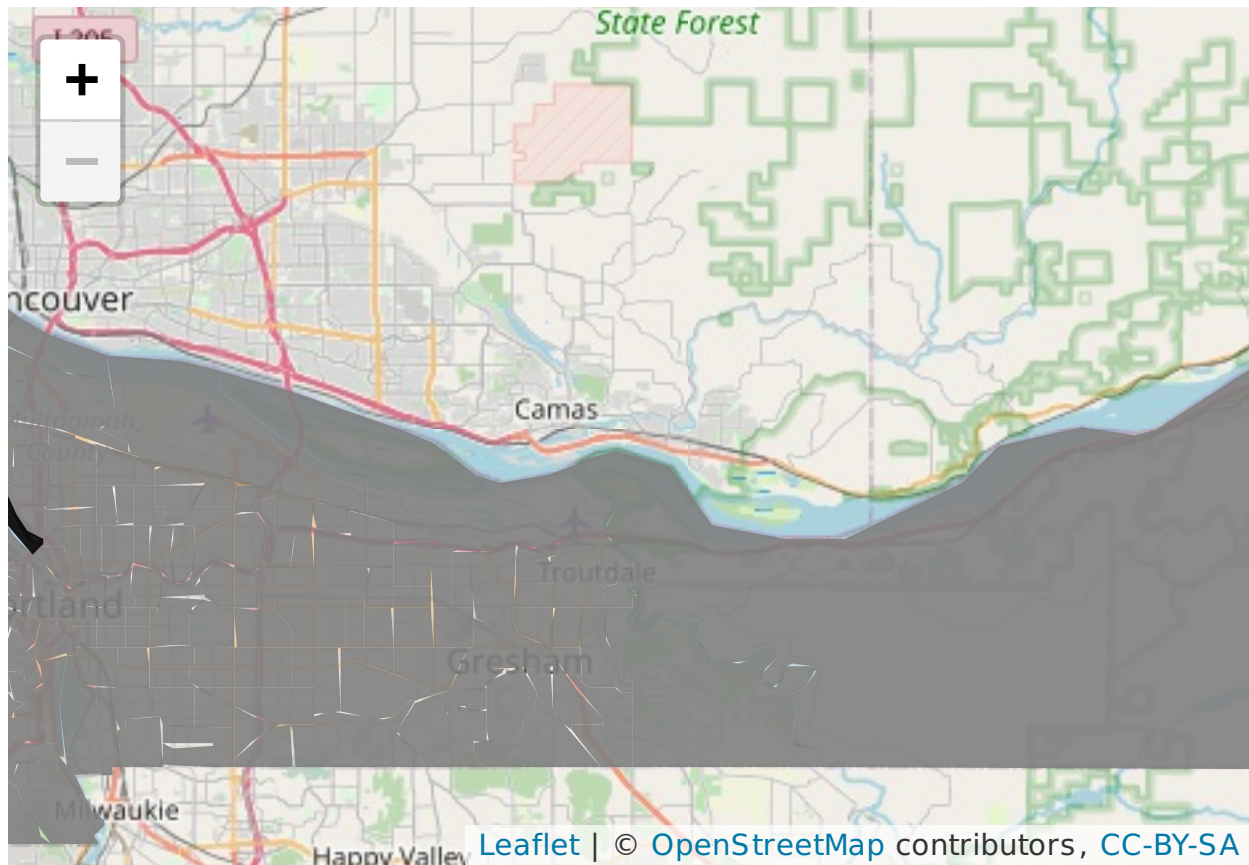
```
ggplot(data = block, mapping = aes(geometry = geometry)) +  
  geom_sf(aes(fill = estimate)) +  
  coord_sf() +  
  scale_fill_viridis_c() + theme_void()
```



We see that the eastern side of the multomah county has the lowest median gross rent or no data in block group resolution. This is probably because not many people live there. We see in the first map that the western places have higher median gross rent than the rest. In the tract resolution, around Portland seems to have the highest median gross rent. Except for that, the median gross rent is pretty variable across the county regions. Tract seems to be the most useful resolution for this variable because the county subdivision is pretty broad category and has just five divisions and a lot of data is missing (colored grey) in the block group resolution.

c. Making the tract map interactive:

```
tract %>%
  leaflet(options = leafletOptions(minZoom = 10, maxZoom = 15)) %>%
  addTiles() %>%
  addPolygons(popup = ~NAME, color = ~pal(estimate),
              stroke = FALSE, fillOpacity = 0.9)
```

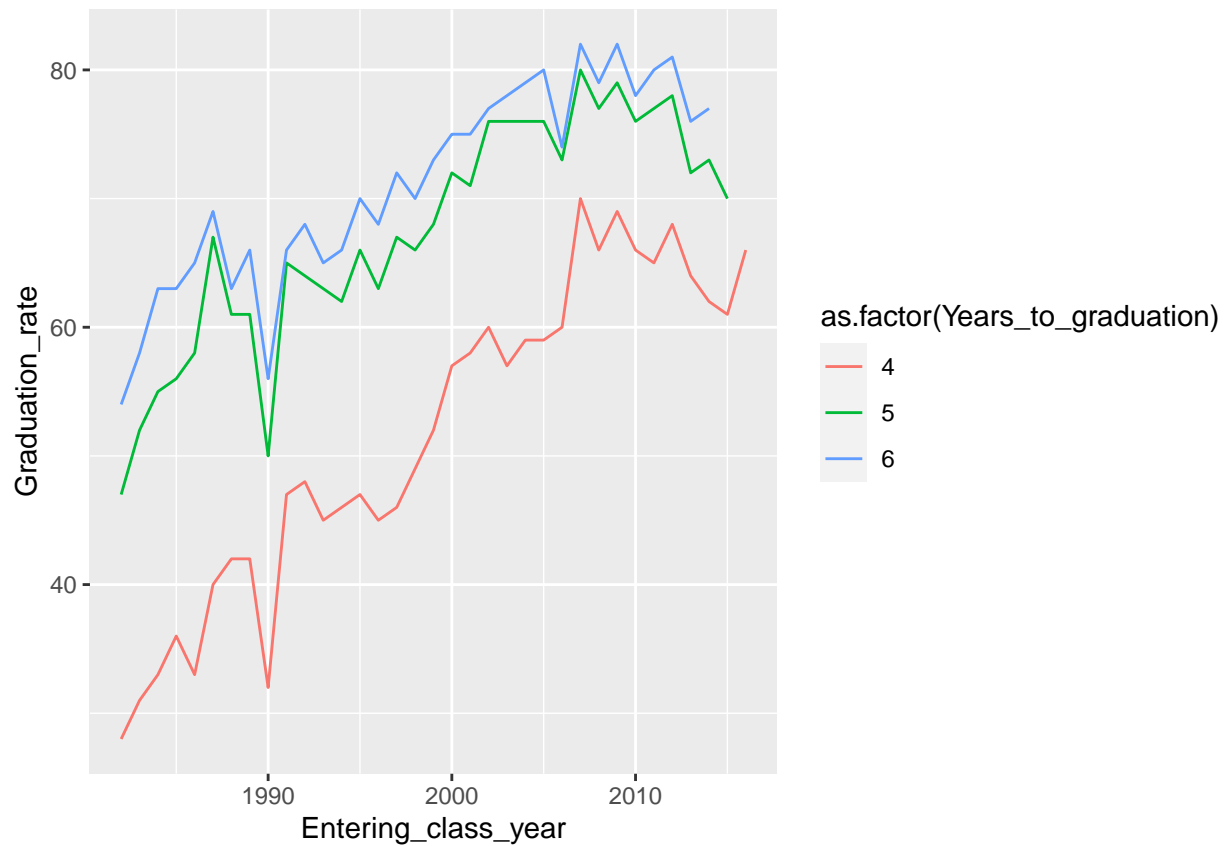


Problem 3: Take a Static Plot and Animate It!

a. Recreating graduation rate over year graph:

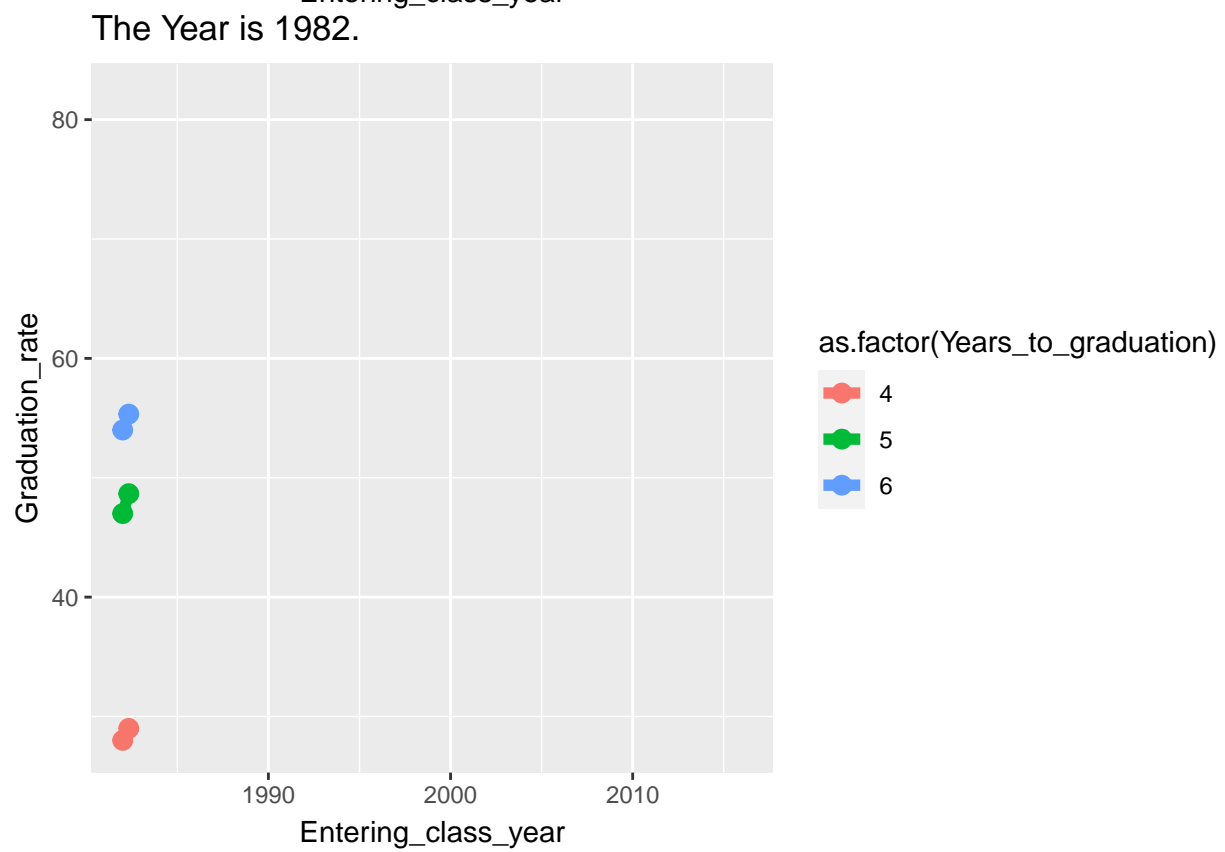
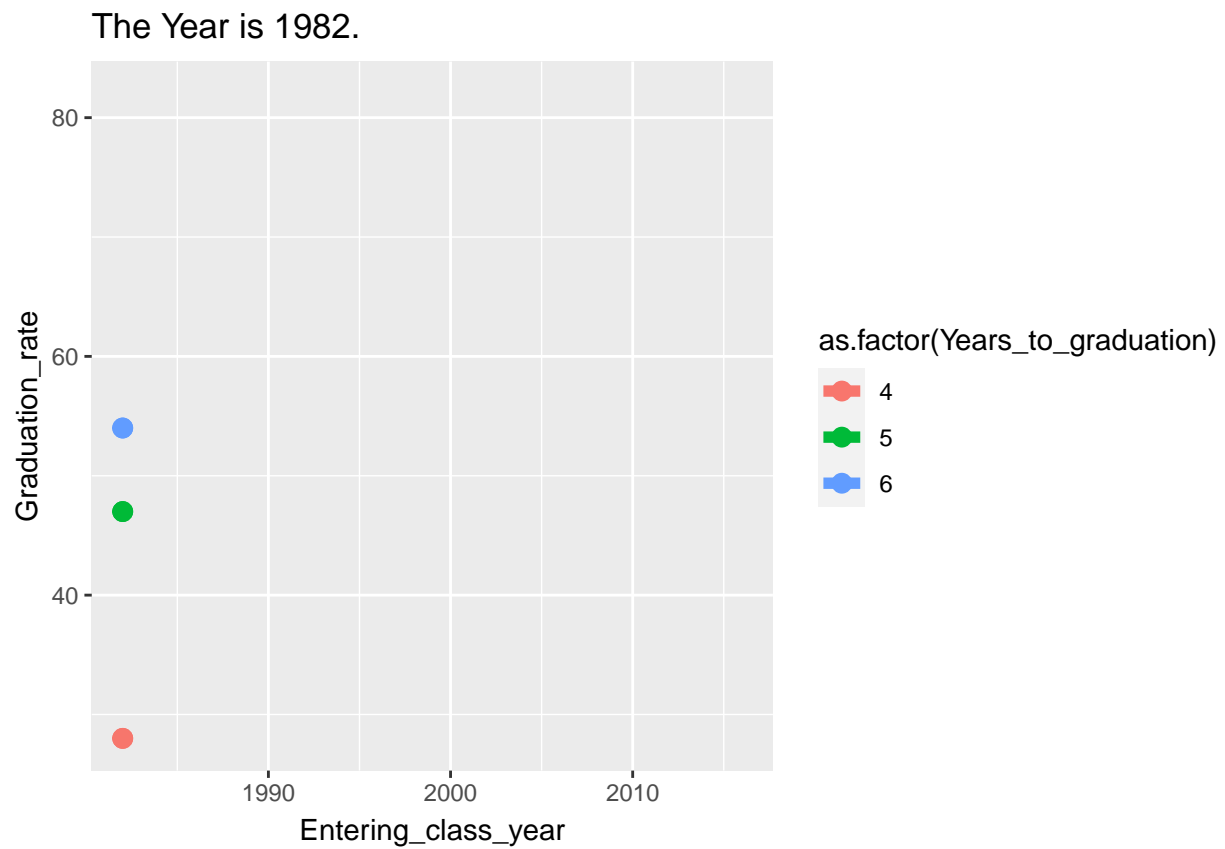
```
grad <- read_html("https://www.reed.edu/ir/gradrateshist.html")%>%
  html_nodes("table") %>%
  html_table(fill = TRUE)
grad1 <- grad[[1]]
colnames(grad1) <- c("Entering class year", "Number in cohort", "4", "5", "6")
grad1 <- grad1 %>%
  filter(row_number() > 1)
grad1 <- pivot_longer(grad1, cols = c("4", "5", "6"),
                      names_to = "Years to graduation",
                      values_to = "Graduation rate") %>%
  mutate(Entering_class_year = as.integer('Entering class year'),
         Cohort_size = as.integer('Number in cohort'),
         Years_to_graduation = as.integer('Years to graduation'),
         Graduation_rate = parse_number('Graduation rate')) %>%
  select(Entering_class_year, Cohort_size, Years_to_graduation, Graduation_rate)
```

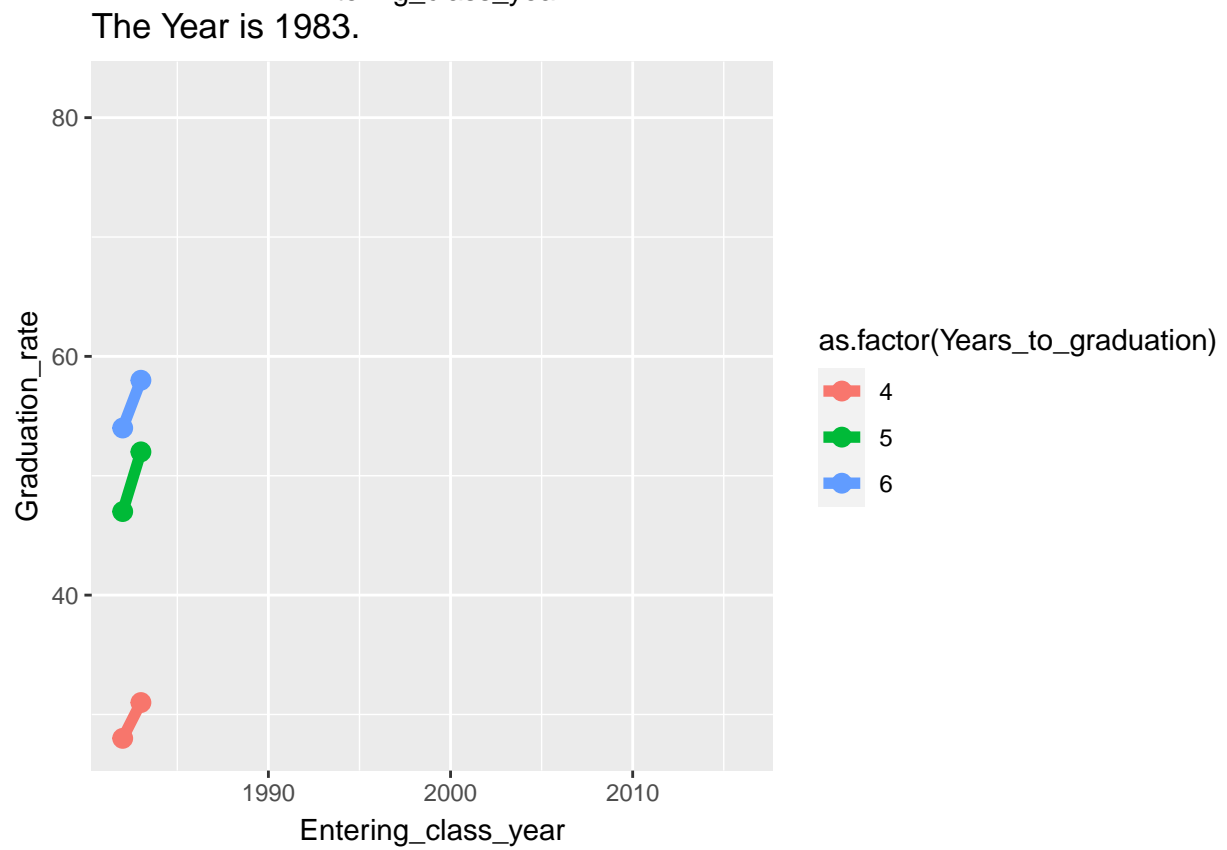
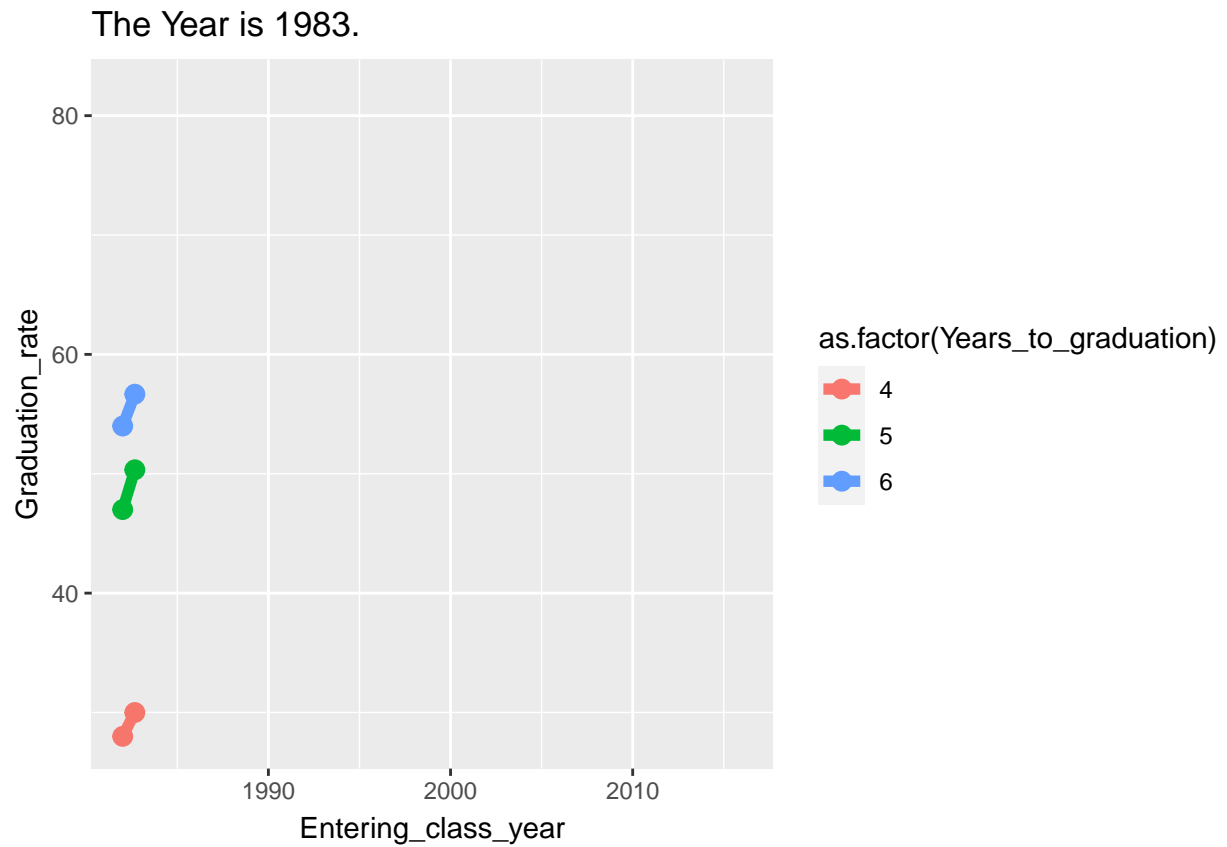
```
plot <- grad1 %>% ggplot(aes(x = Entering_class_year,
  y = Graduation_rate,
  color = as.factor(Years_to_graduation))) + geom_line()
plot
```



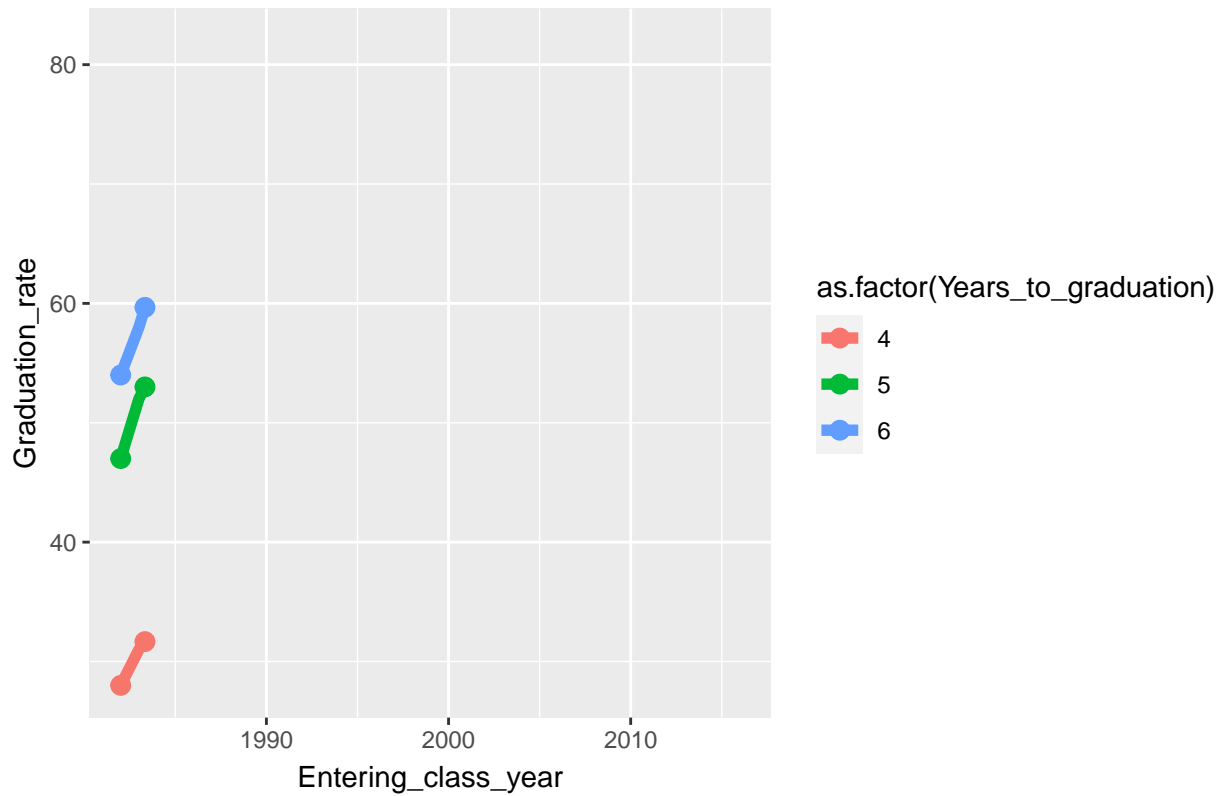
b. With animation:

```
library(gganimate)
animate <- grad1 %>% ggplot(aes(x = Entering_class_year,
  y = Graduation_rate,
  color = as.factor(Years_to_graduation))) + geom_line(size=2)+geom_point(size=3) +
  transition_reveal(Entering_class_year)+
  labs(title = "The Year is {round(frame_along, 0)}.")
animate
```

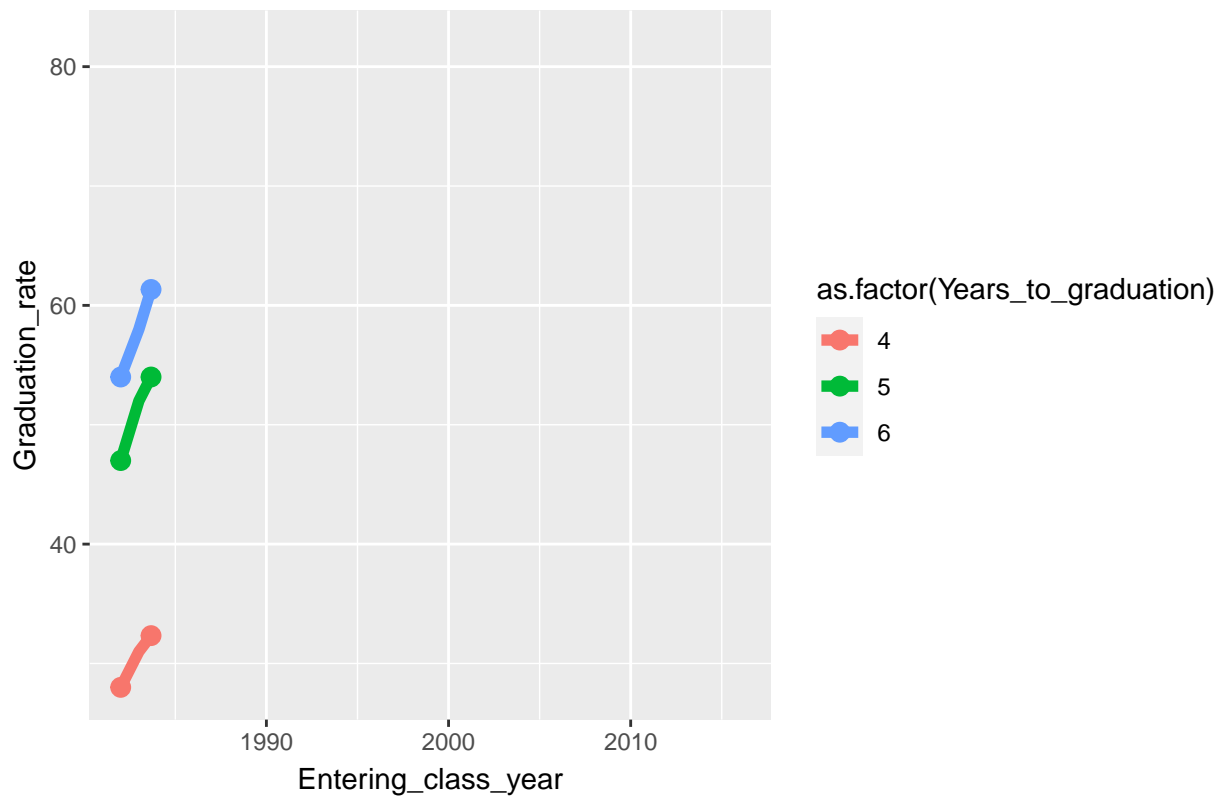


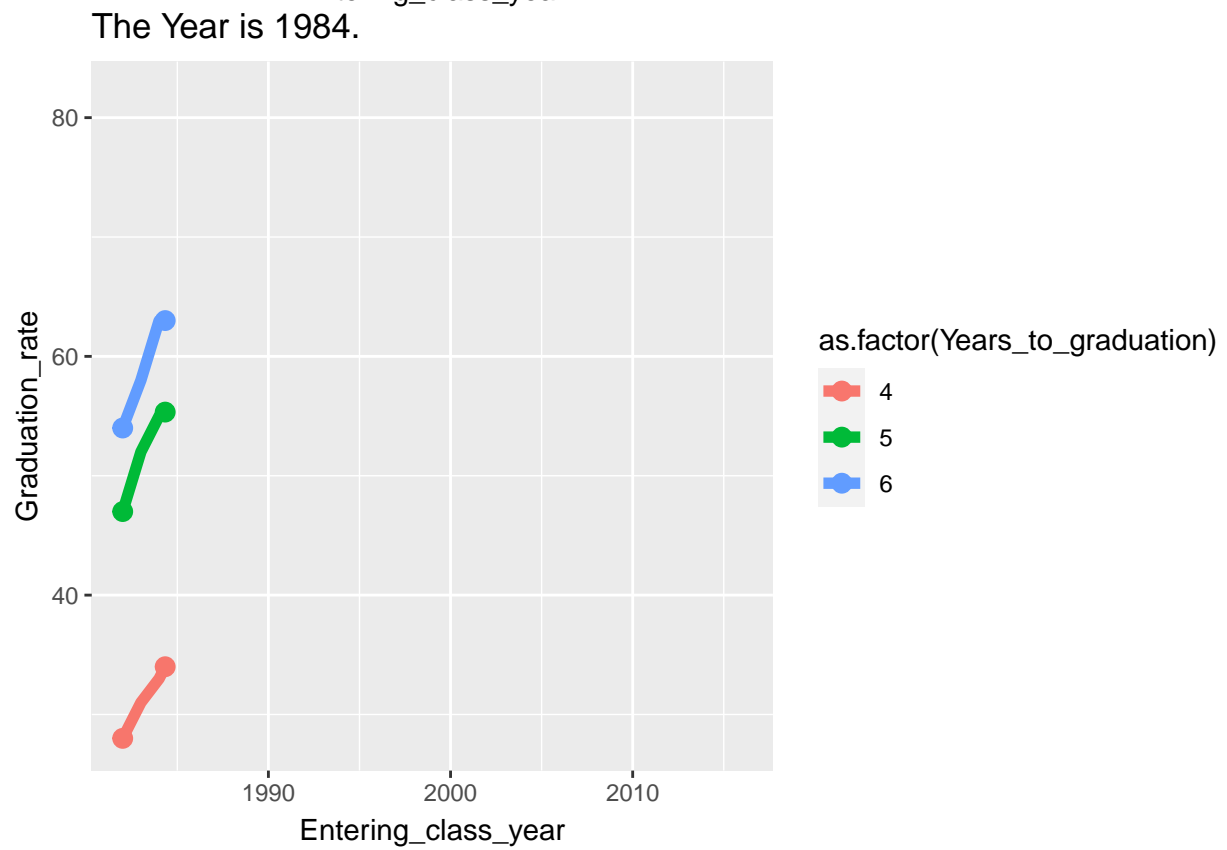
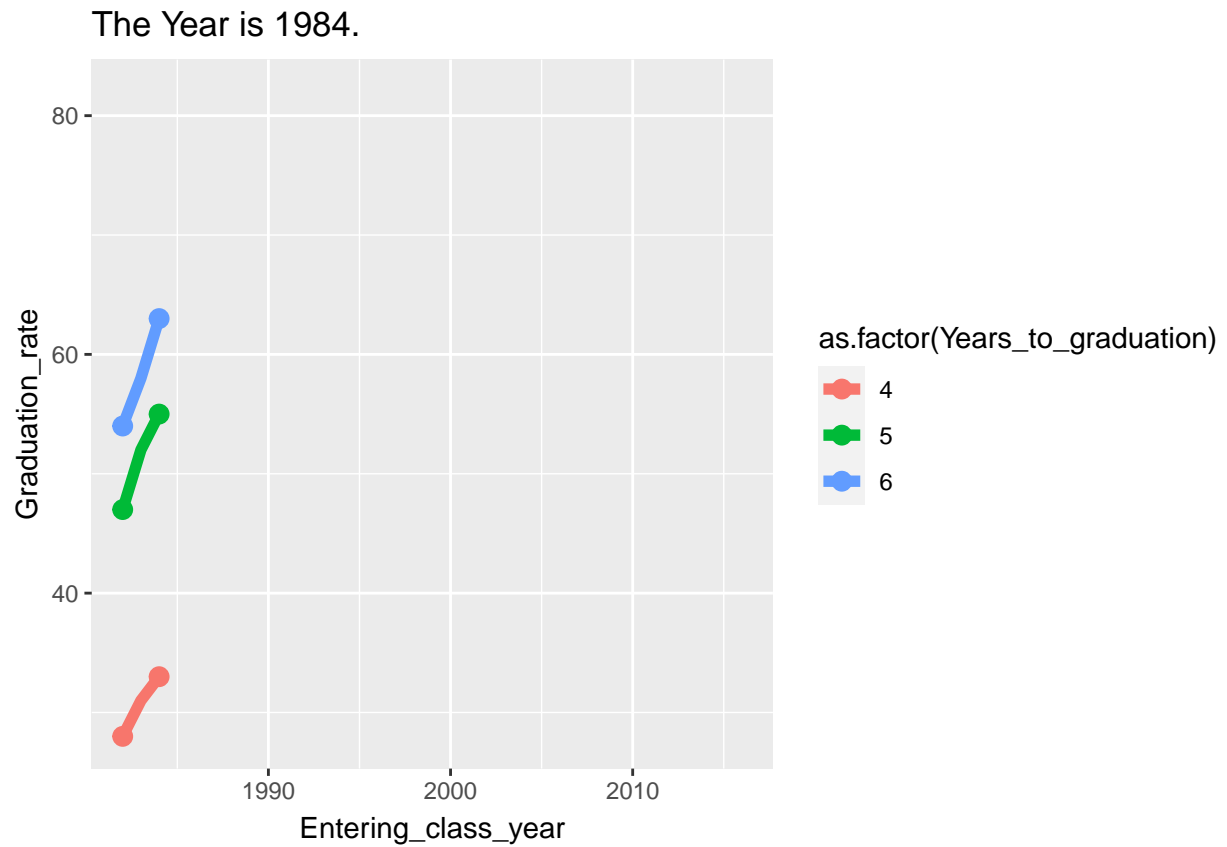


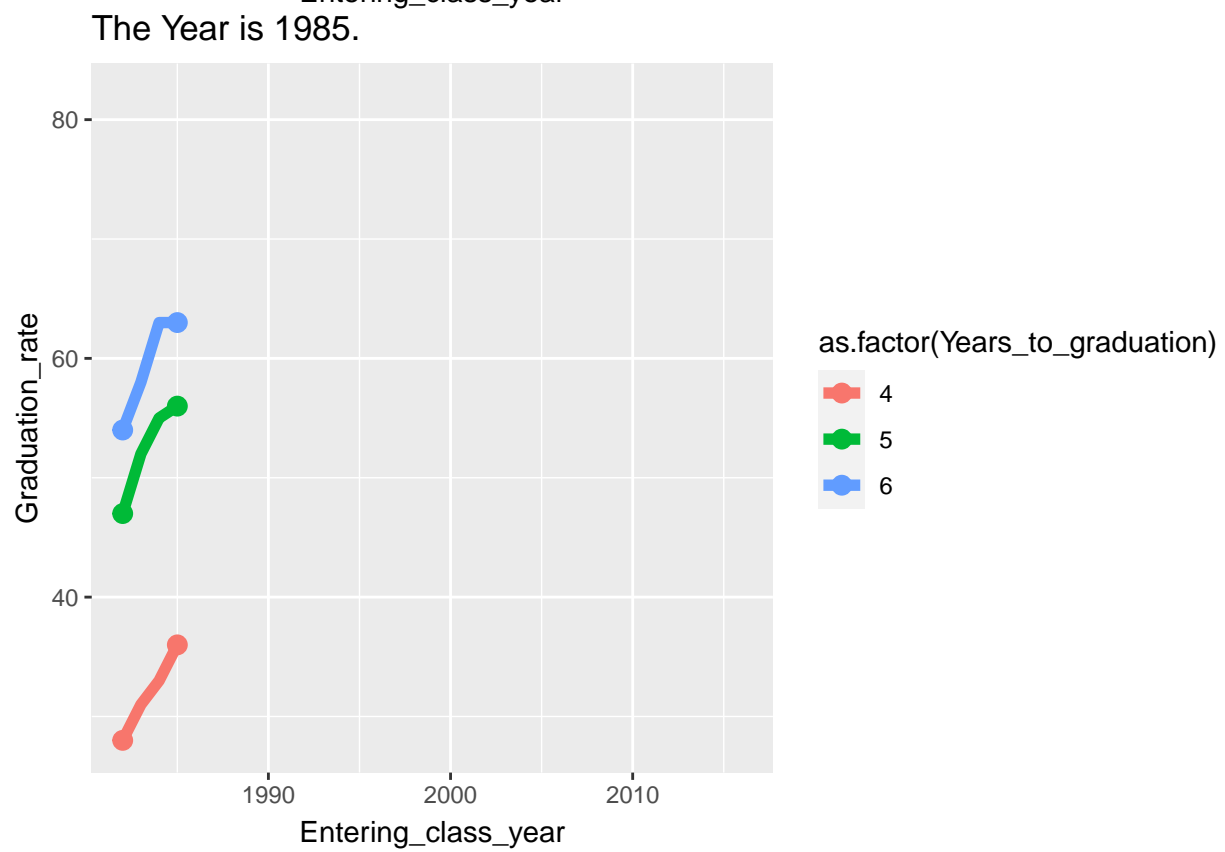
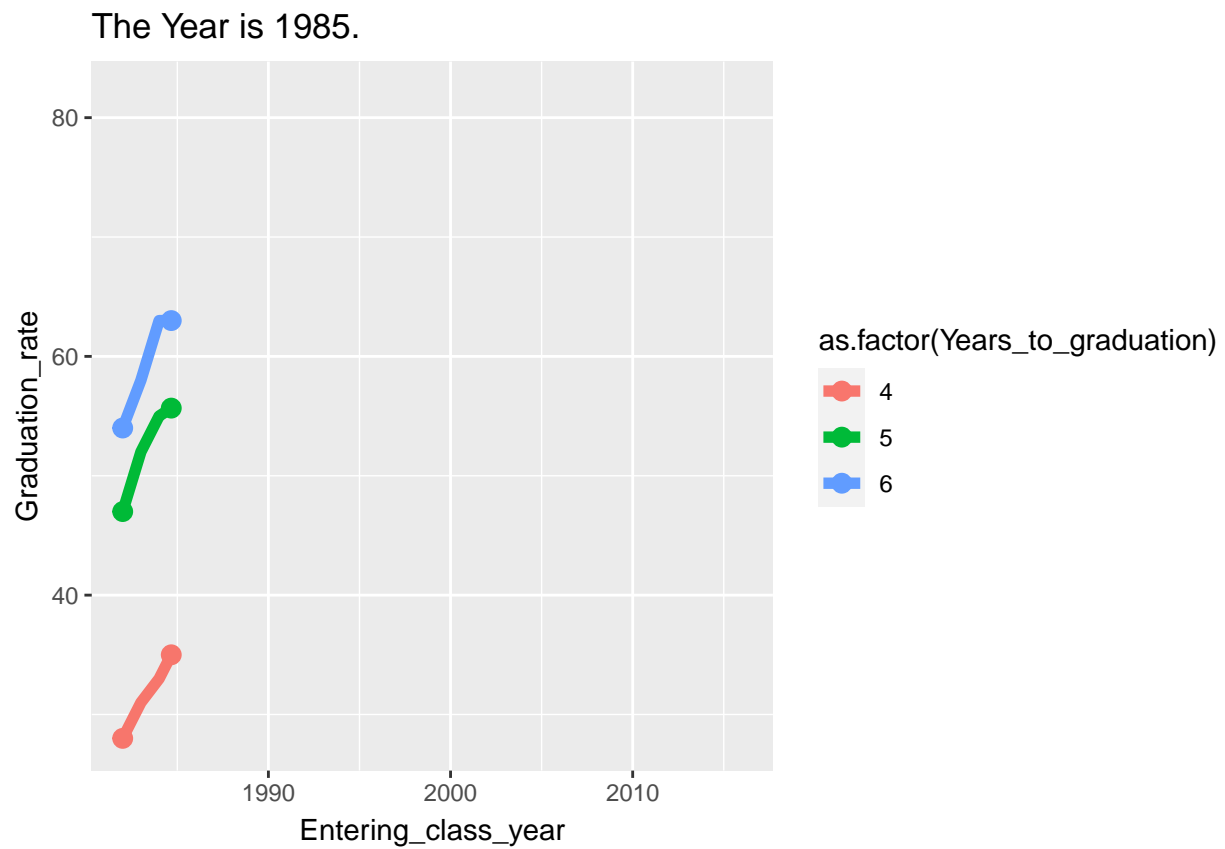
The Year is 1983.



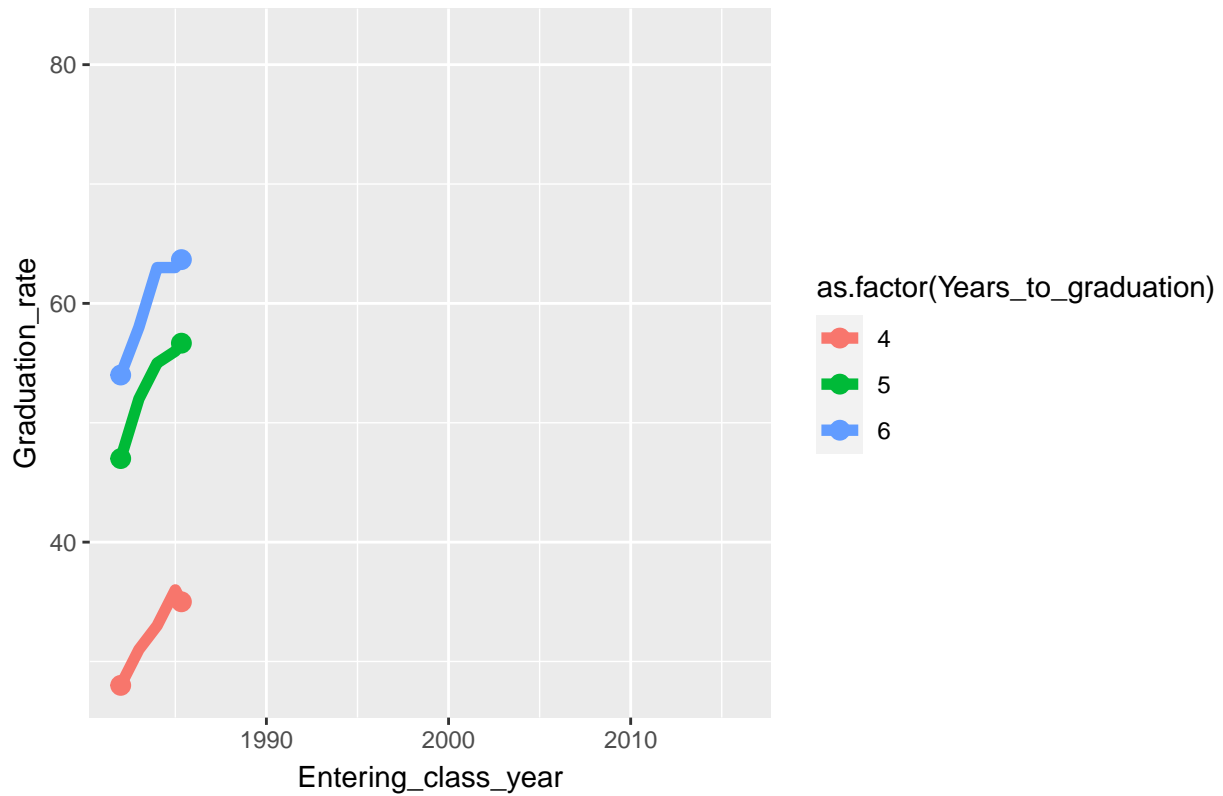
The Year is 1984.







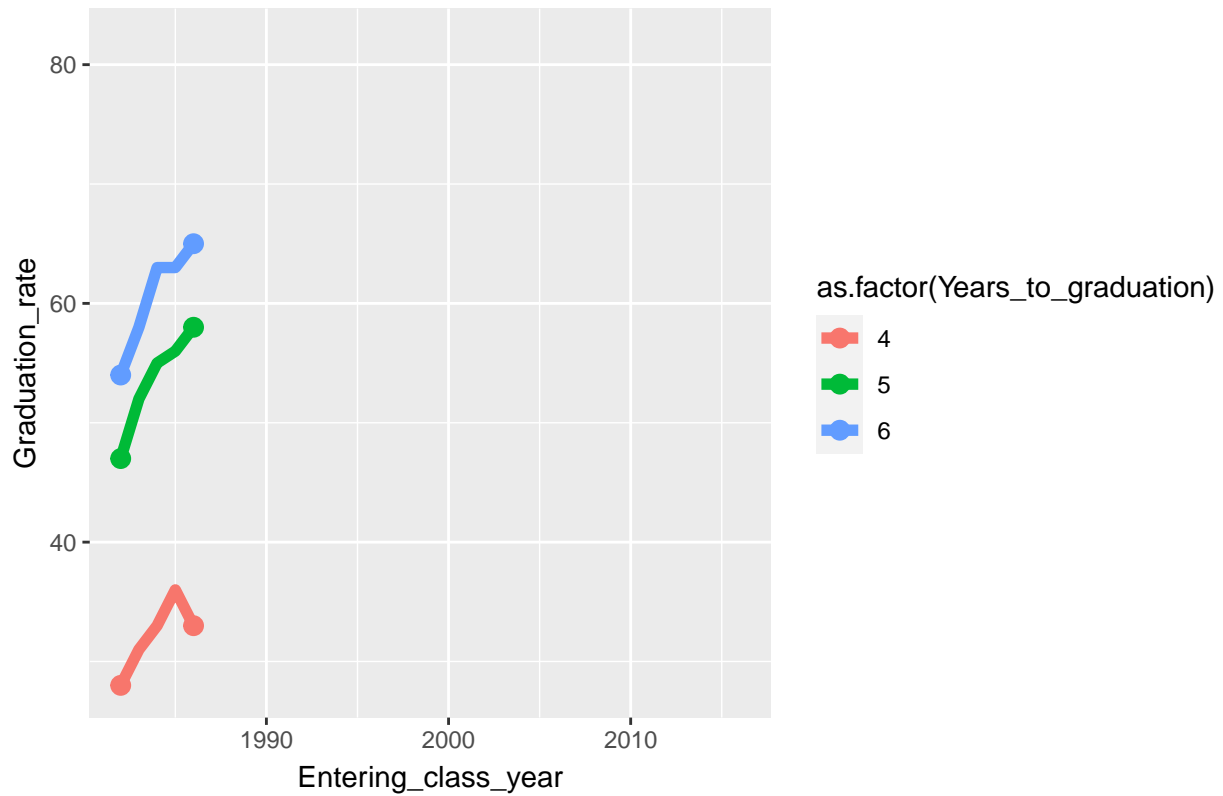
The Year is 1985.



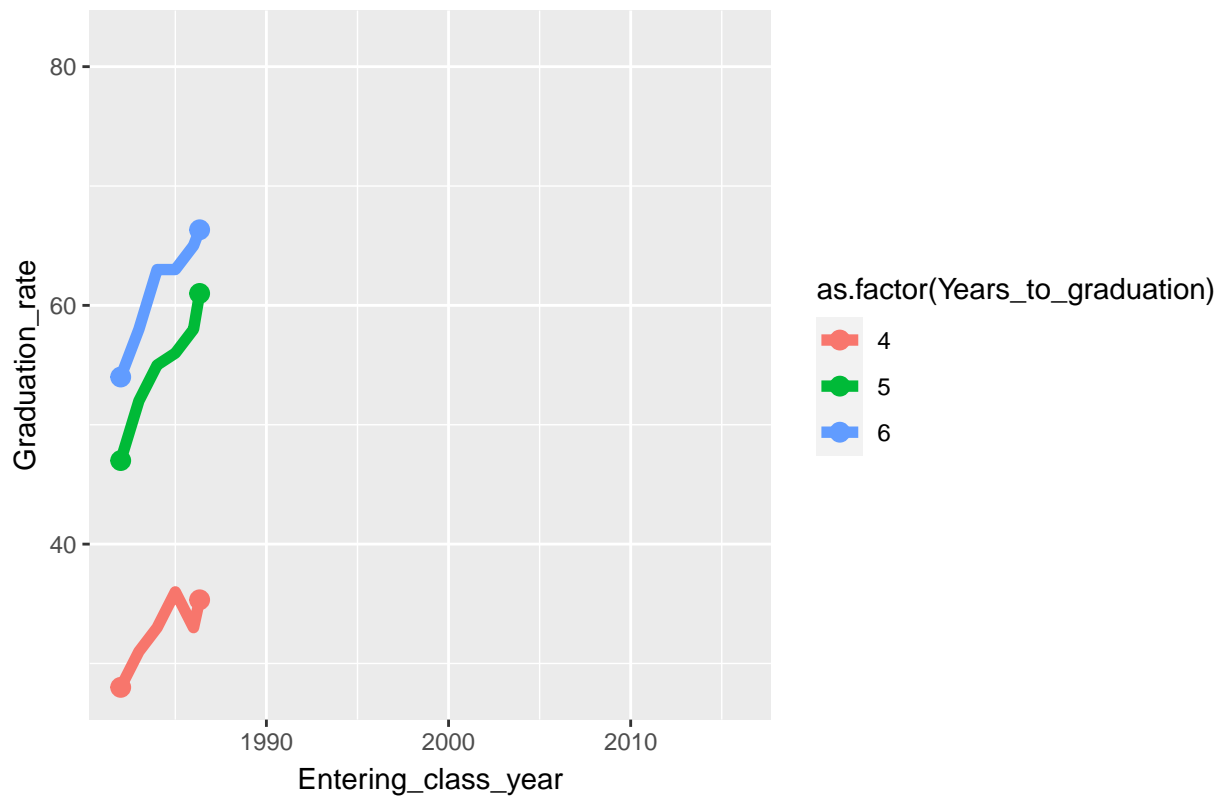
The Year is 1986.

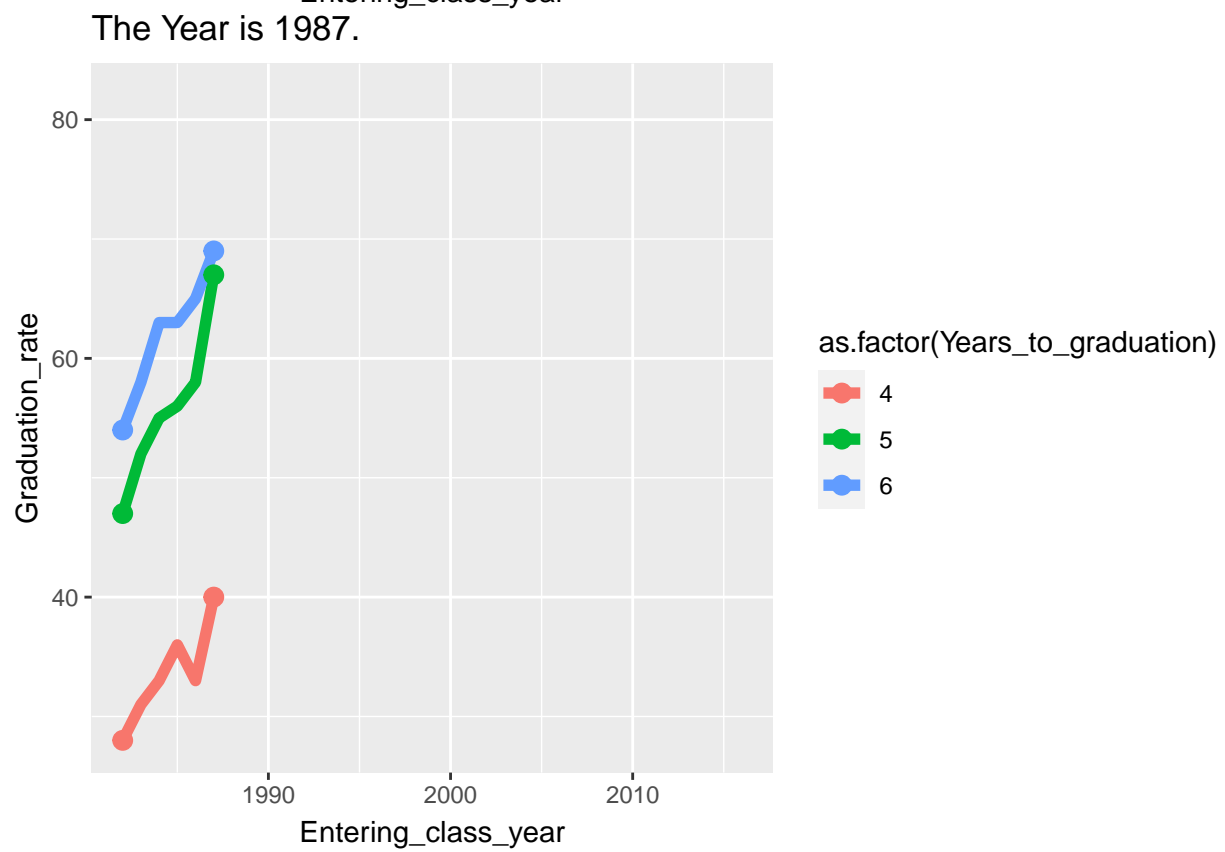
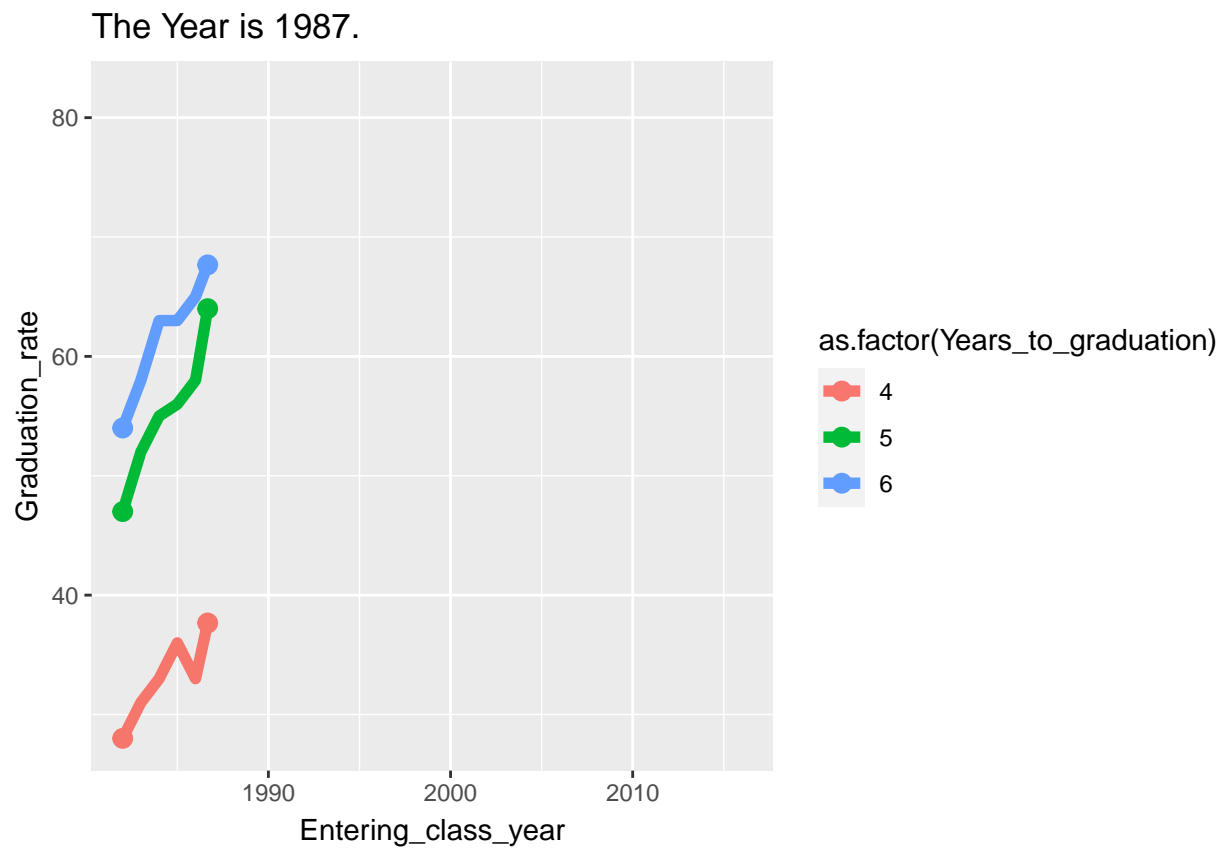


The Year is 1986.

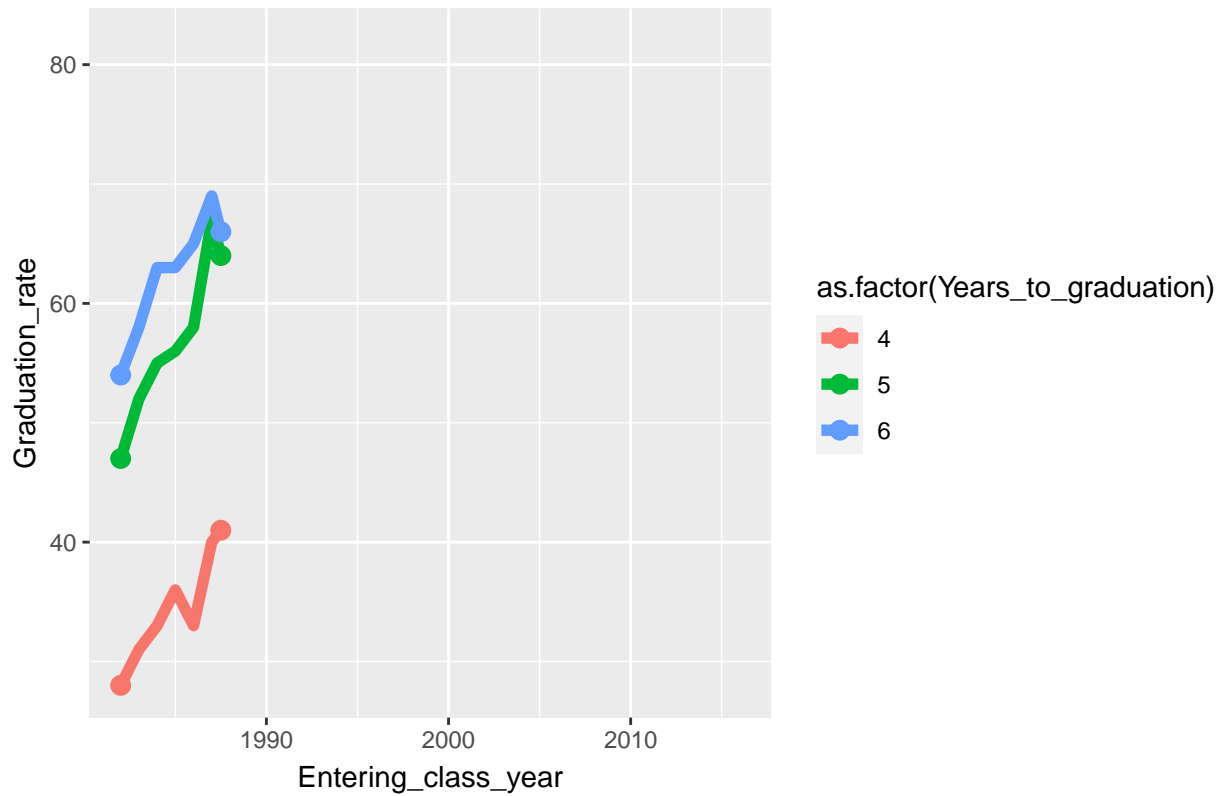


The Year is 1986.





The Year is 1987.



The Year is 1988.



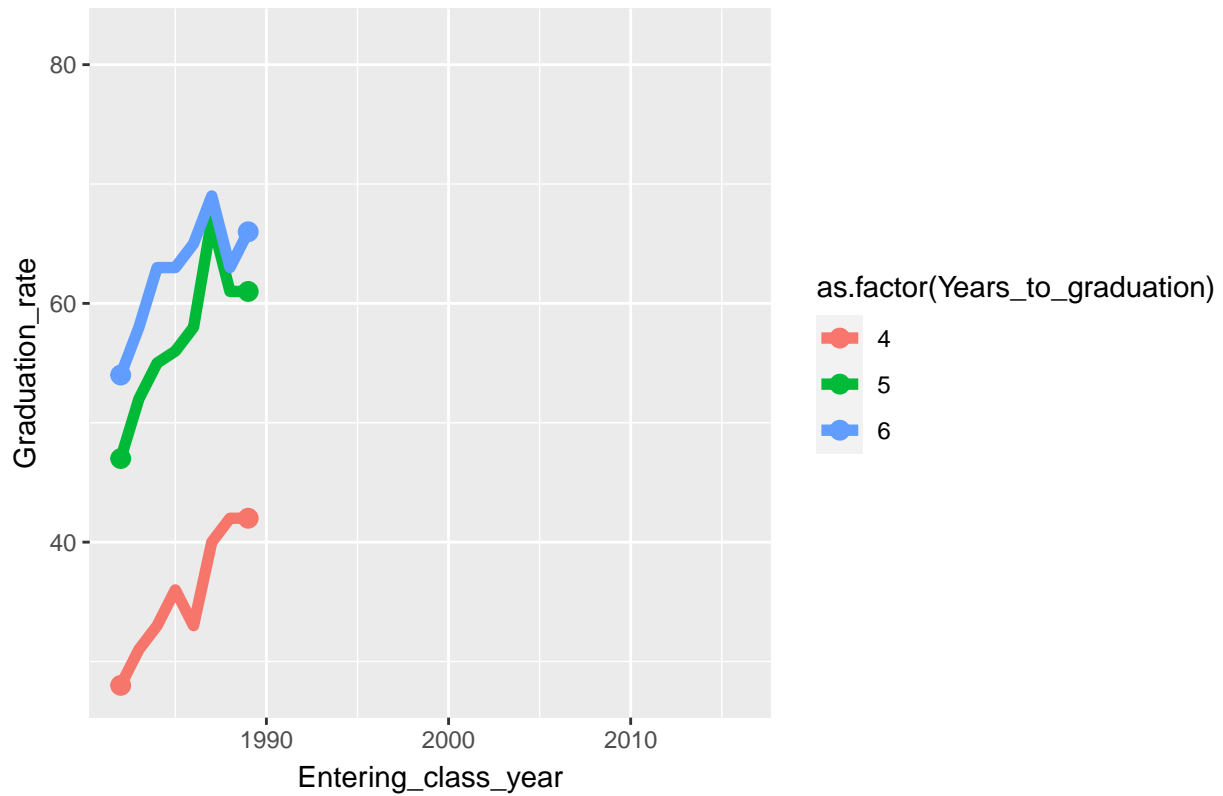
The Year is 1988.



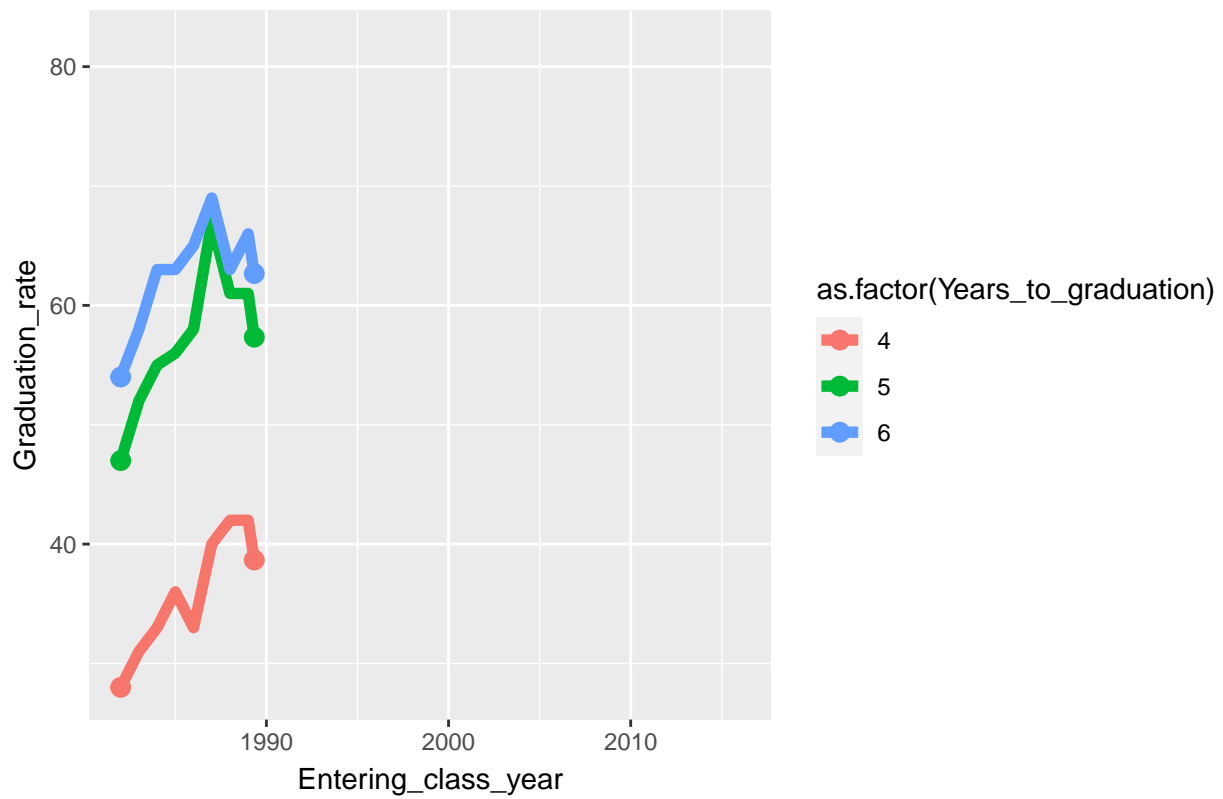
The Year is 1989.

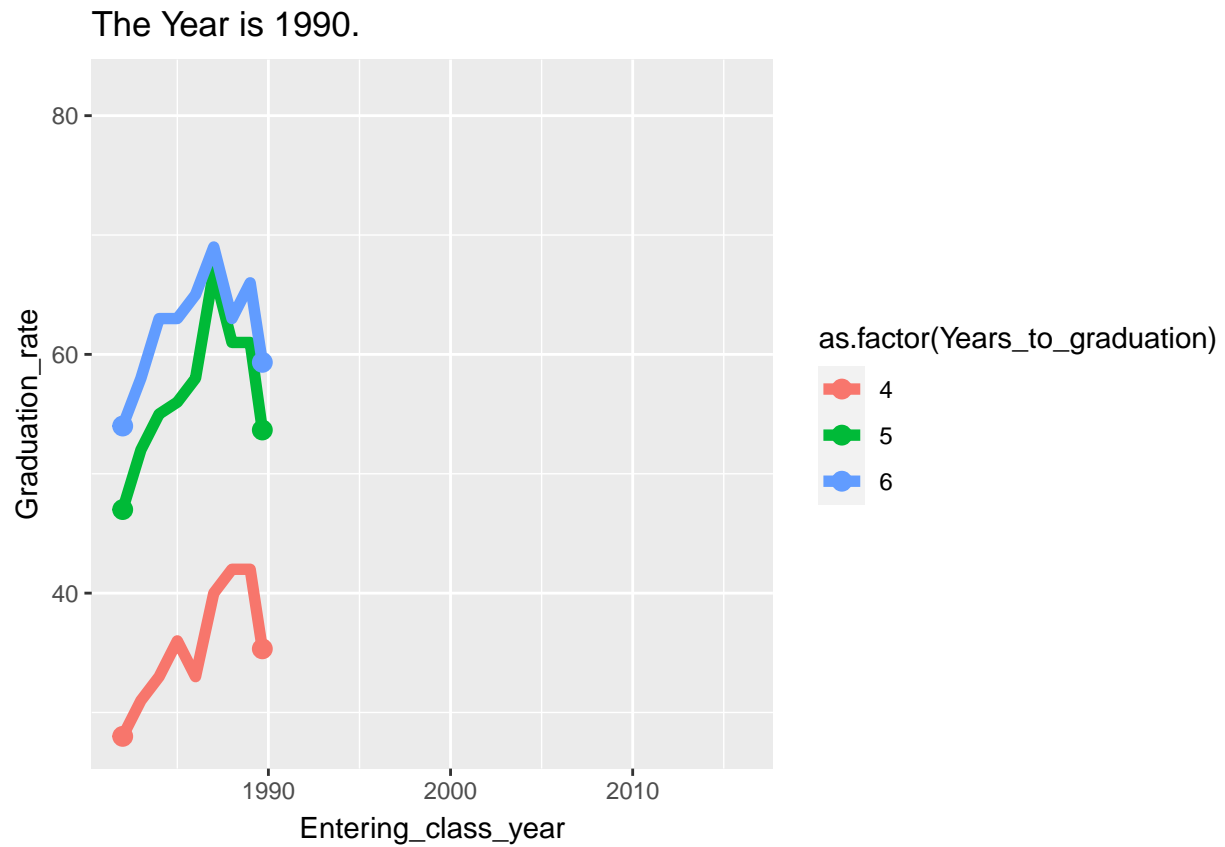


The Year is 1989.



The Year is 1989.



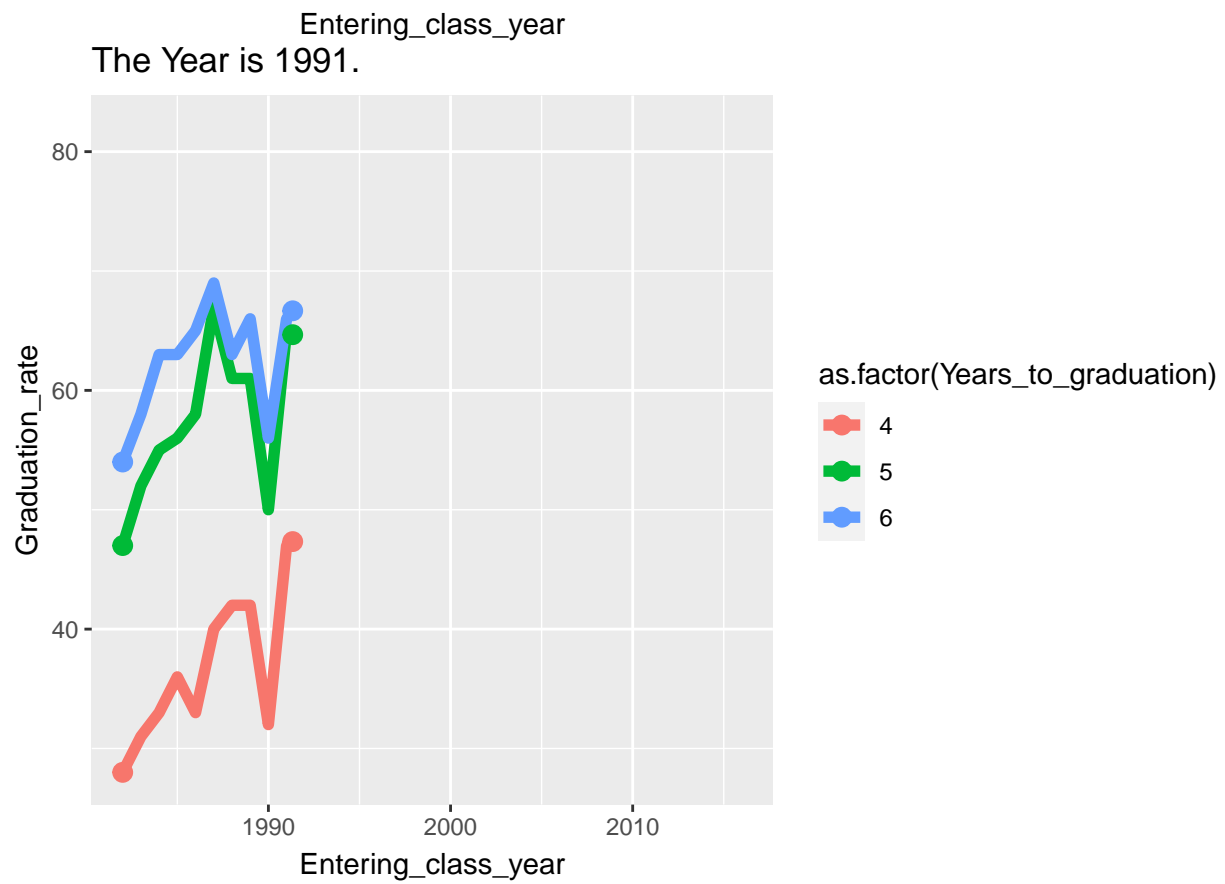
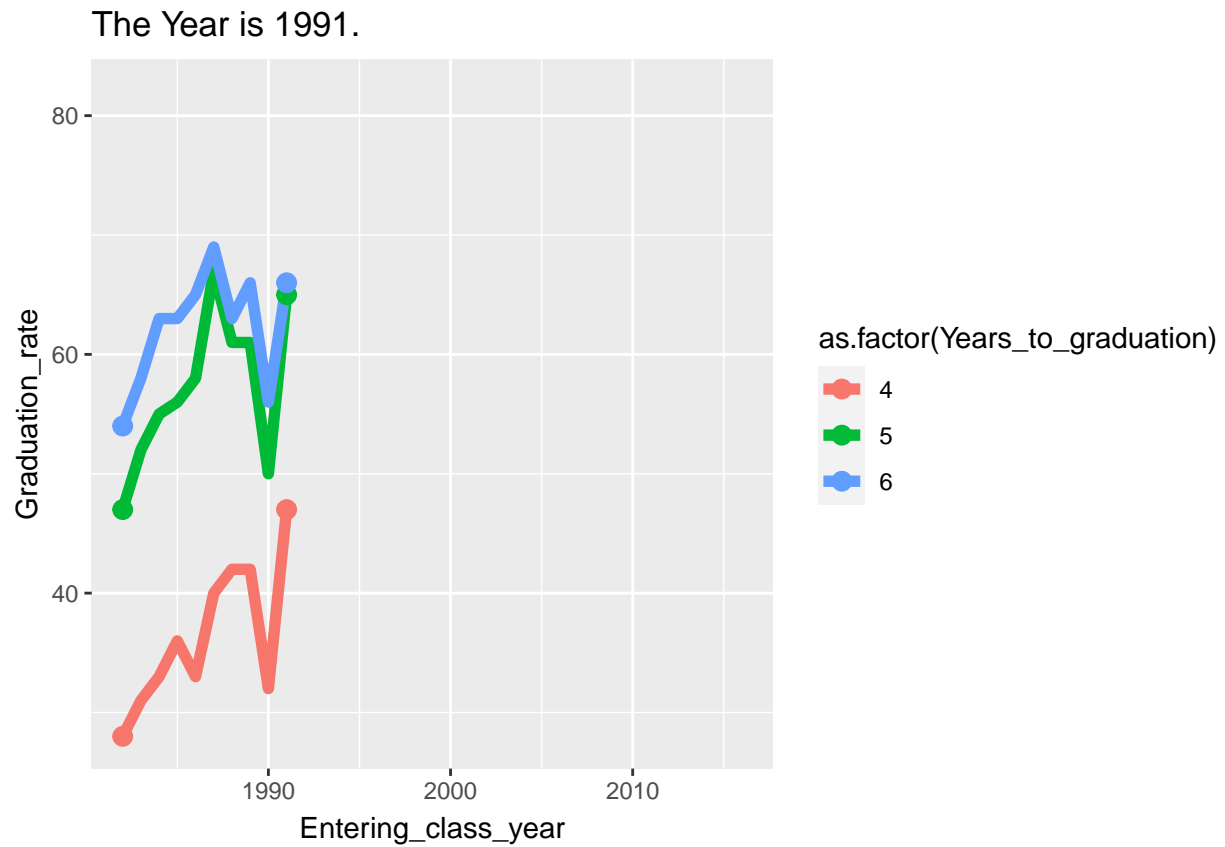


The Year is 1990.

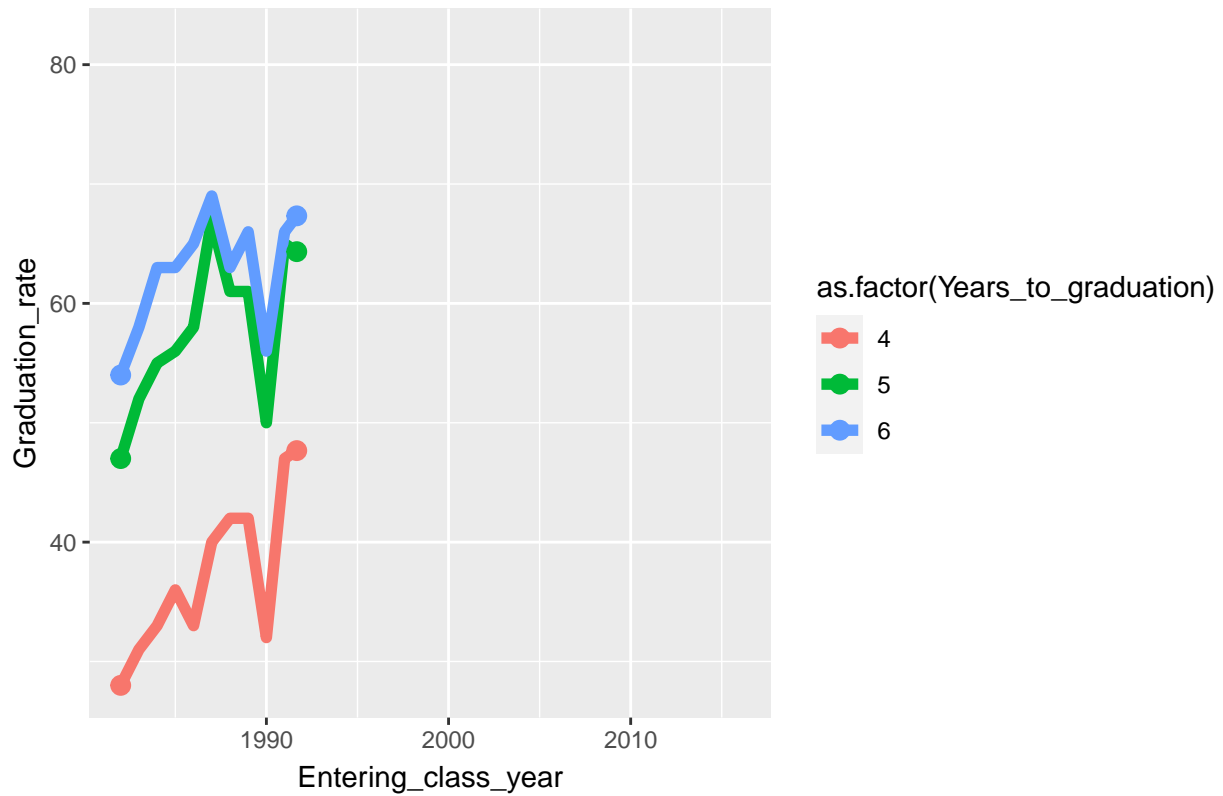


The Year is 1991.





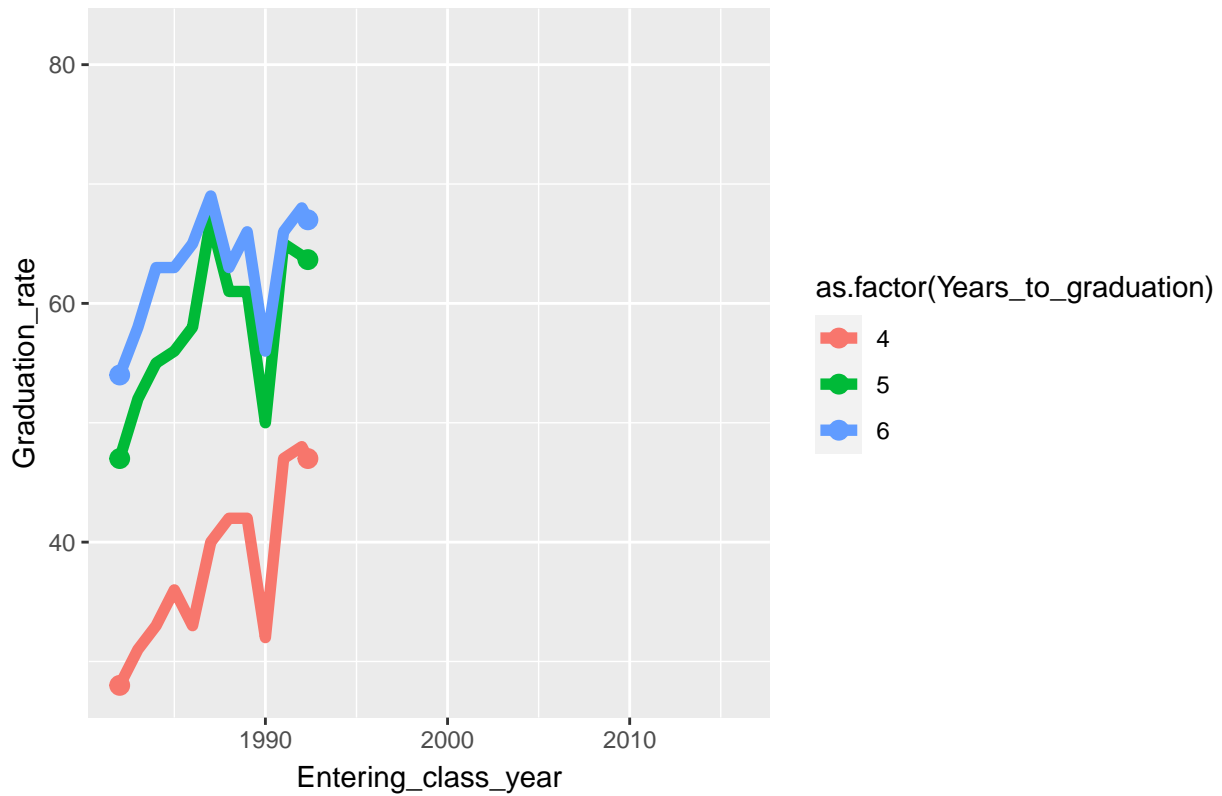
The Year is 1992.



The Year is 1992.



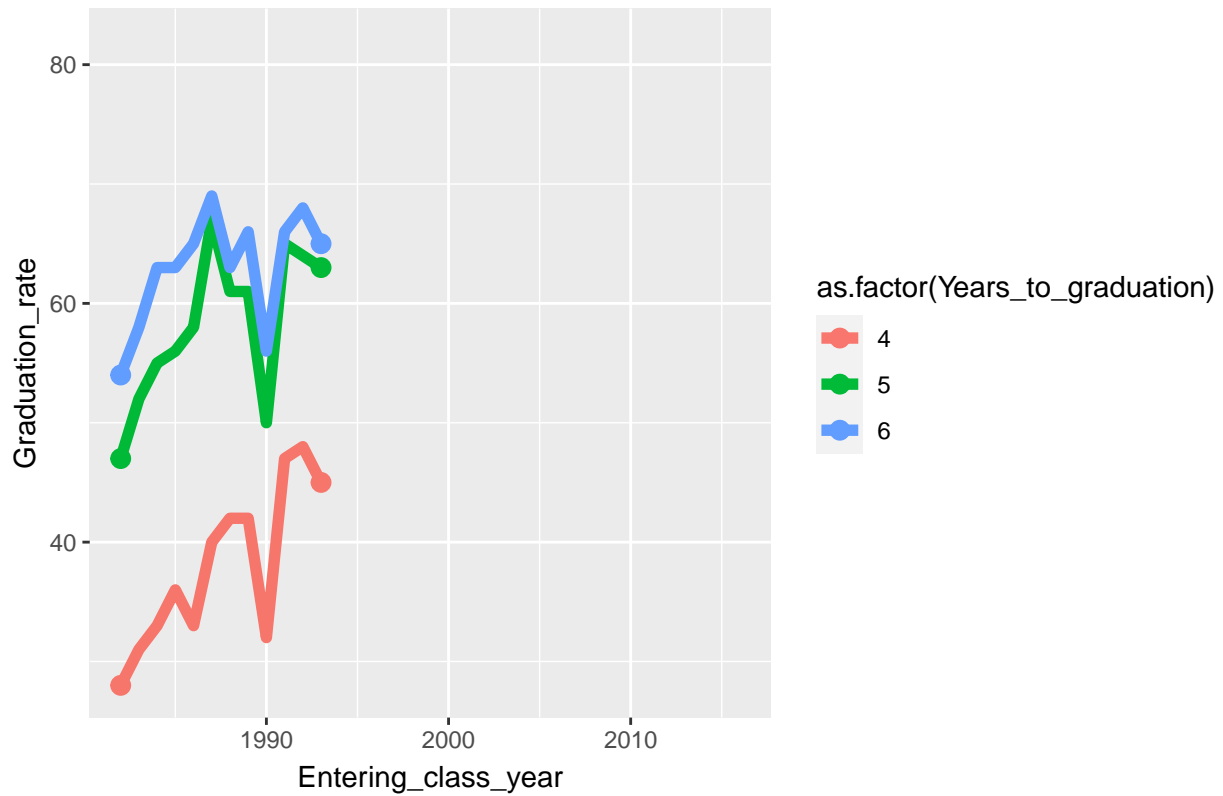
The Year is 1992.



The Year is 1993.



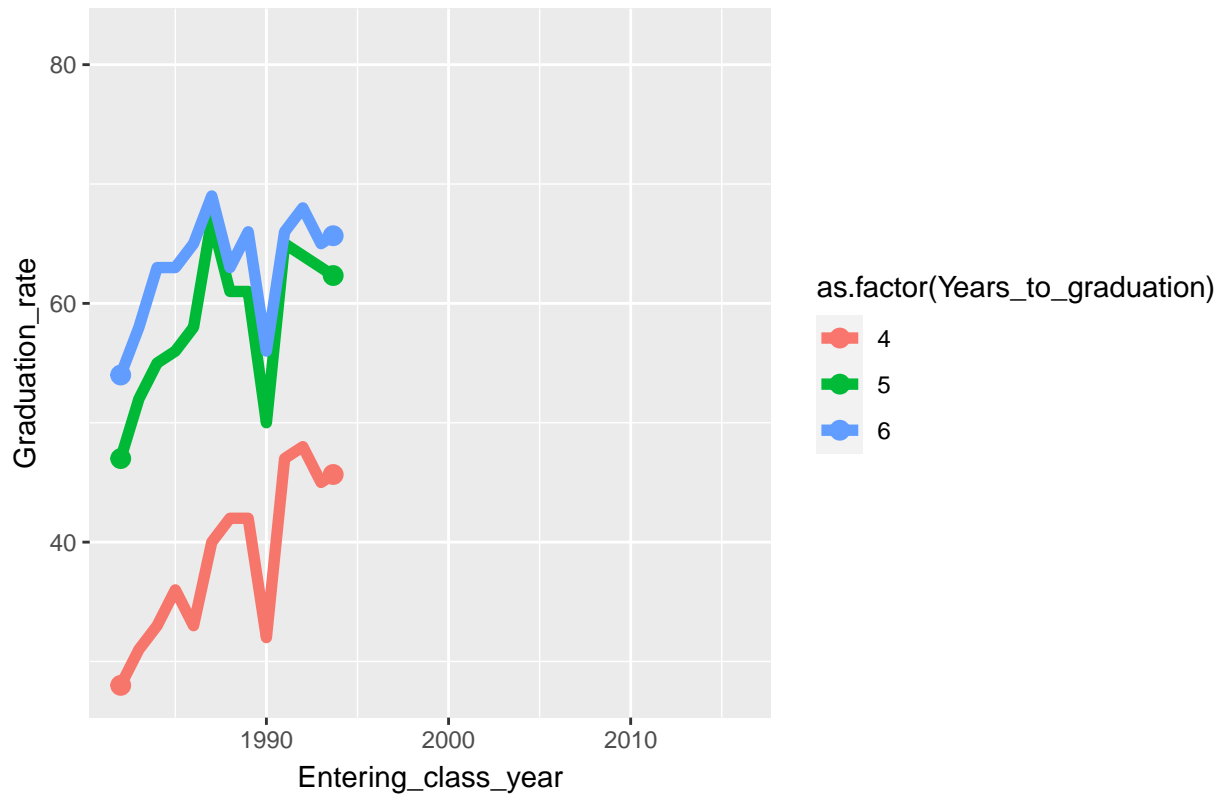
The Year is 1993.



The Year is 1993.



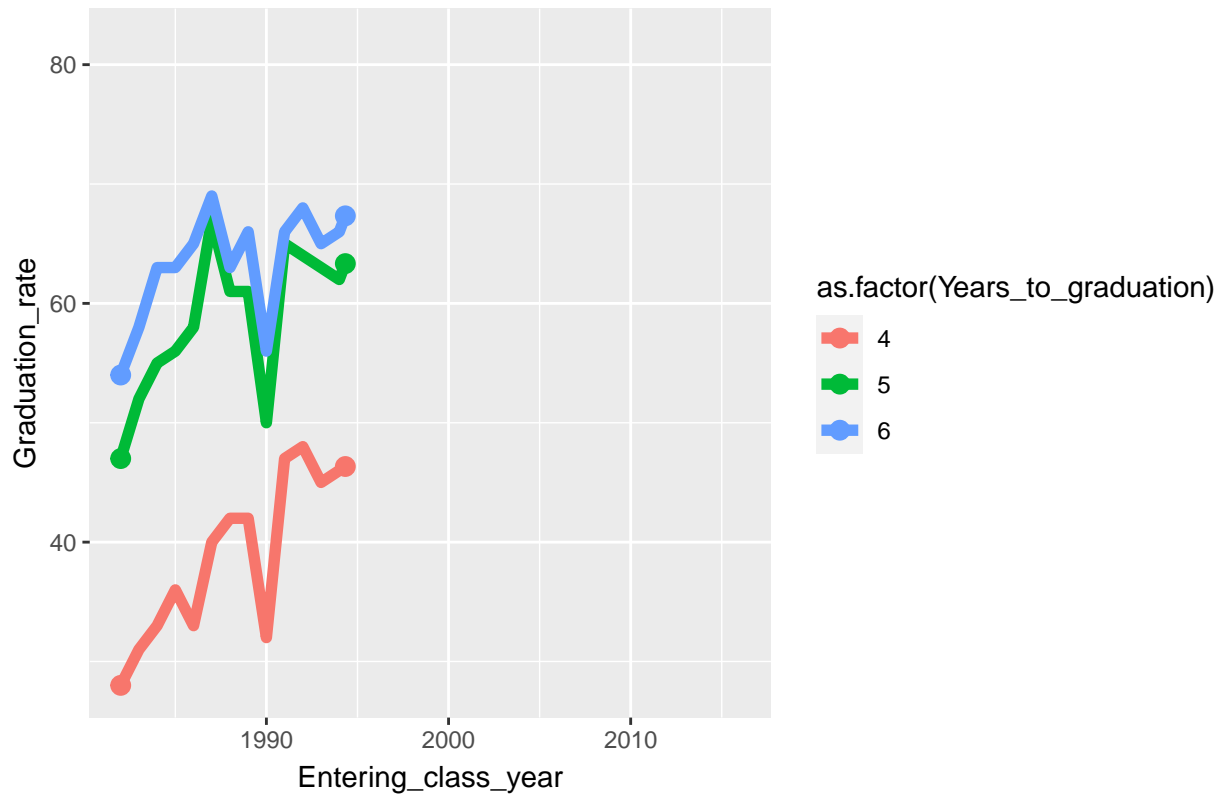
The Year is 1994.



The Year is 1994.



The Year is 1994.



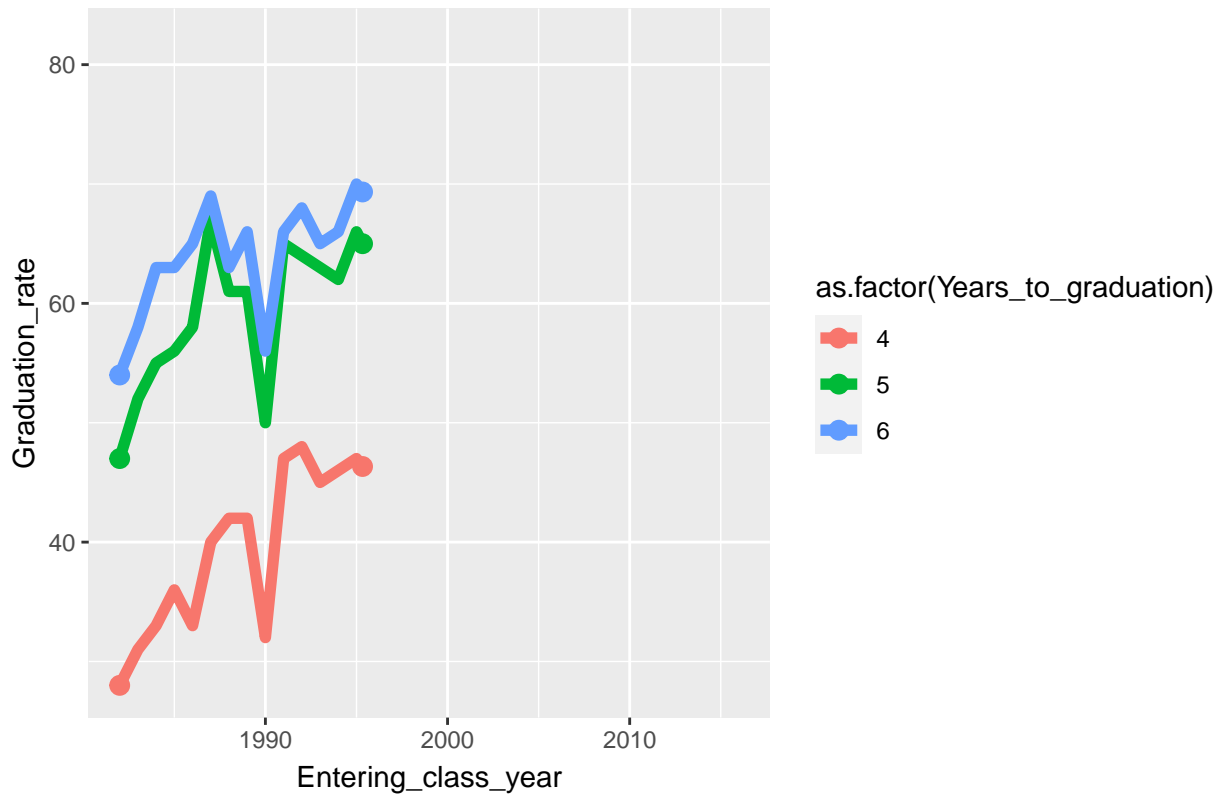
The Year is 1995.



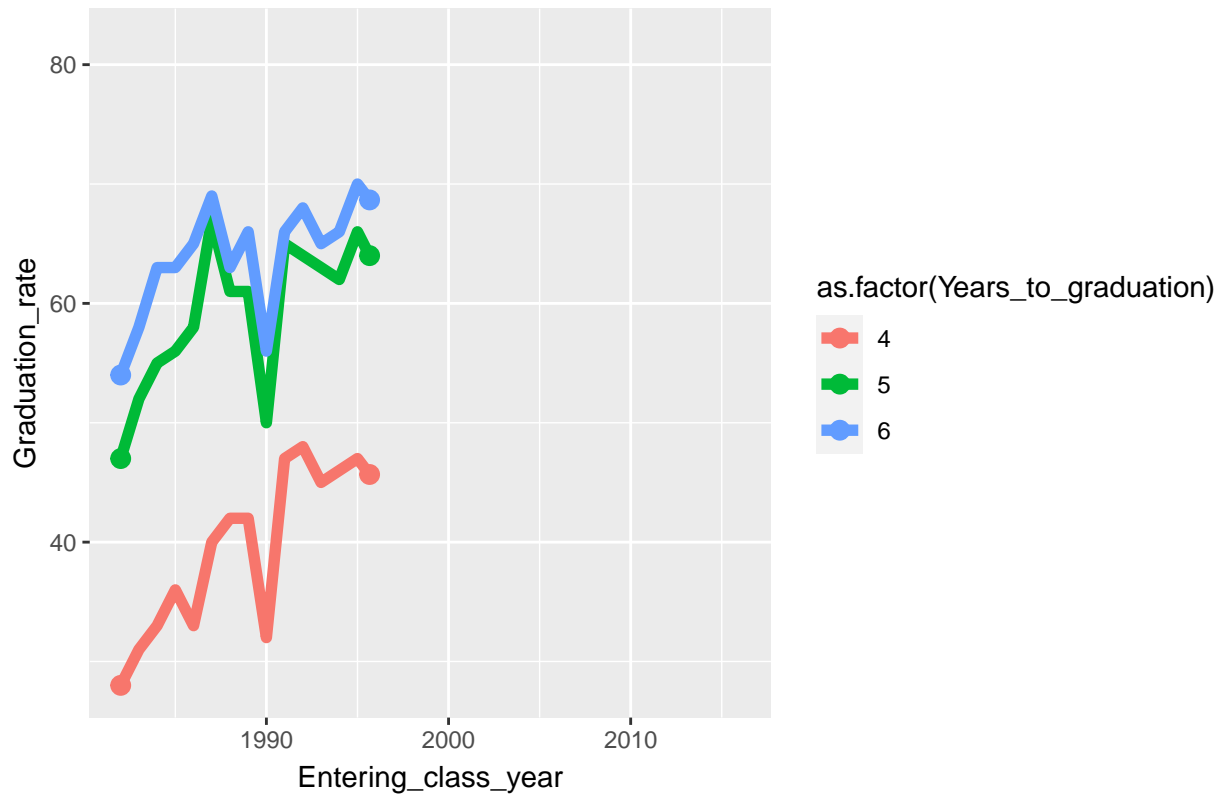
The Year is 1995.



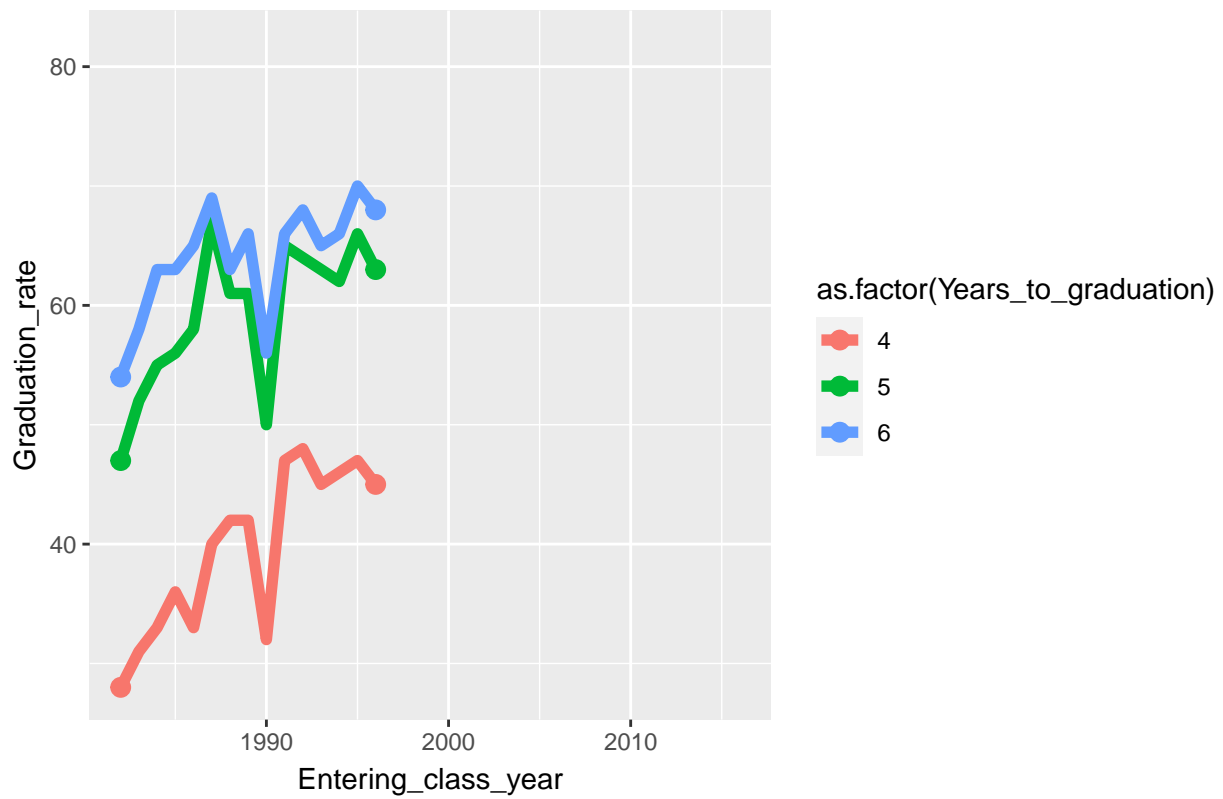
The Year is 1995.



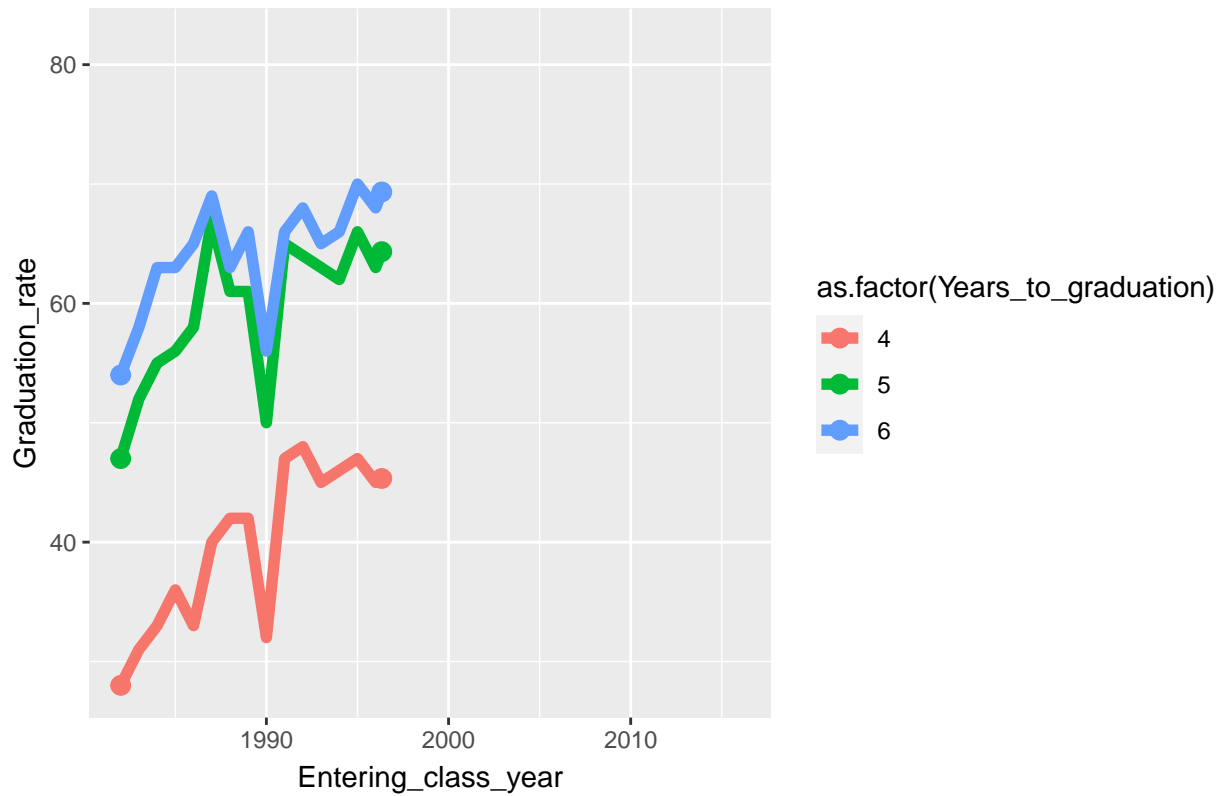
The Year is 1996.



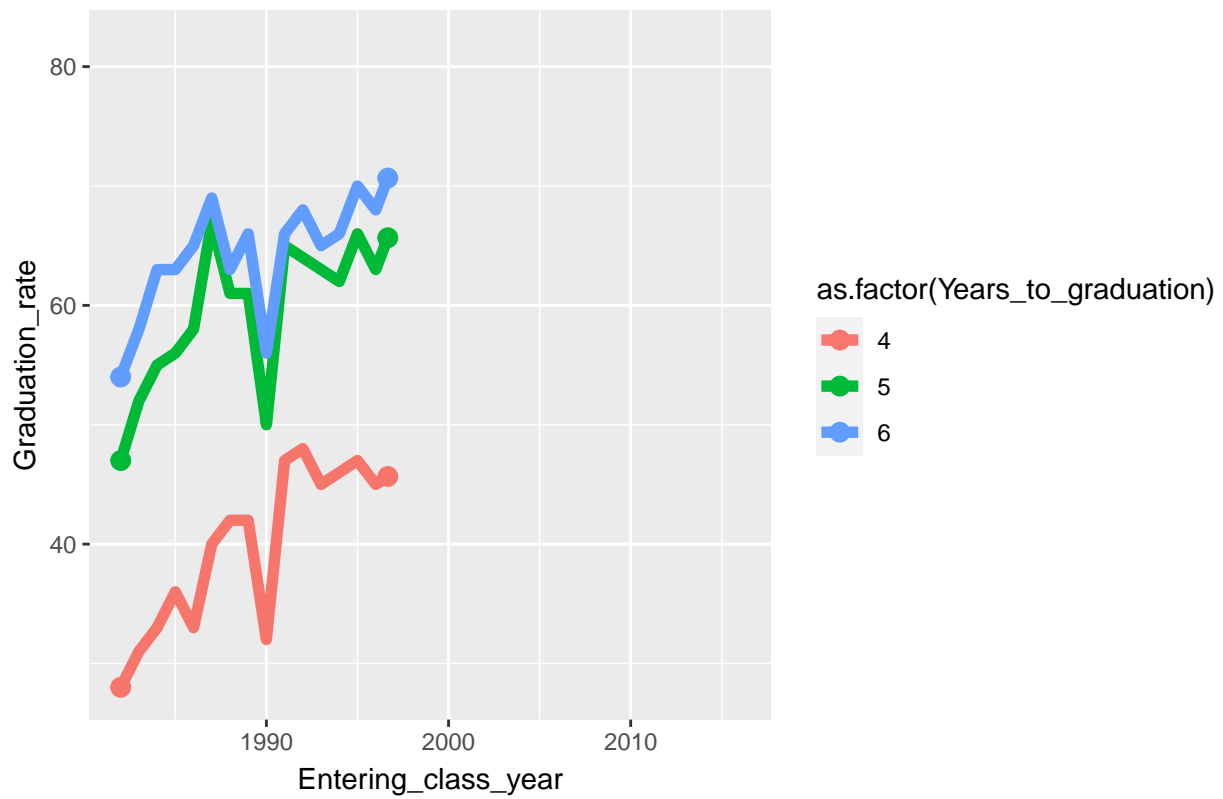
The Year is 1996.



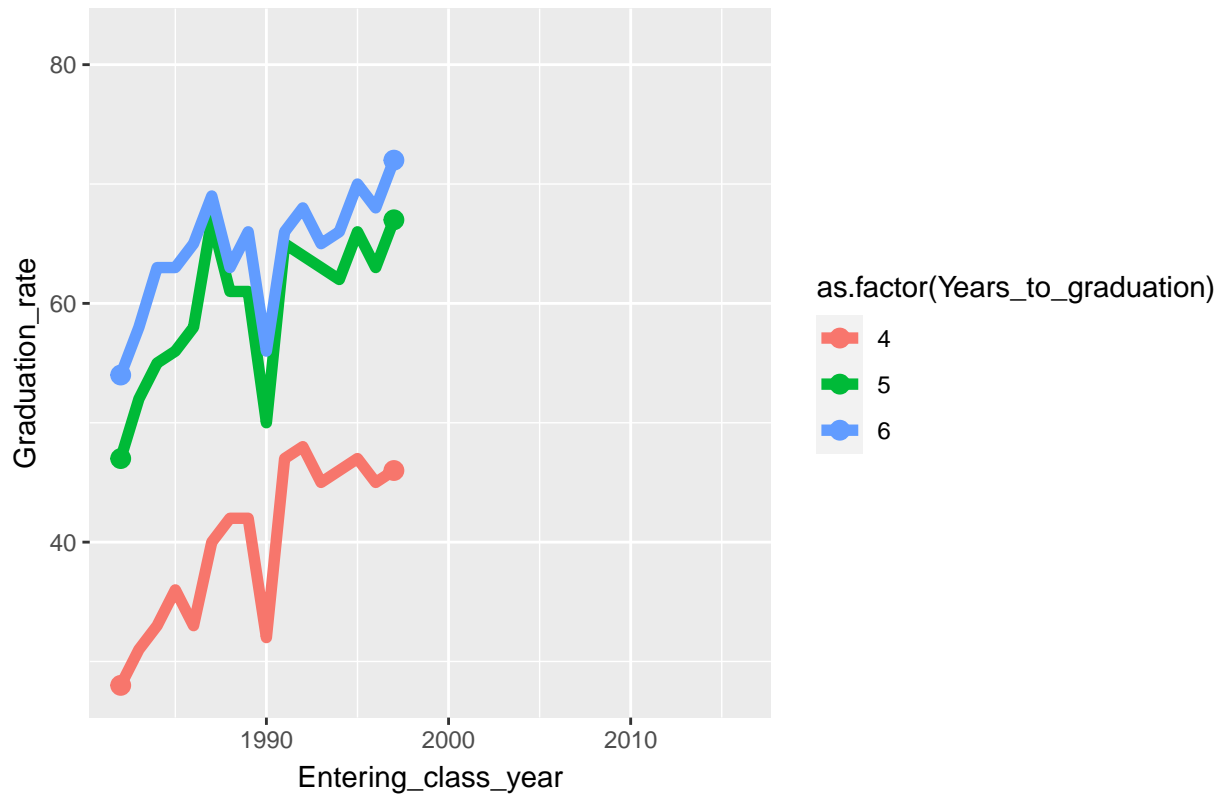
The Year is 1996.



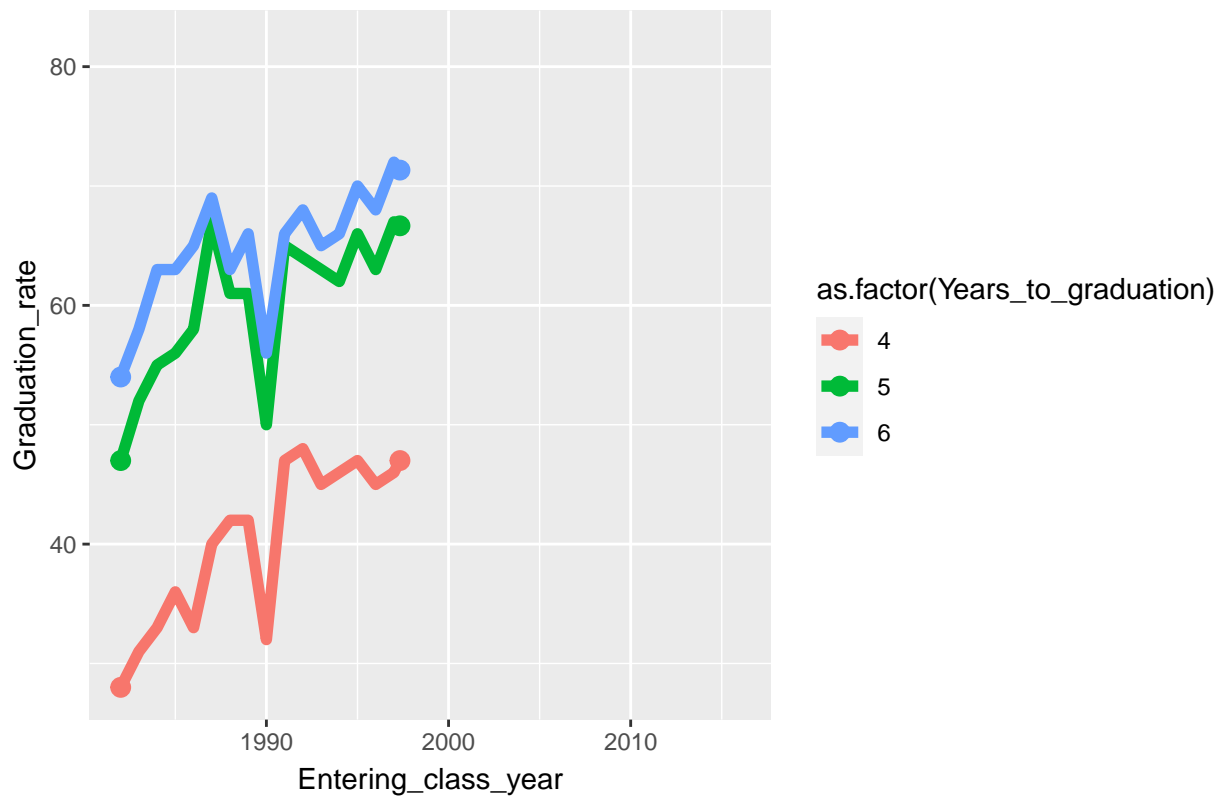
The Year is 1997.



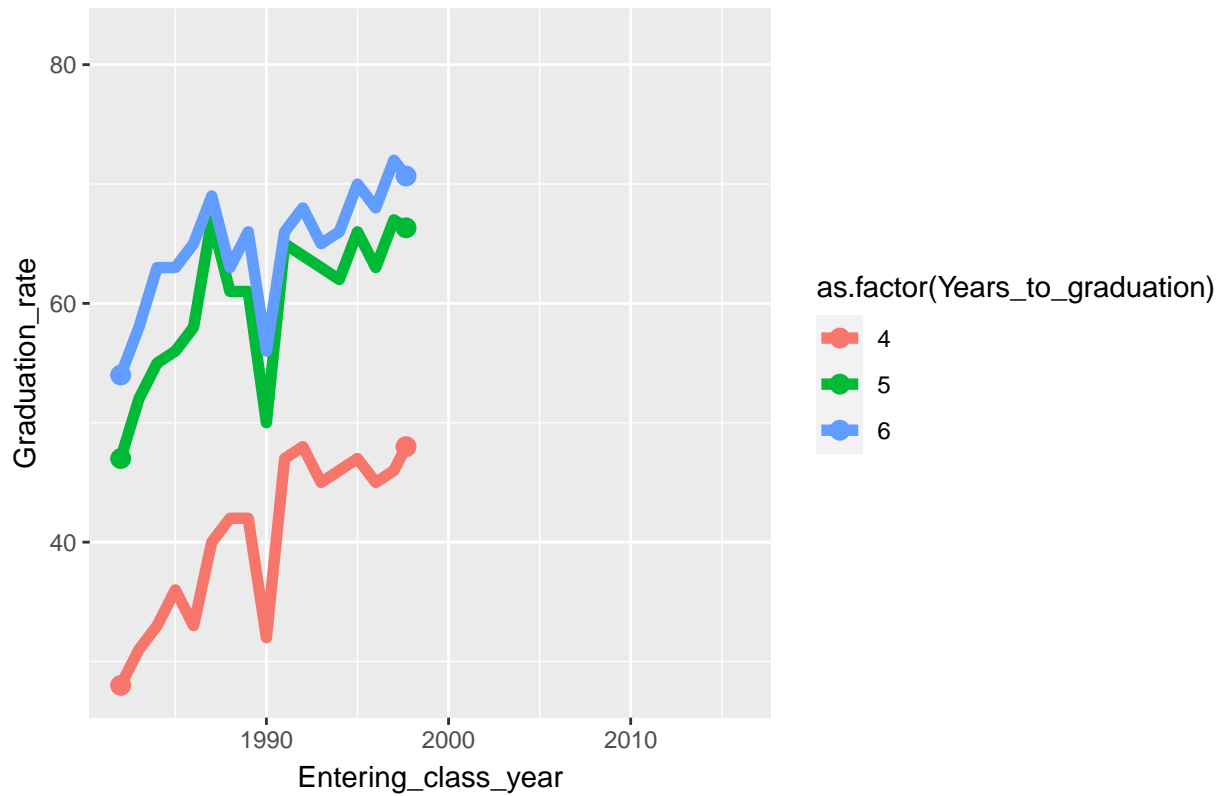
The Year is 1997.



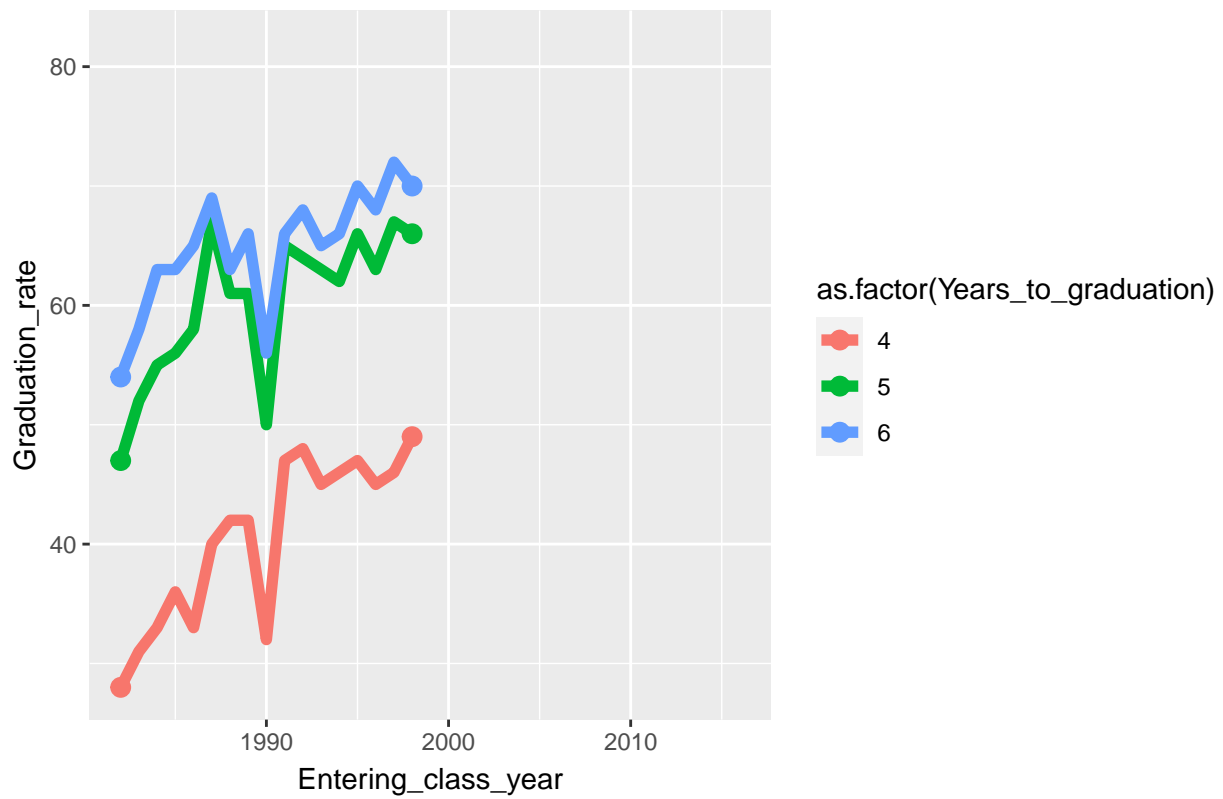
The Year is 1997.



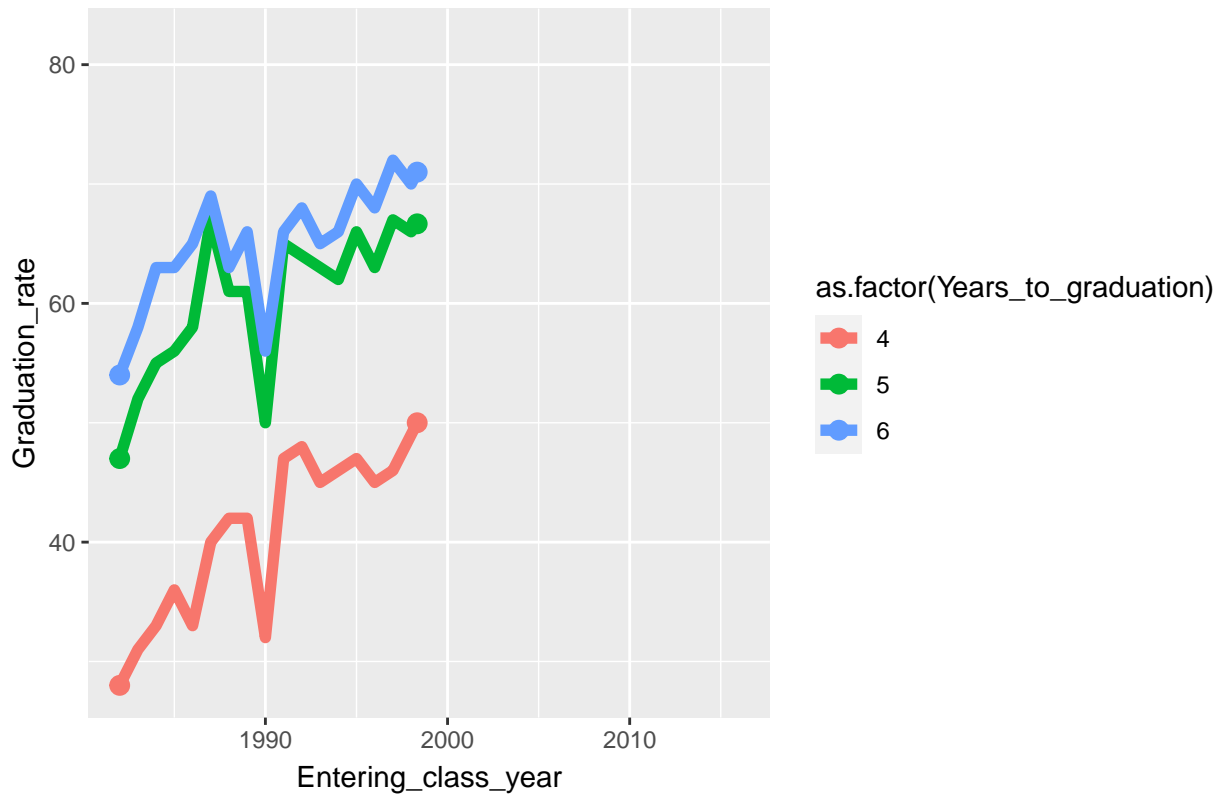
The Year is 1998.



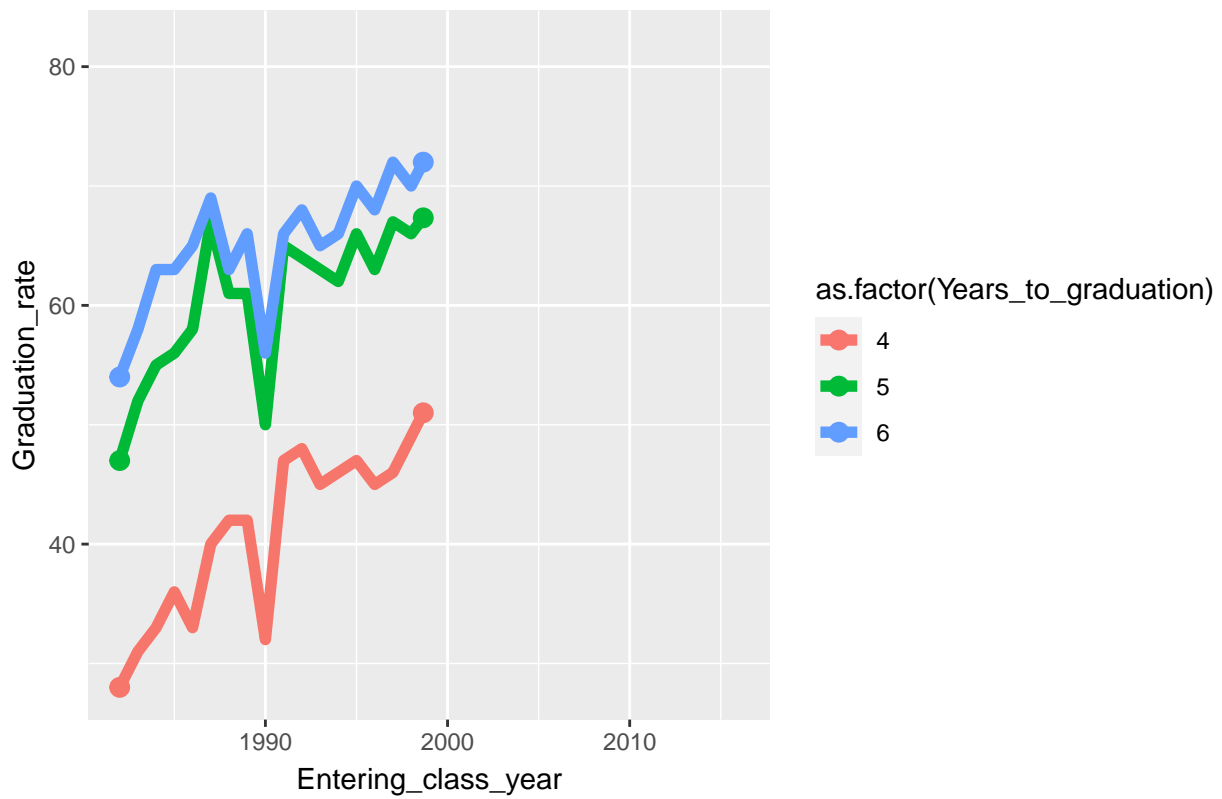
The Year is 1998.



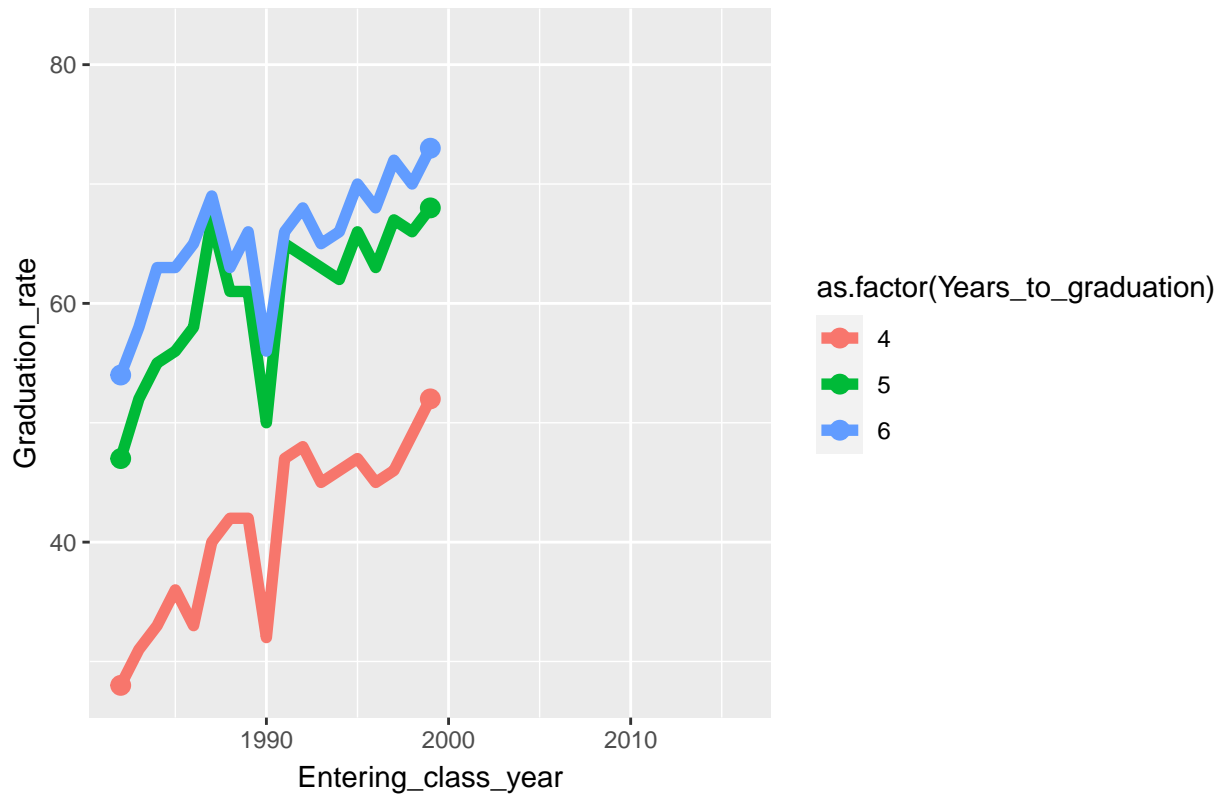
The Year is 1998.



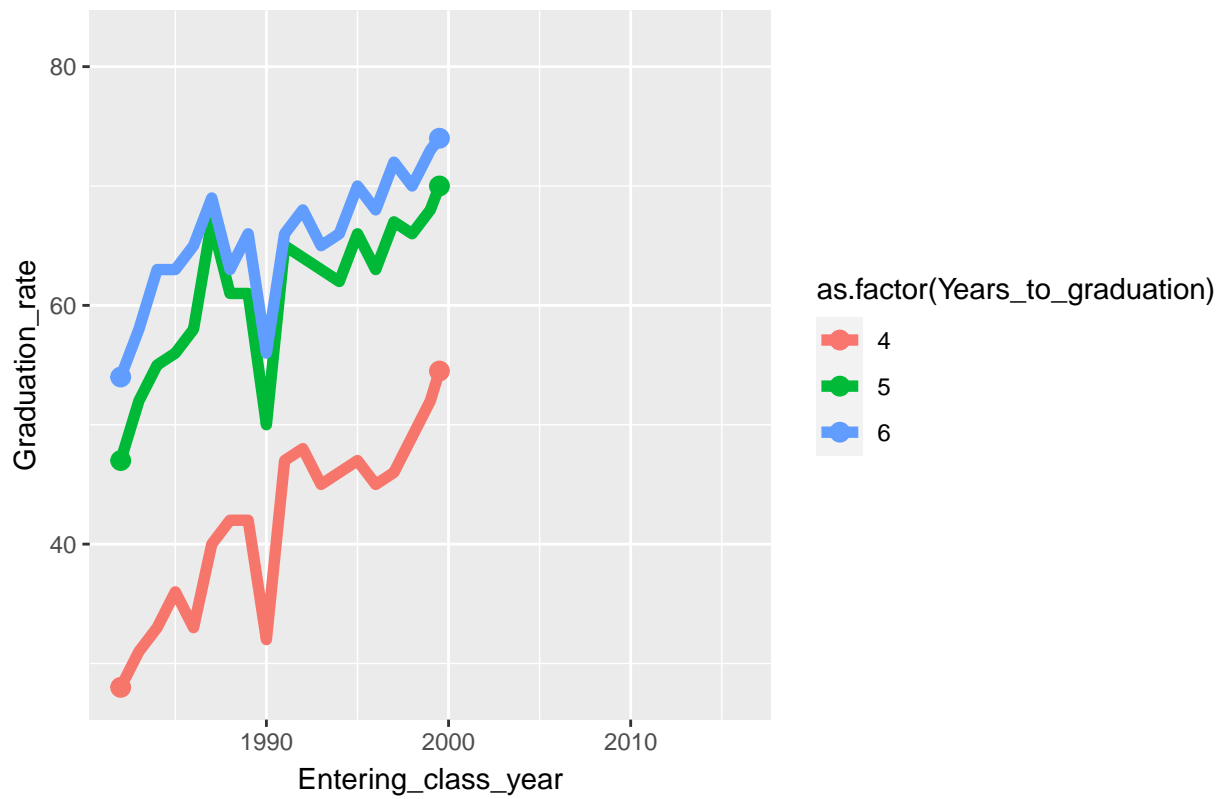
The Year is 1999.



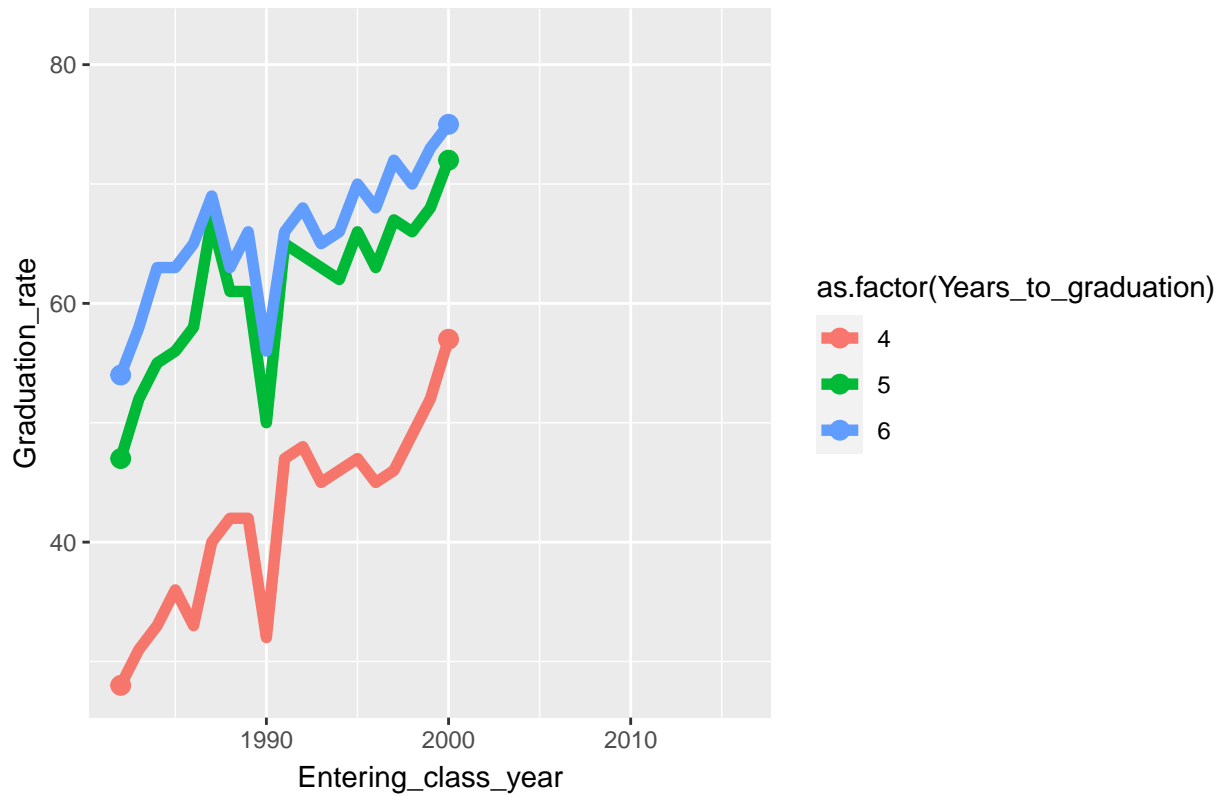
The Year is 1999.



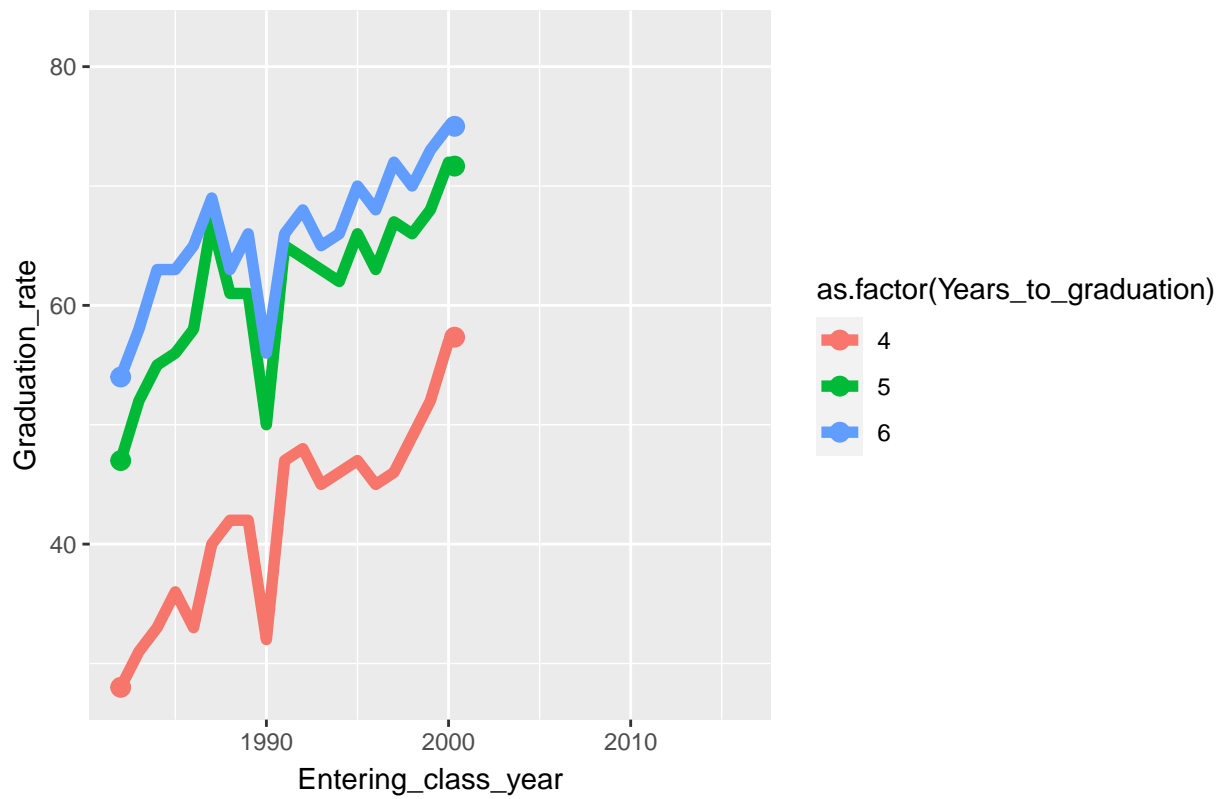
The Year is 2000.



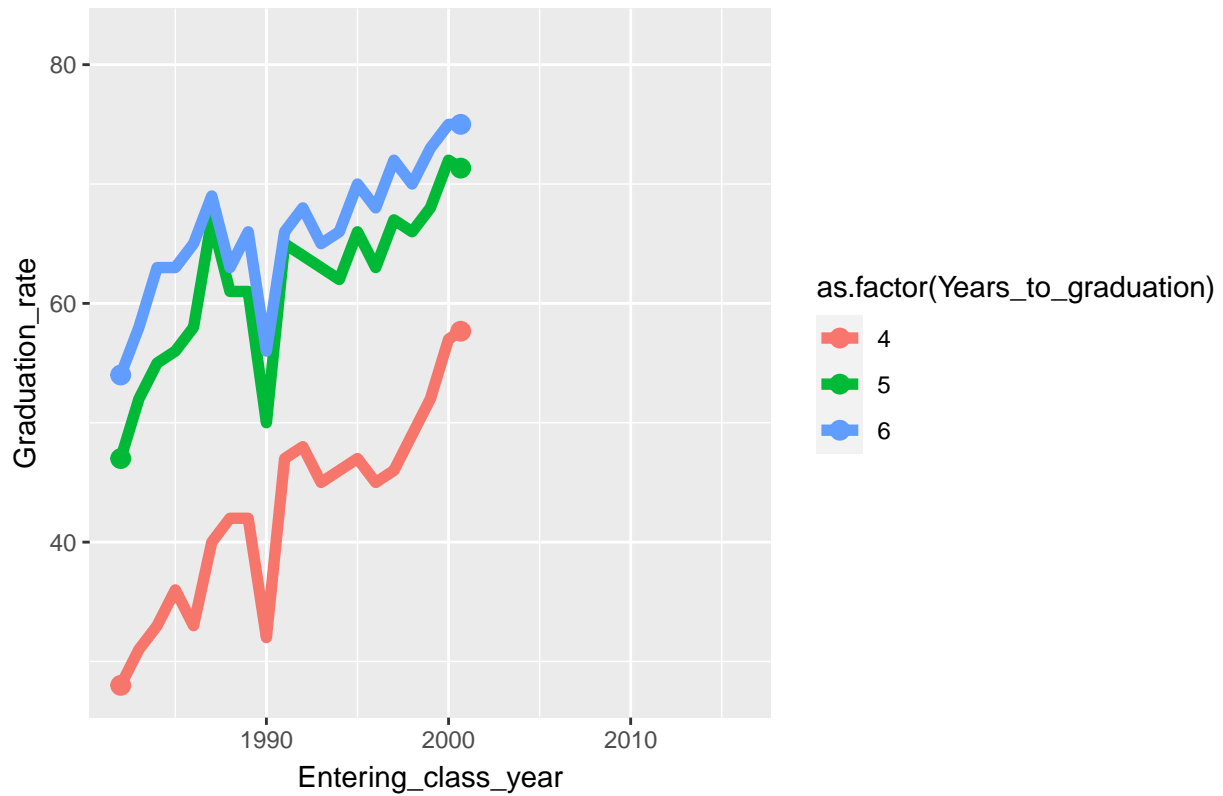
The Year is 2000.



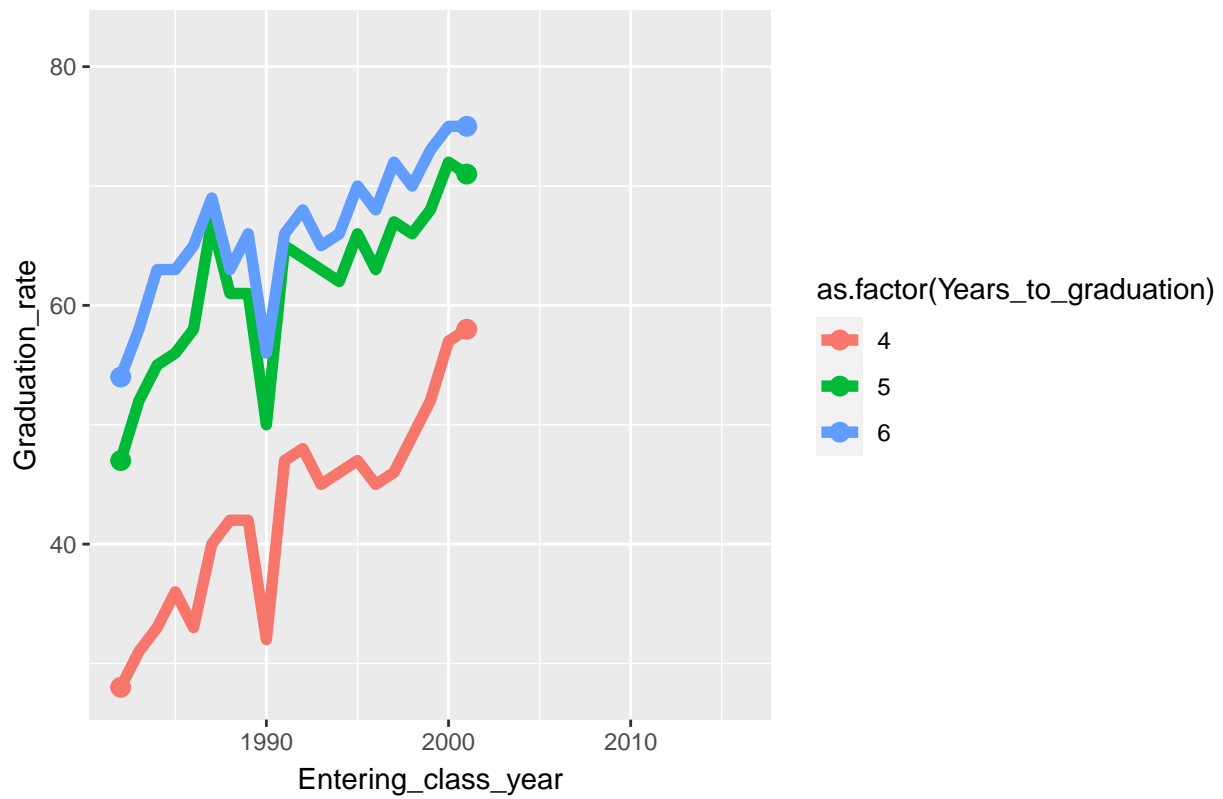
The Year is 2000.



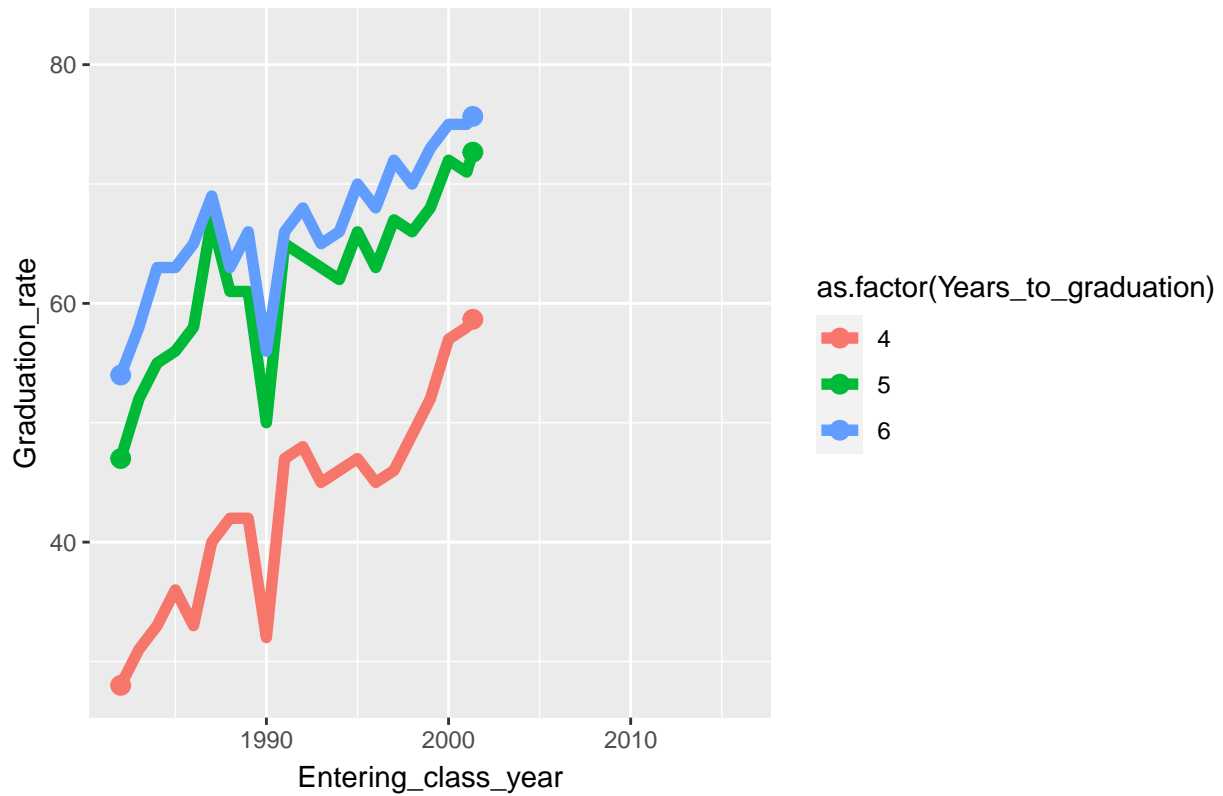
The Year is 2001.



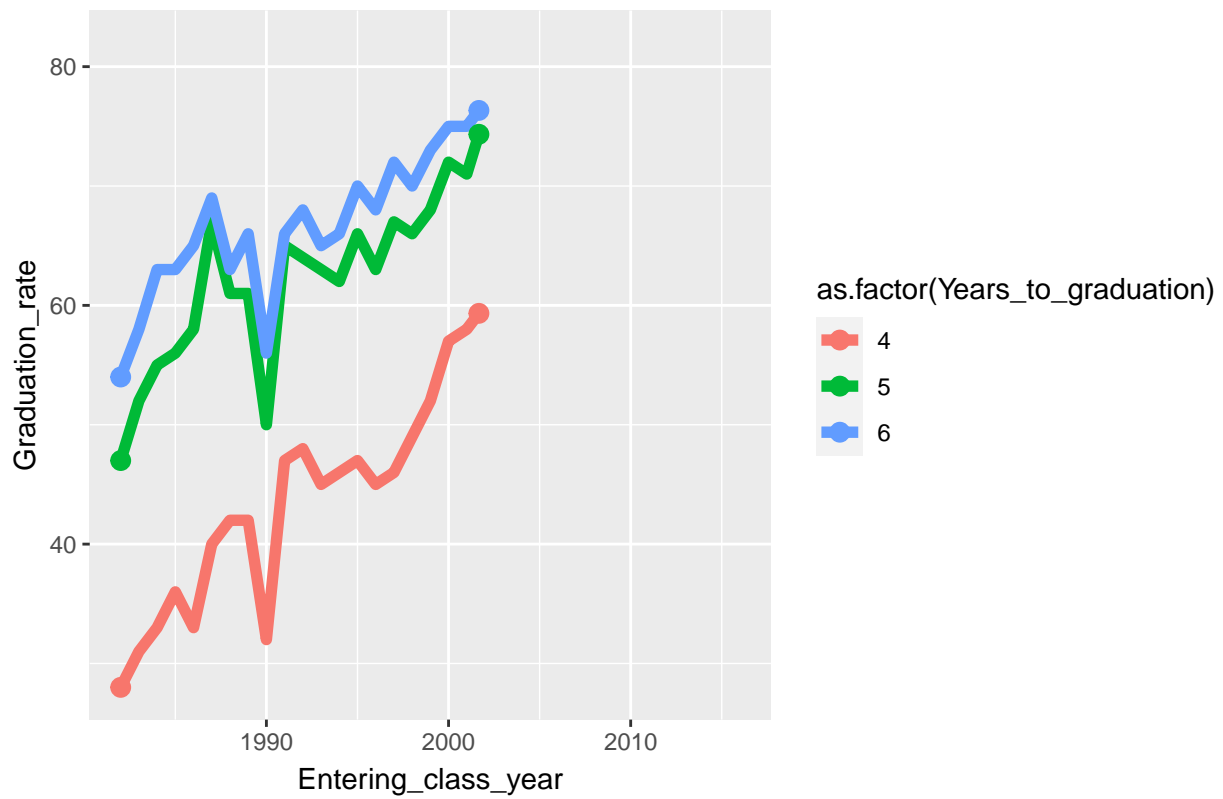
The Year is 2001.



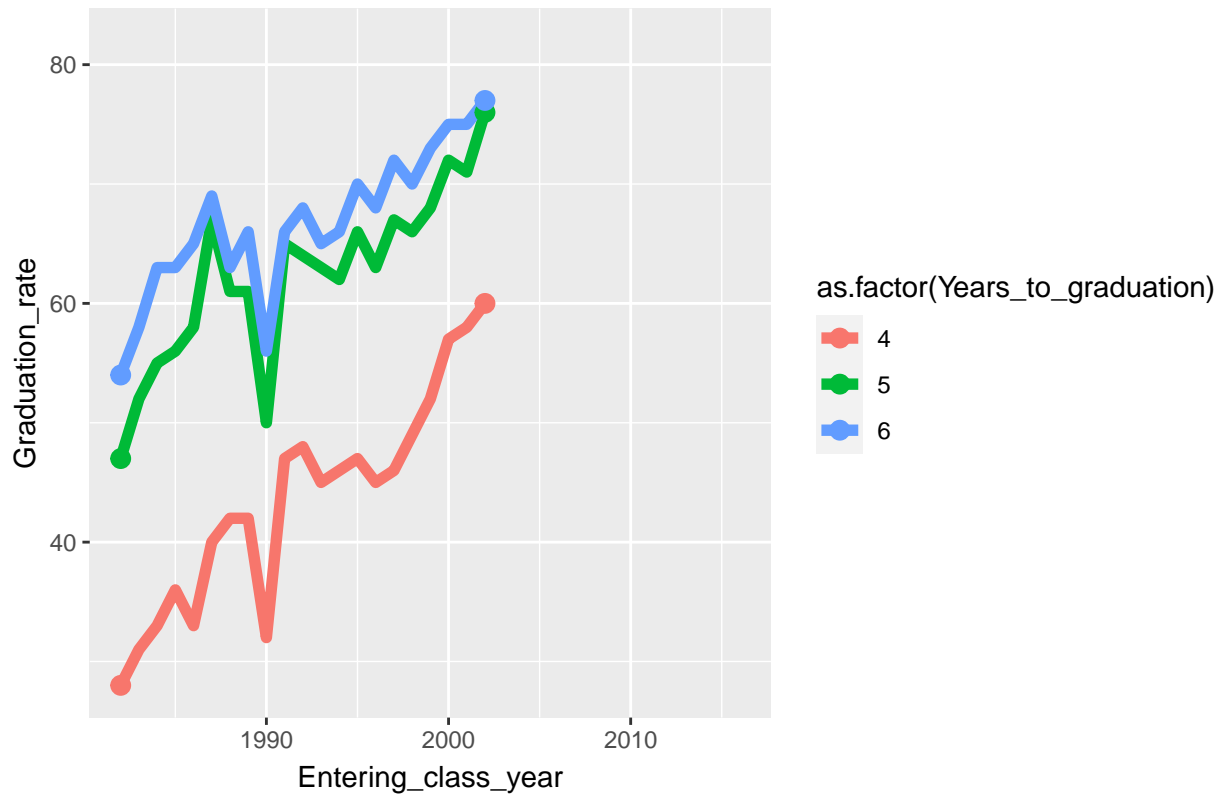
The Year is 2001.



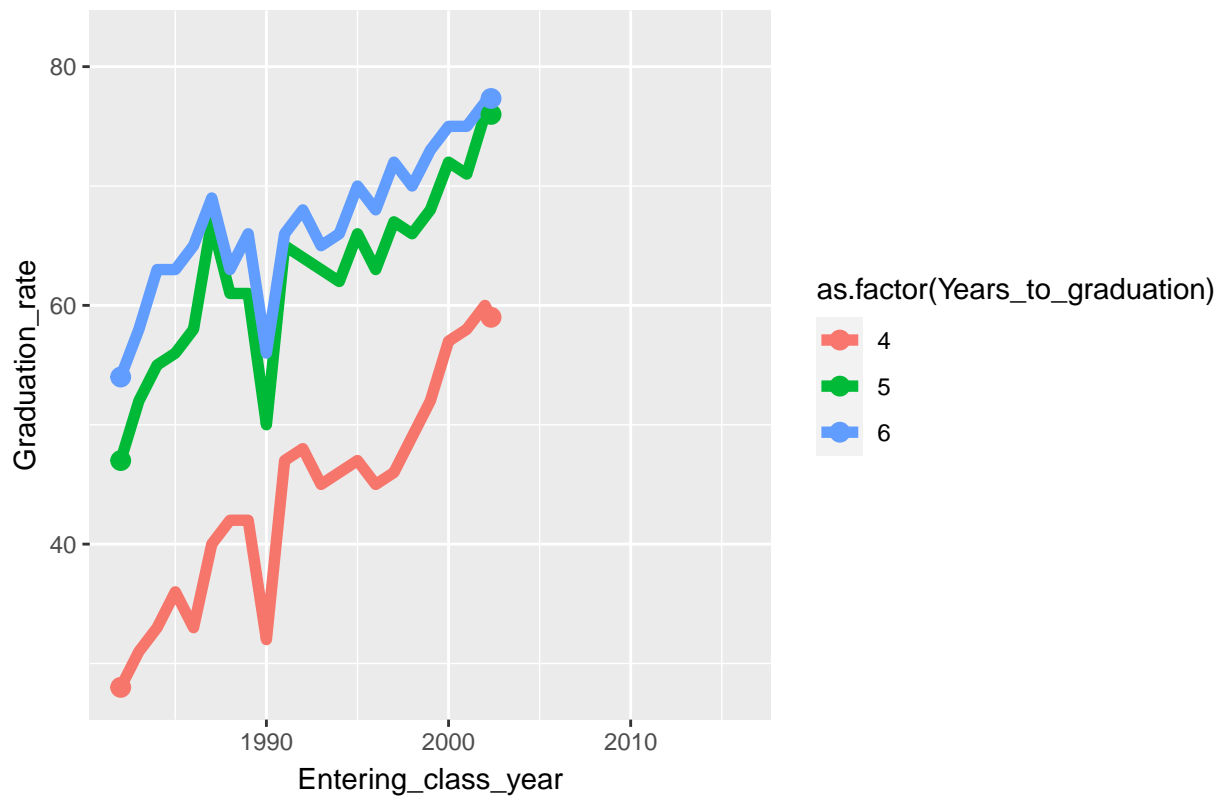
The Year is 2002.



The Year is 2002.



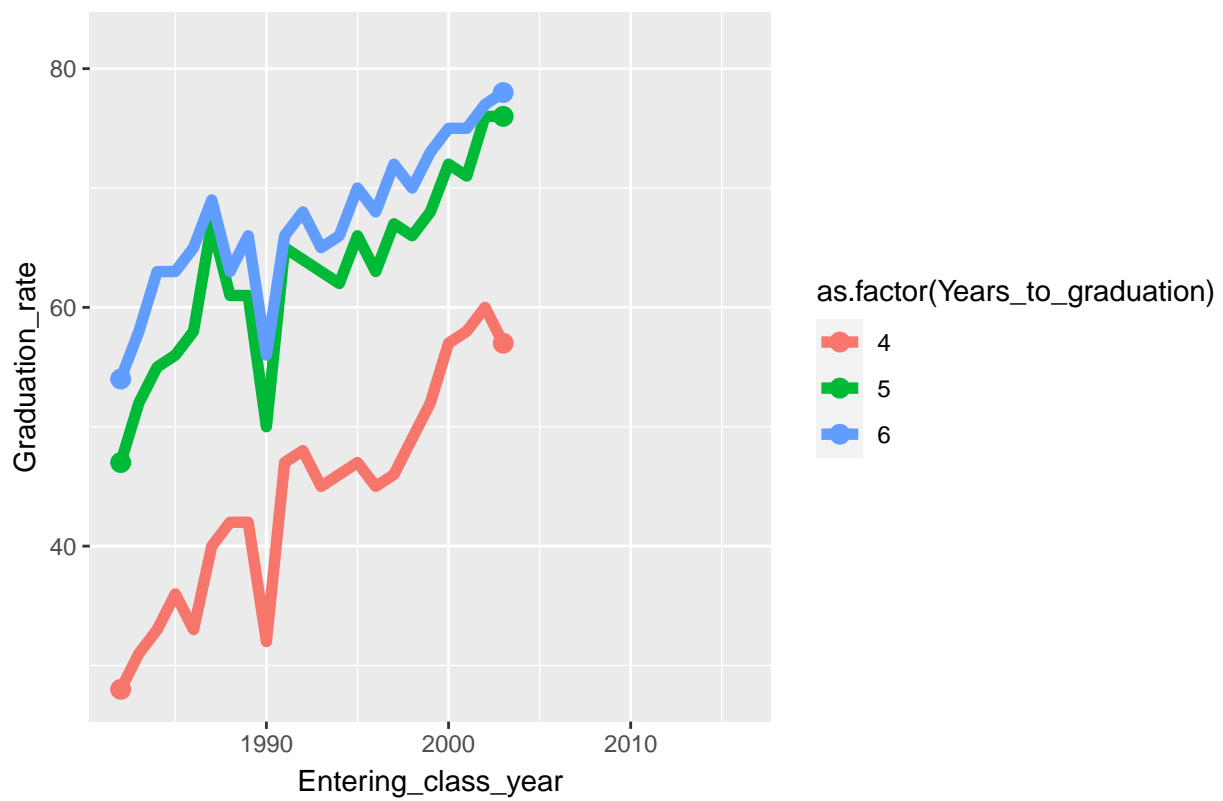
The Year is 2002.



The Year is 2003.



The Year is 2003.



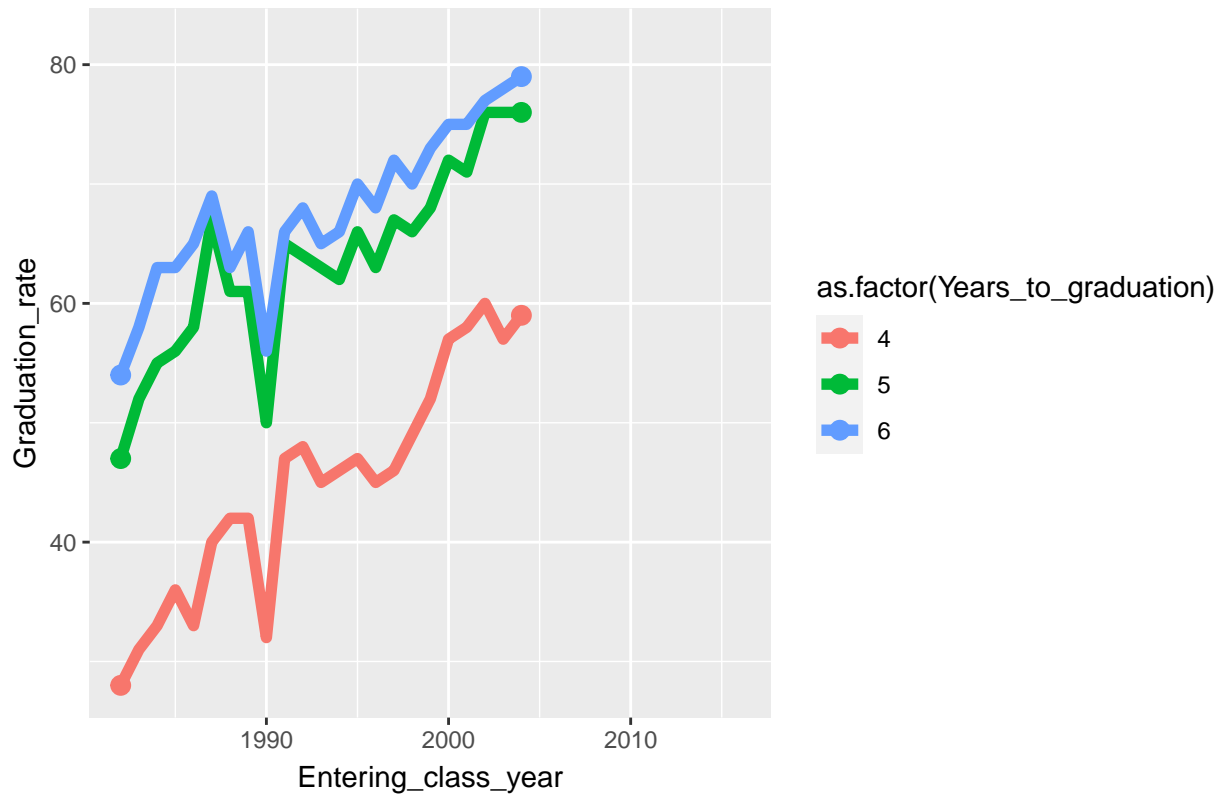
The Year is 2003.



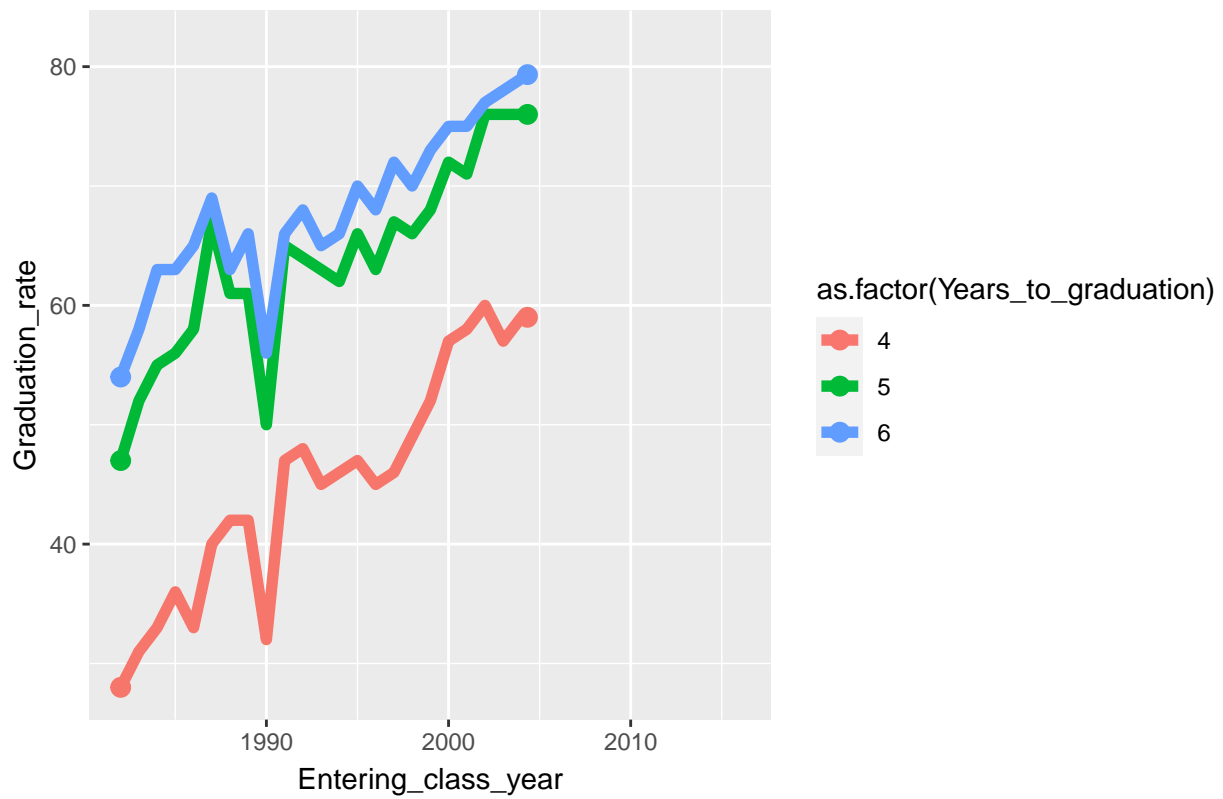
The Year is 2004.



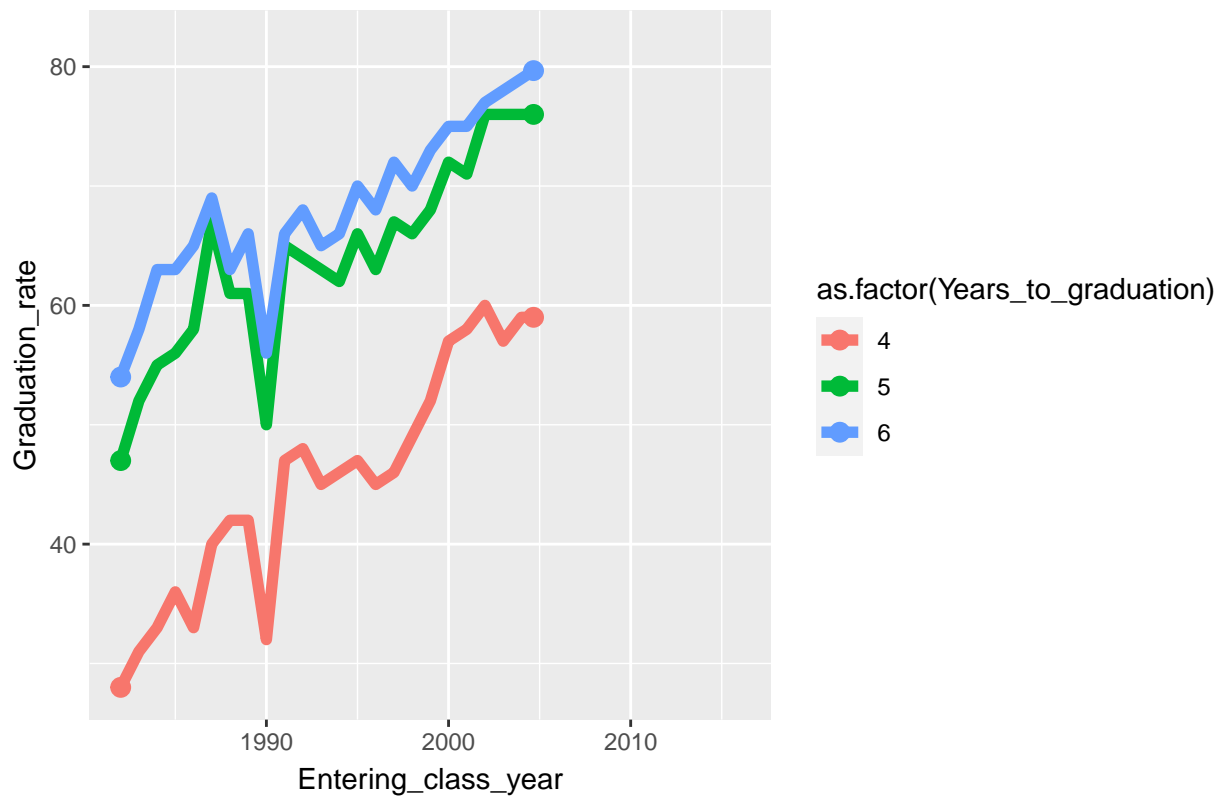
The Year is 2004.



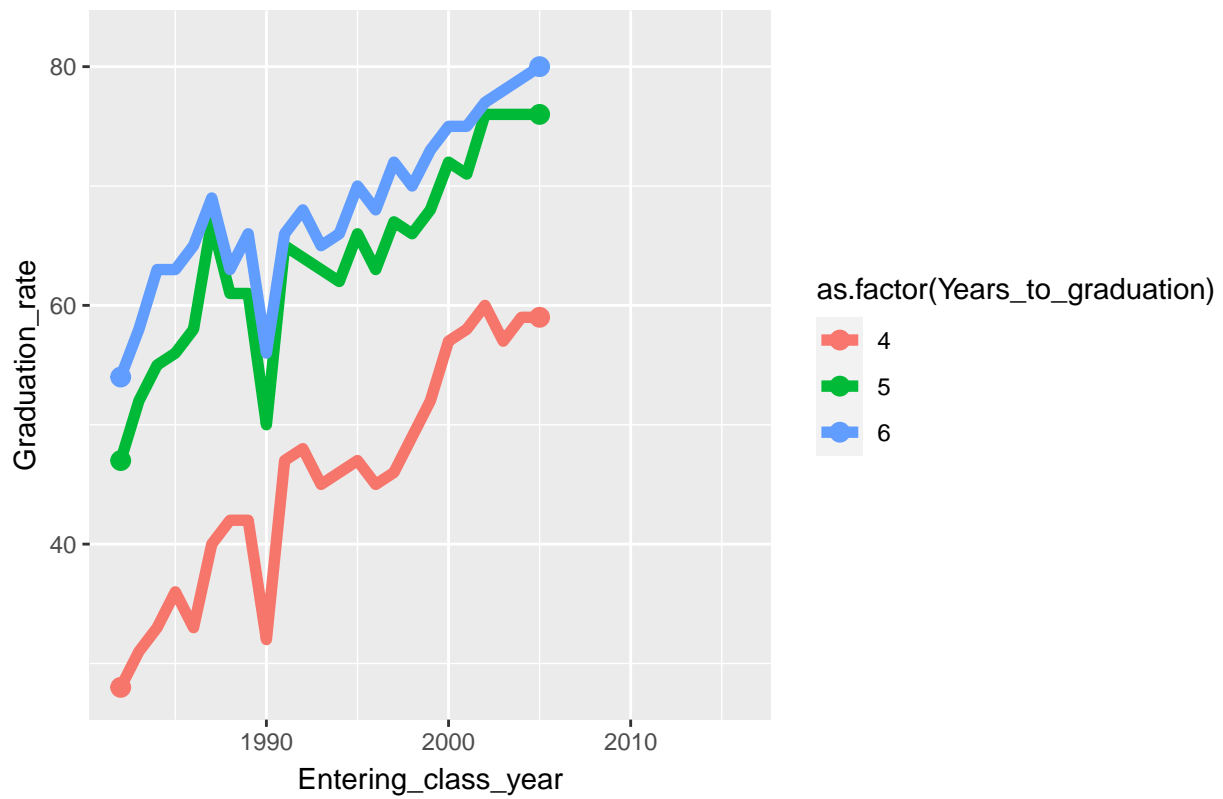
The Year is 2004.



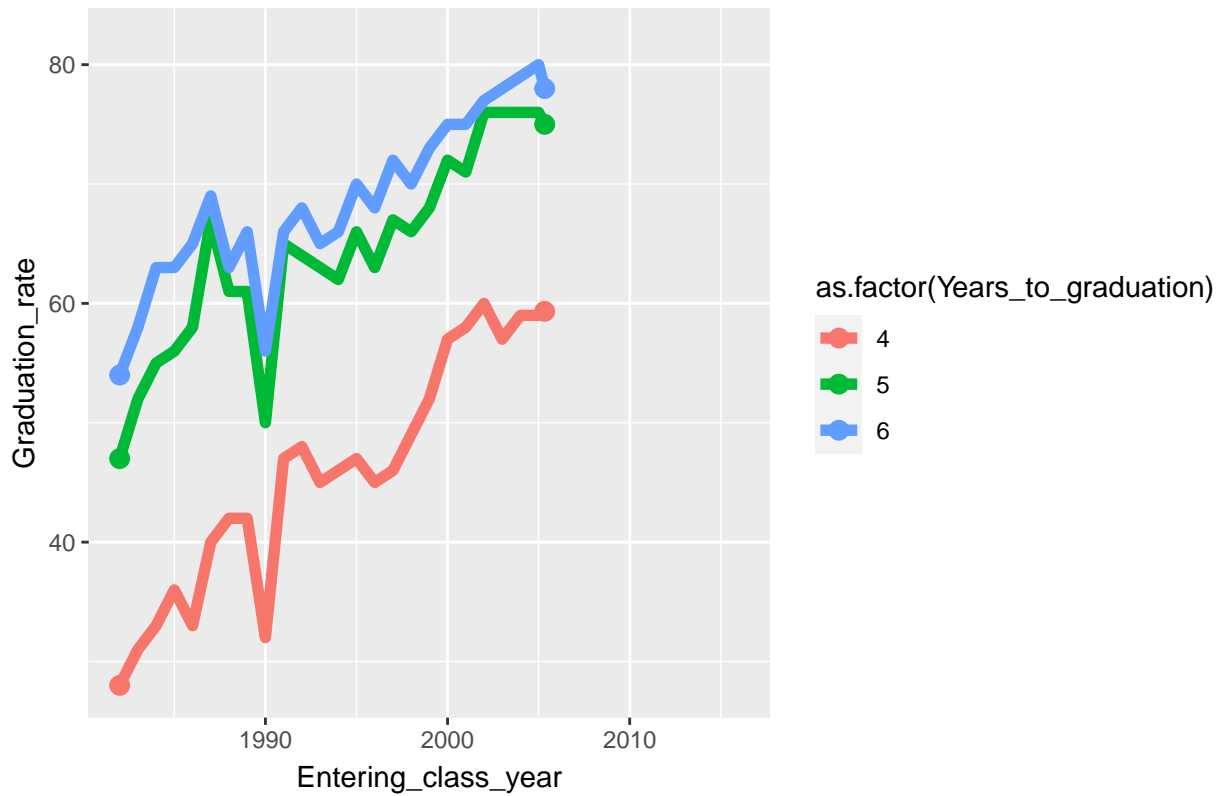
The Year is 2005.



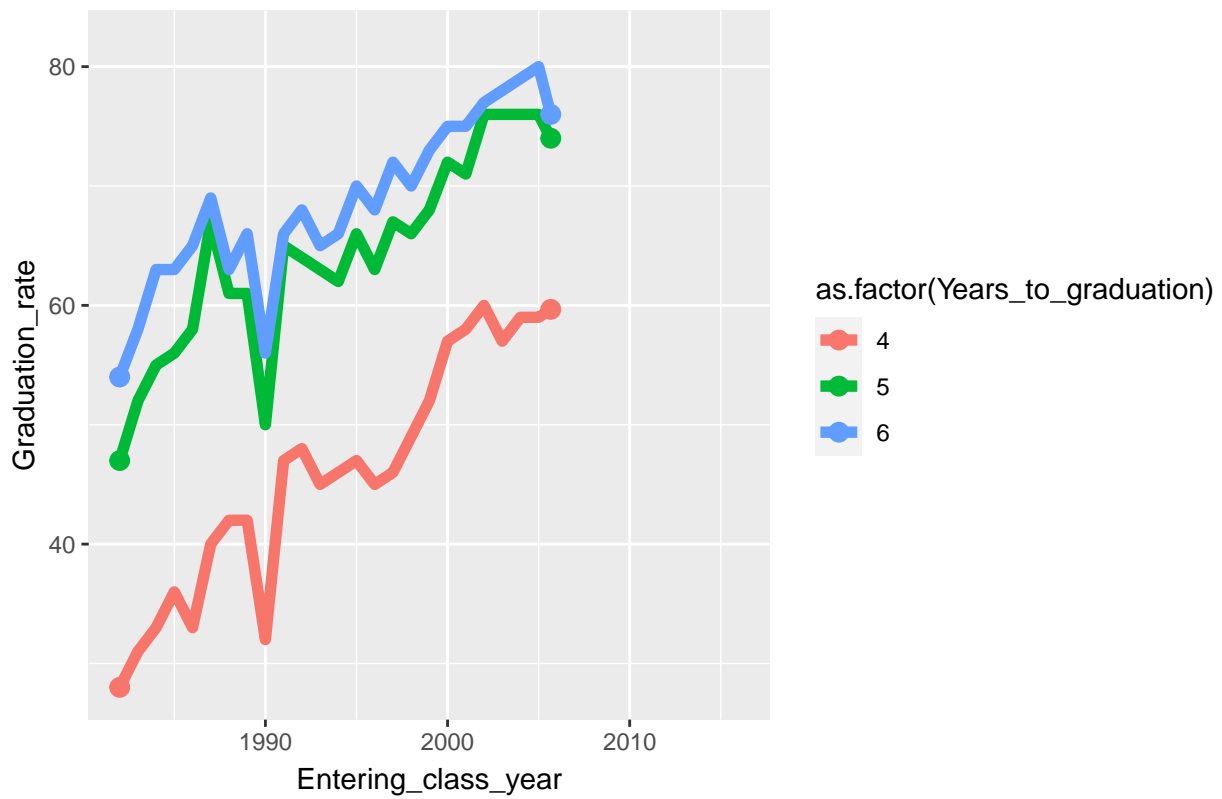
The Year is 2005.



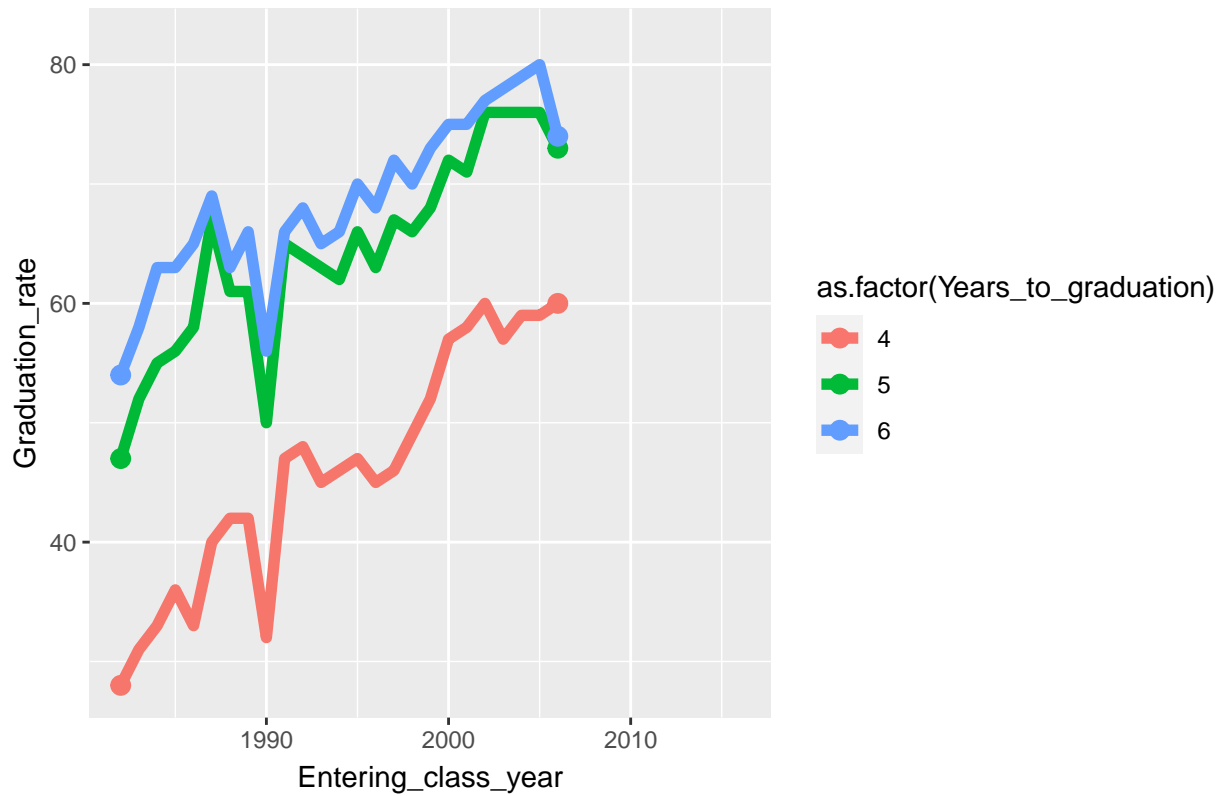
The Year is 2005.



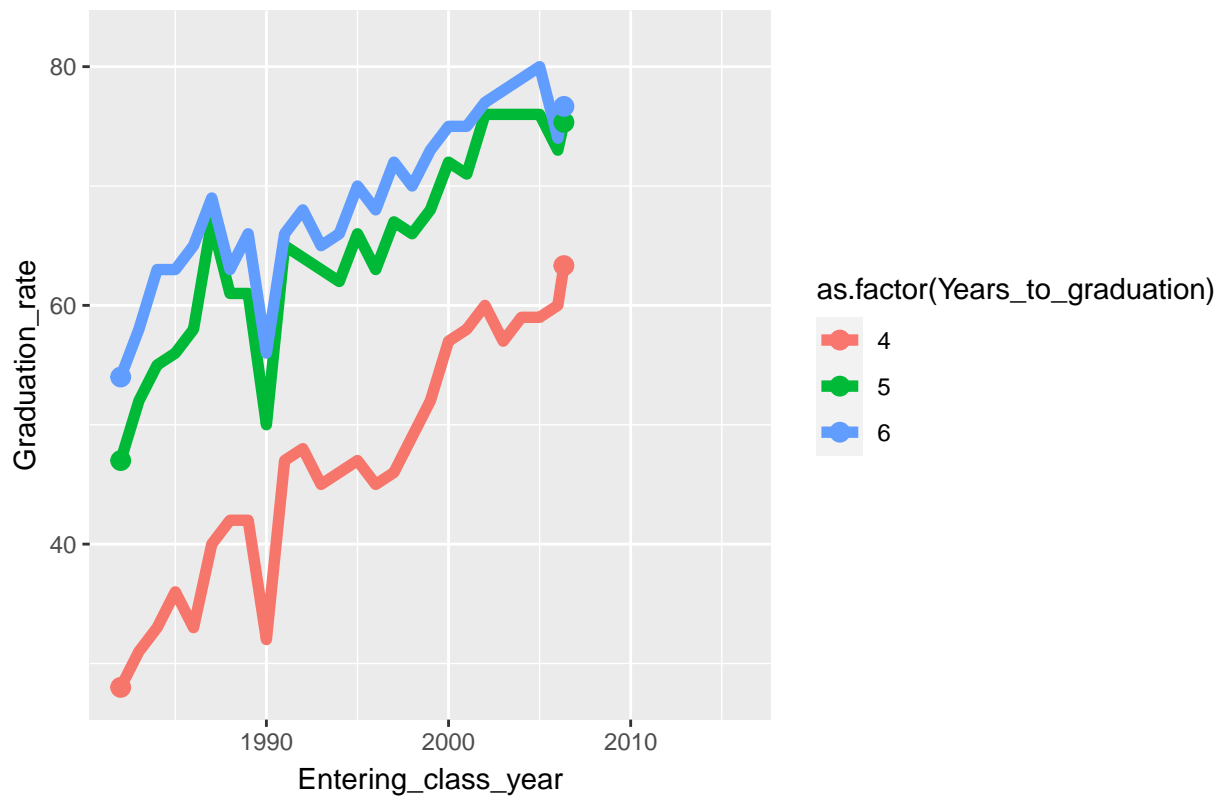
The Year is 2006.



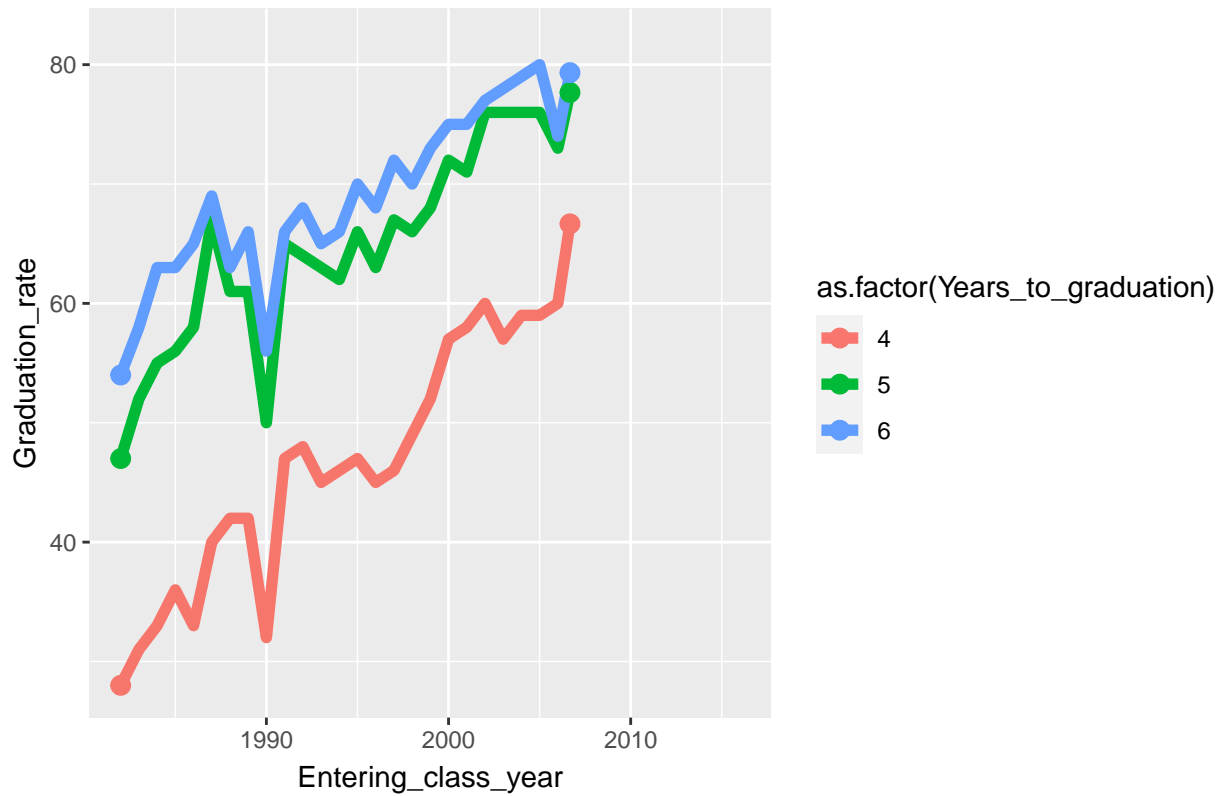
The Year is 2006.



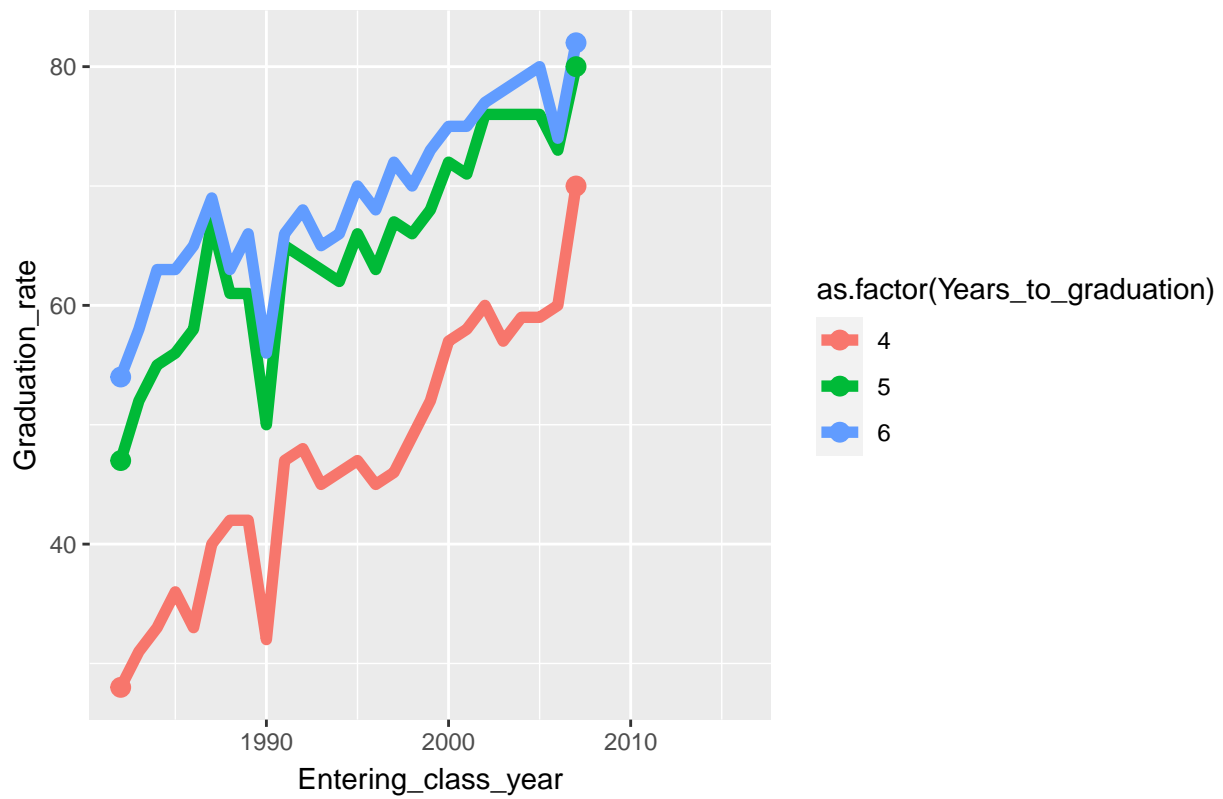
The Year is 2006.



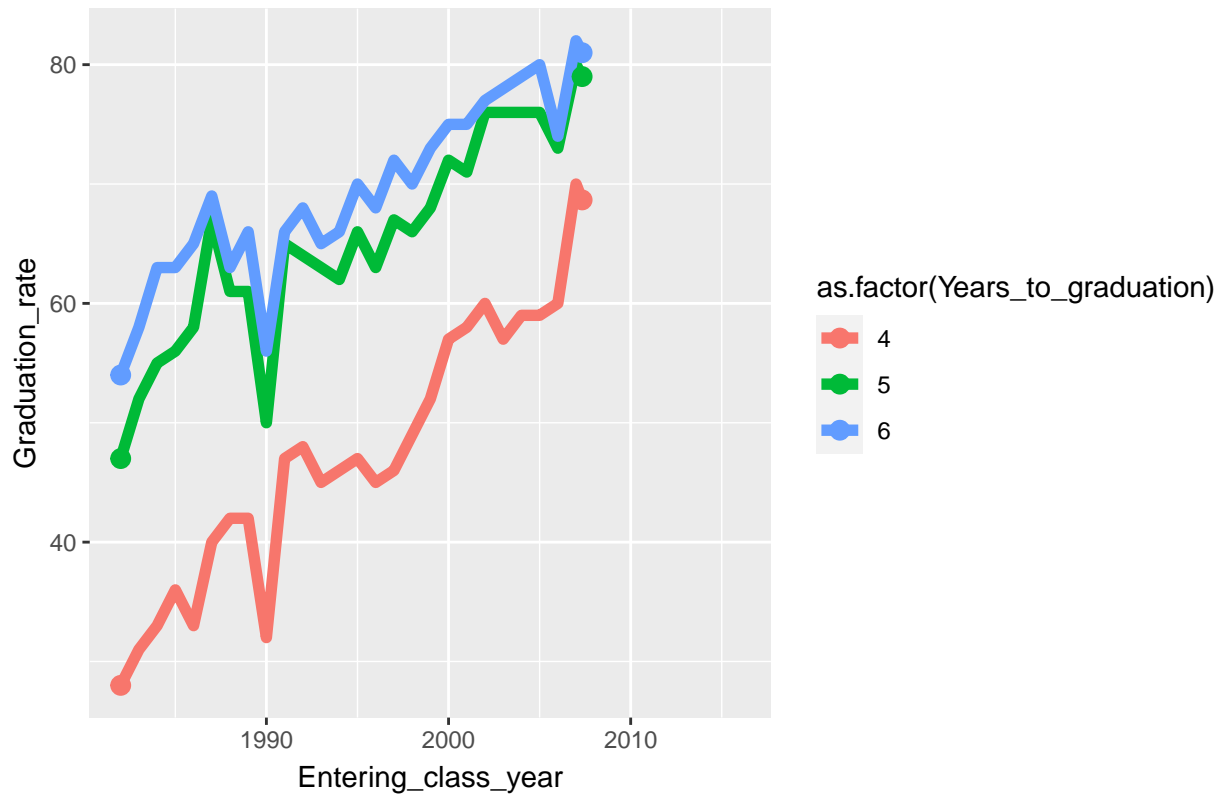
The Year is 2007.



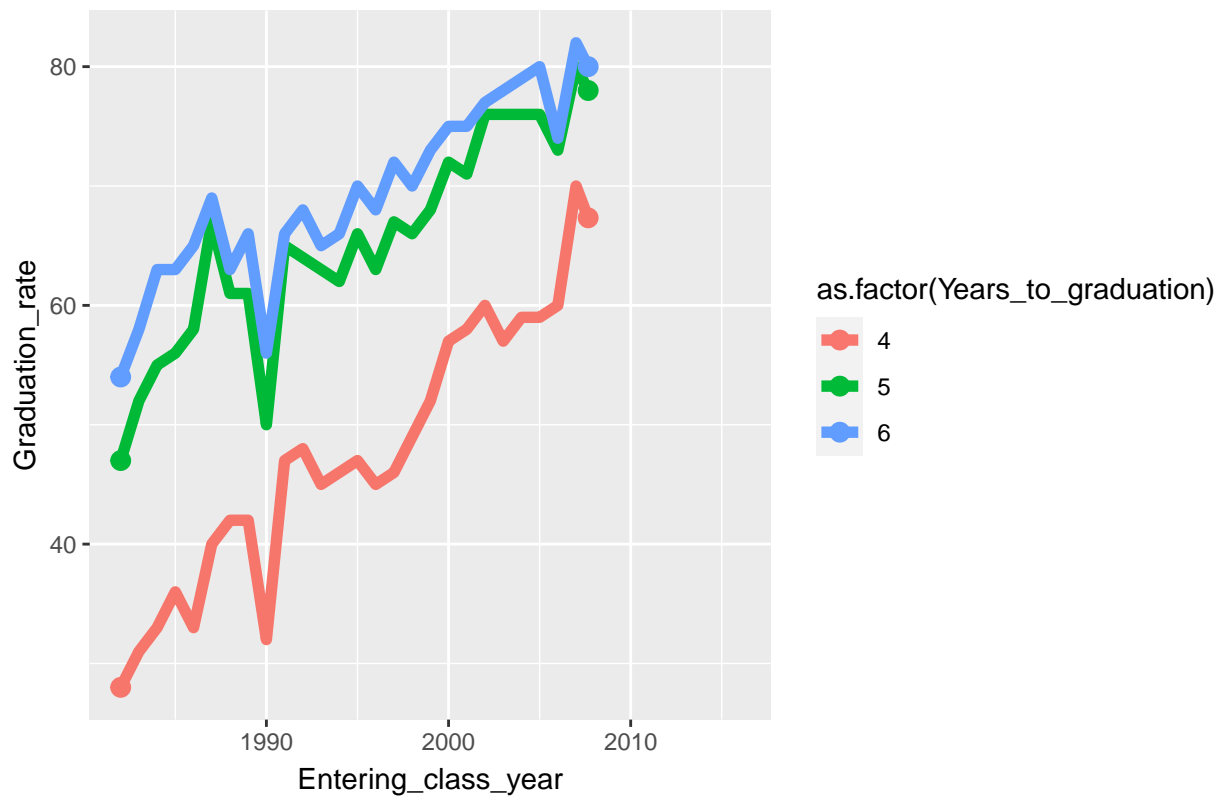
The Year is 2007.



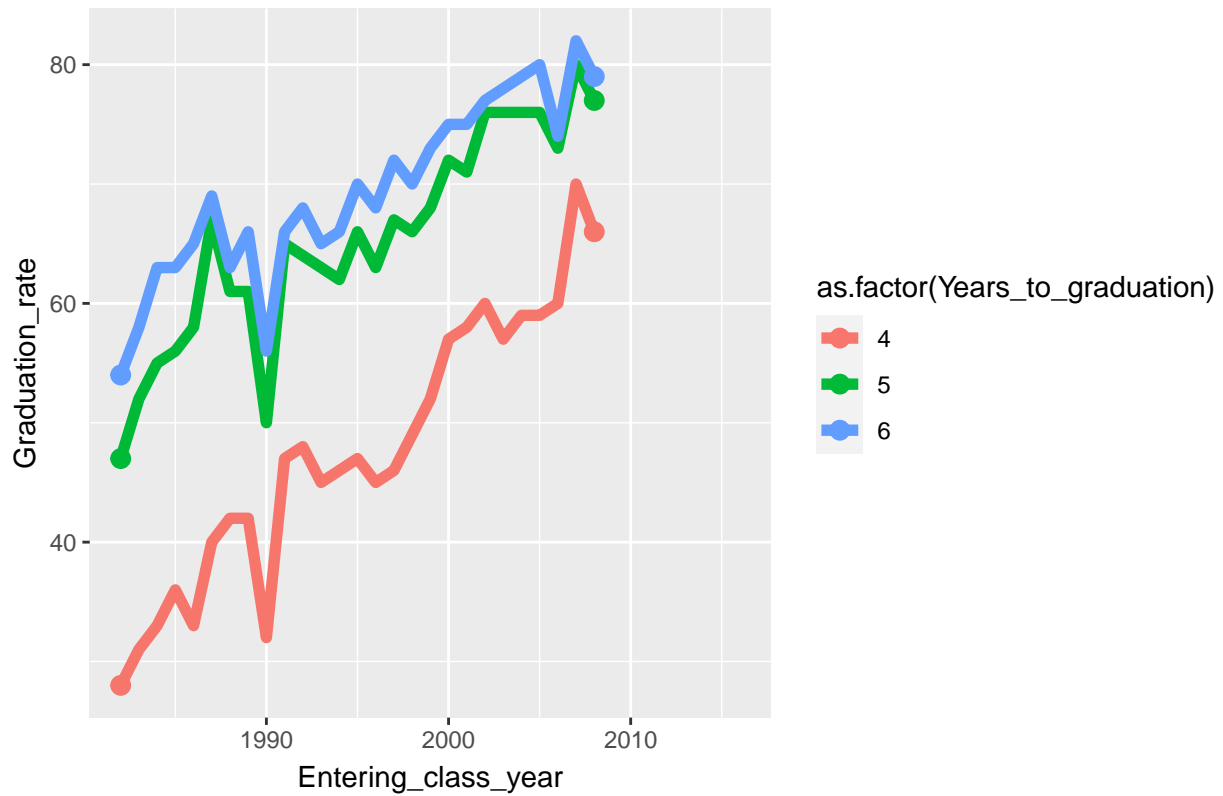
The Year is 2007.



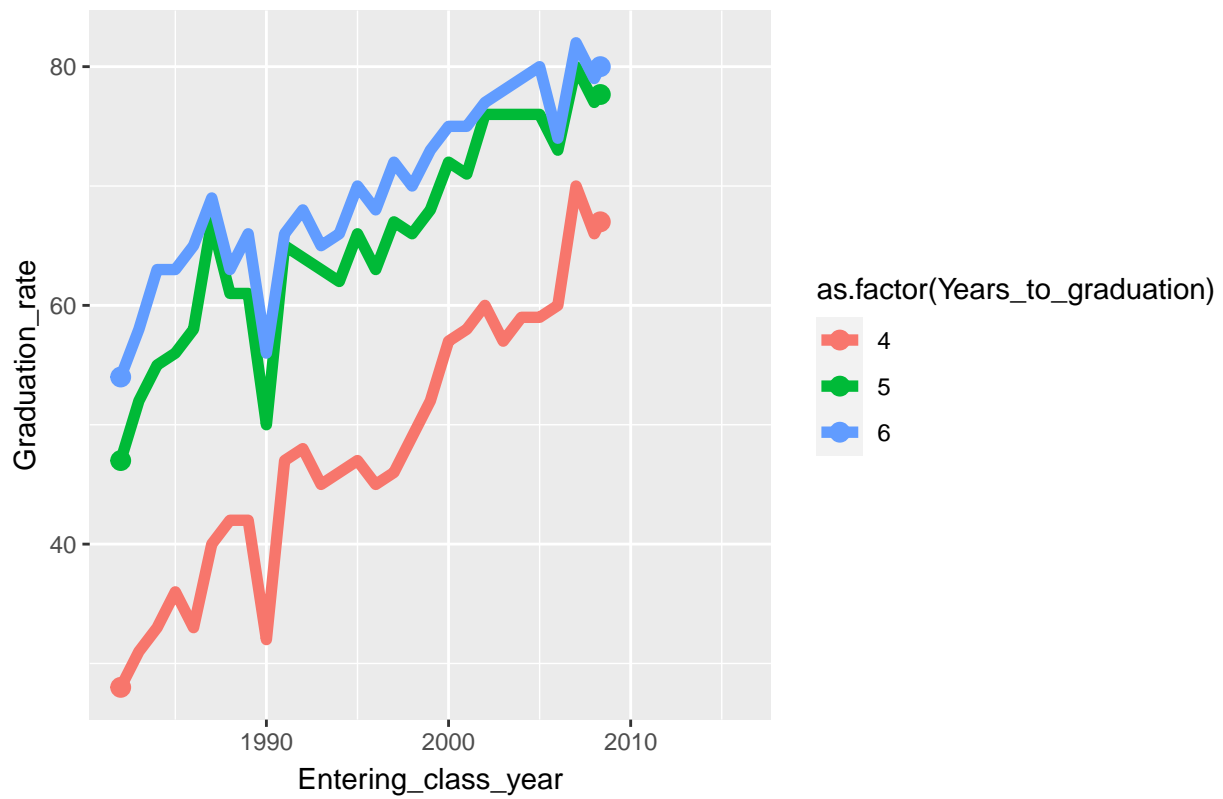
The Year is 2008.



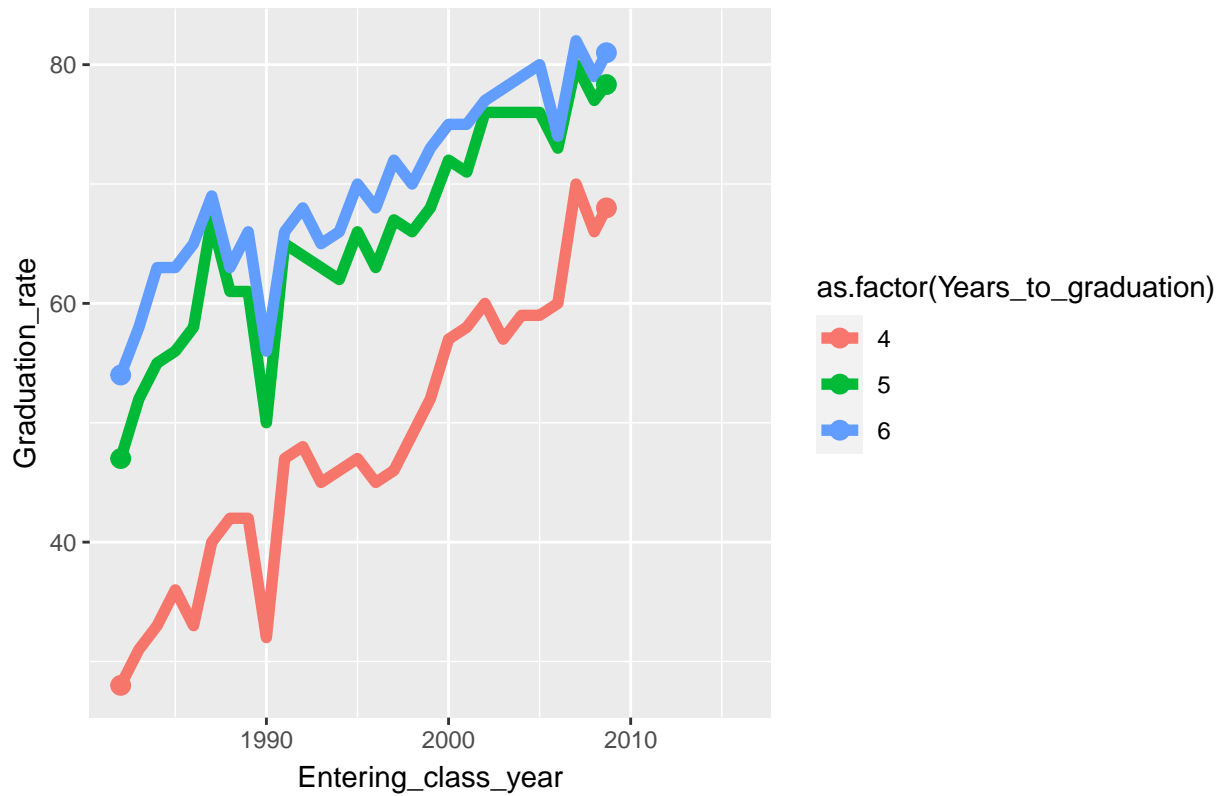
The Year is 2008.



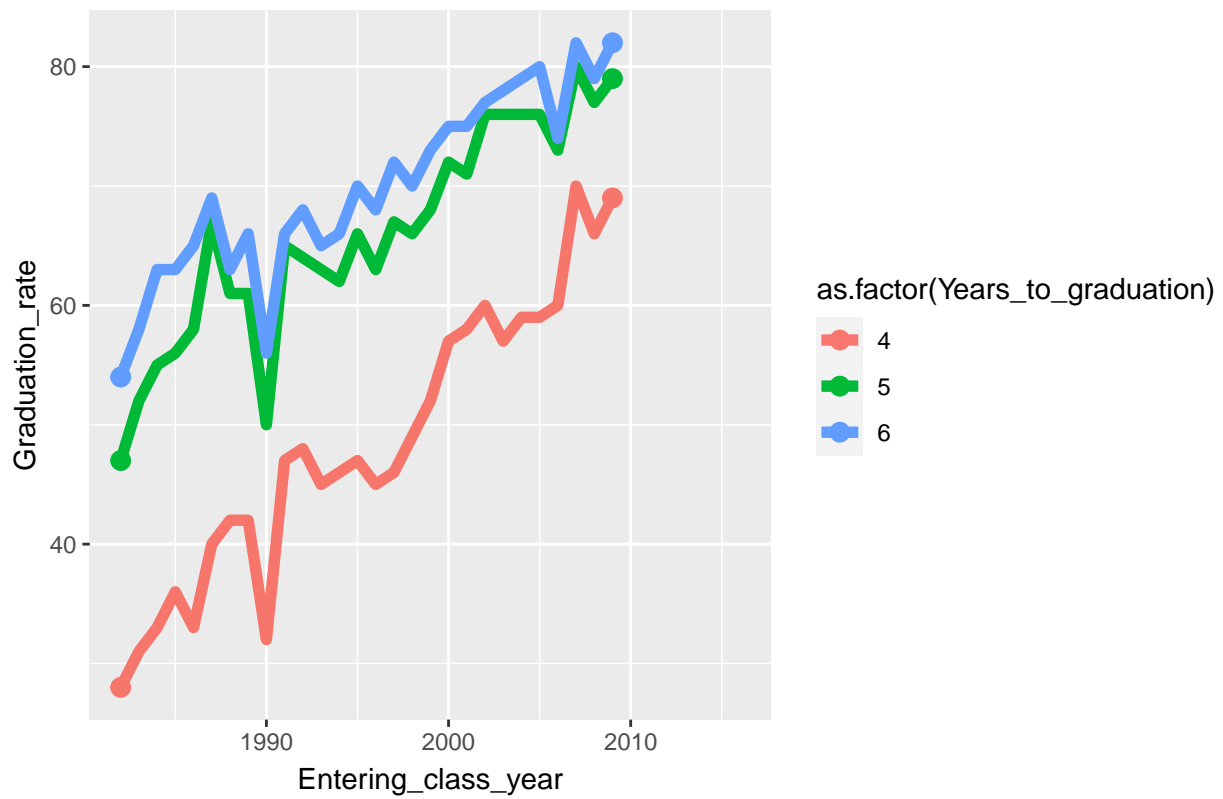
The Year is 2008.



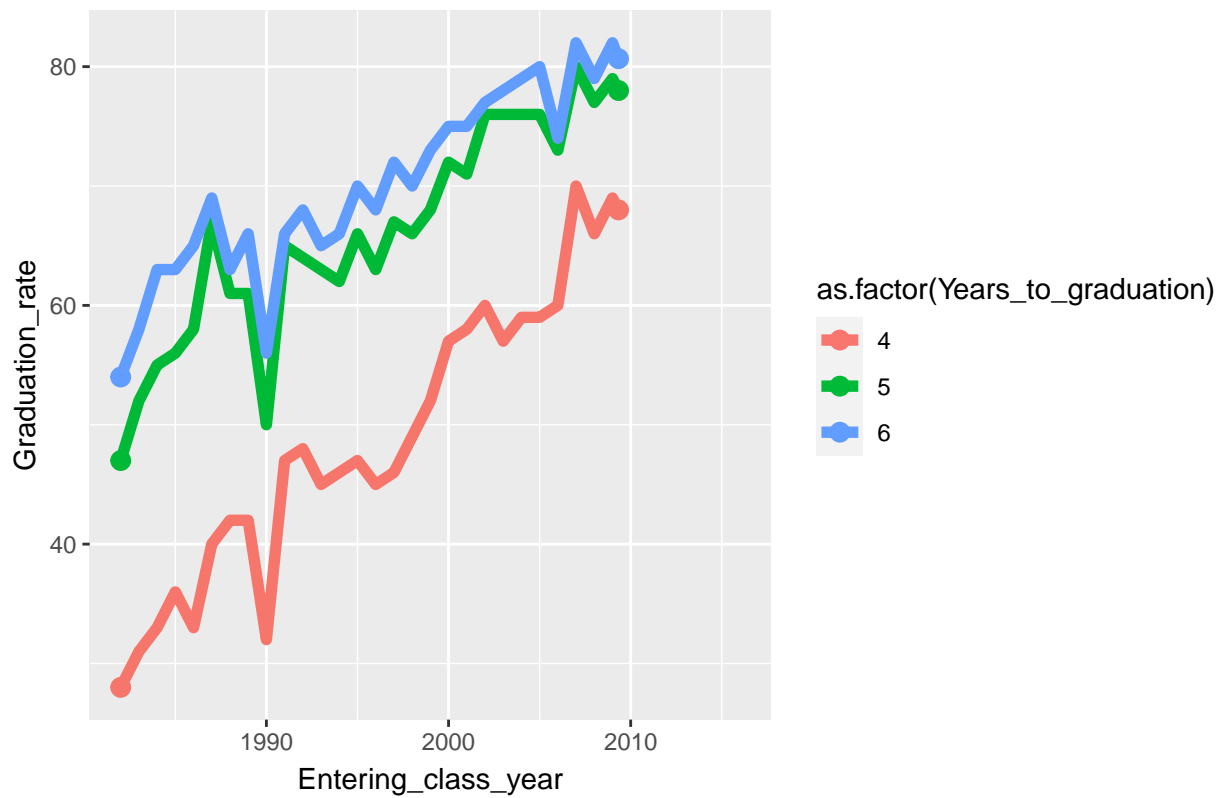
The Year is 2009.



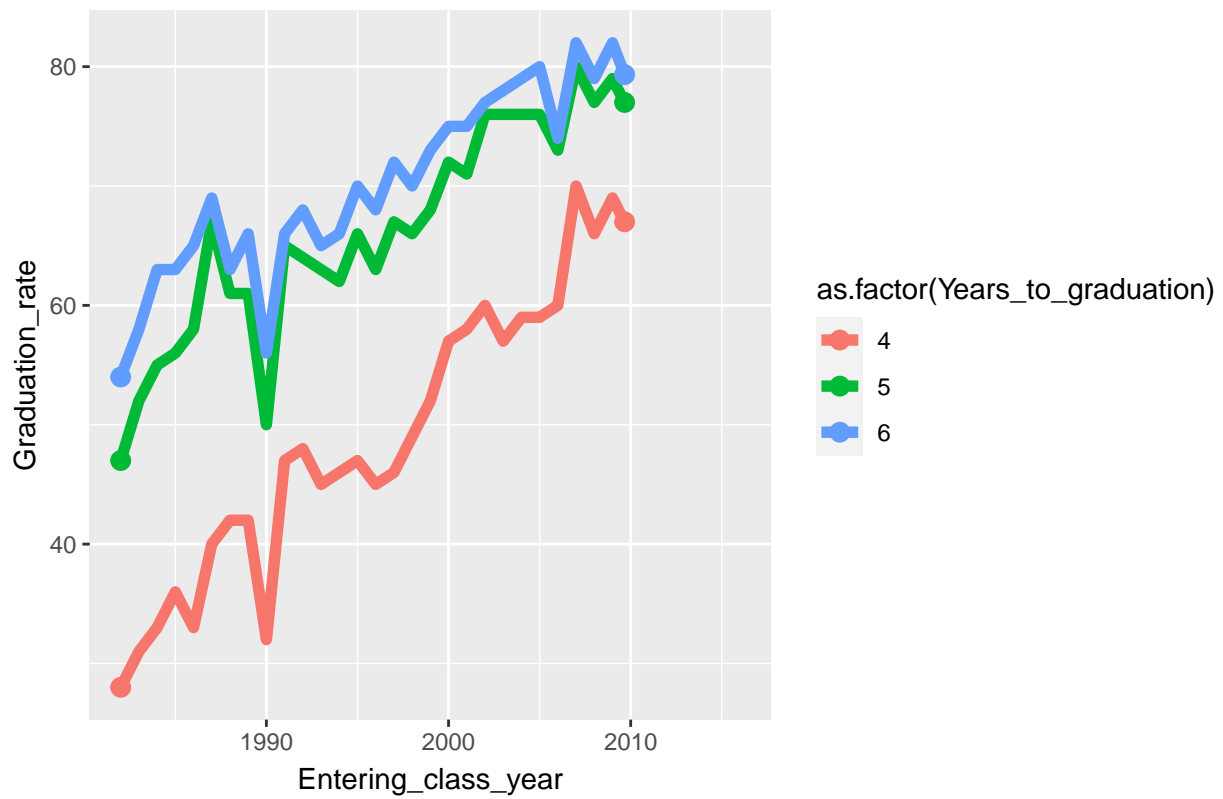
The Year is 2009.



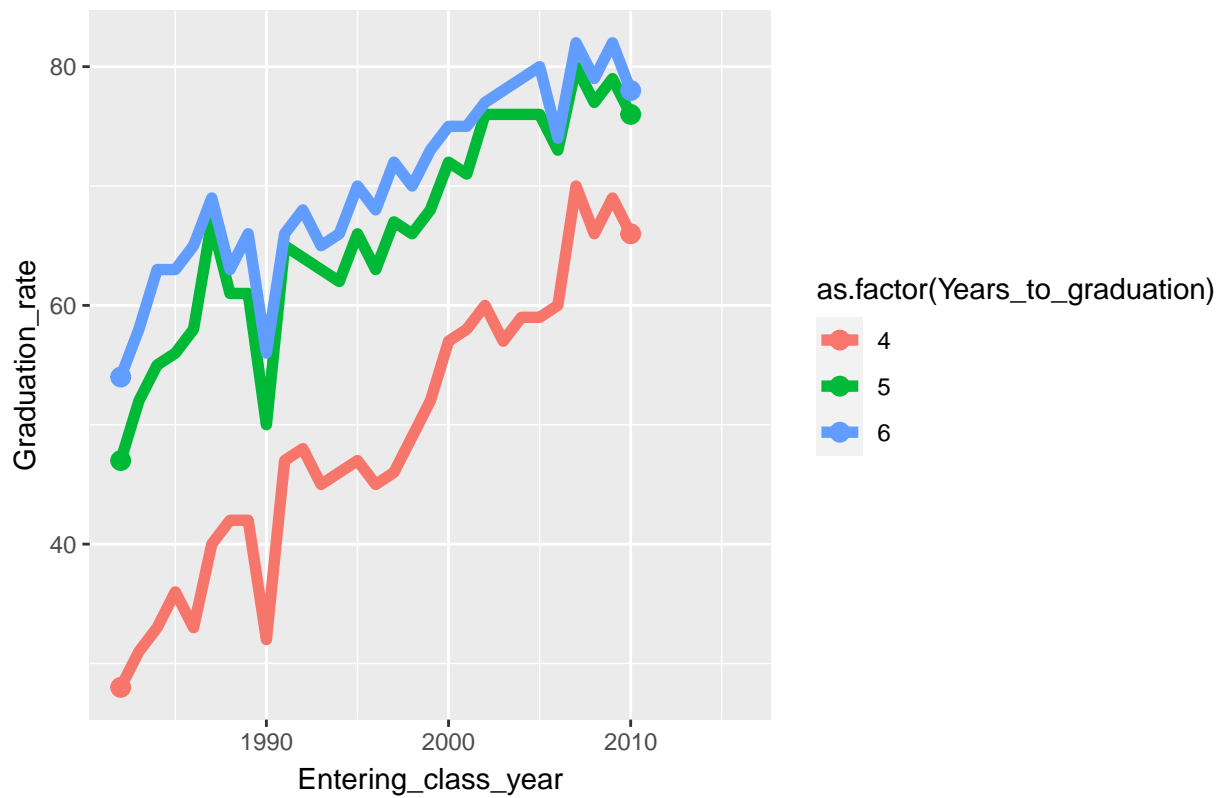
The Year is 2009.



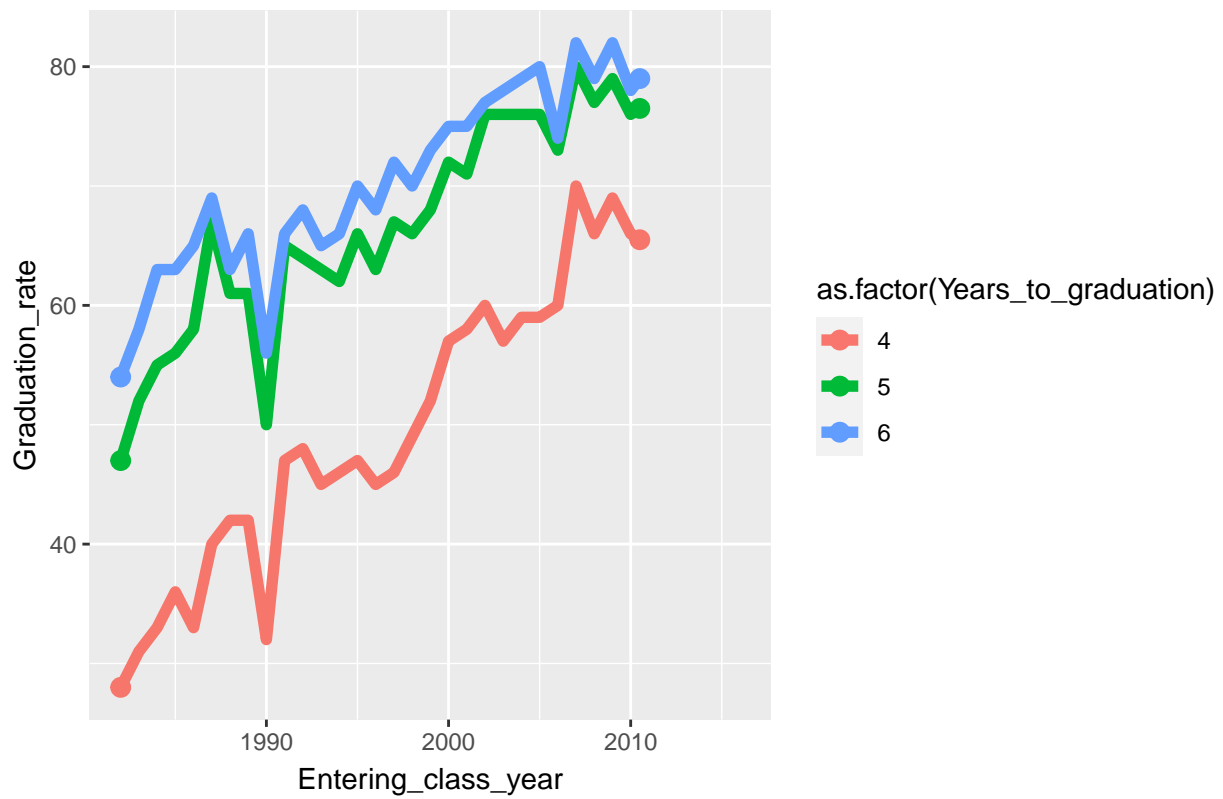
The Year is 2010.



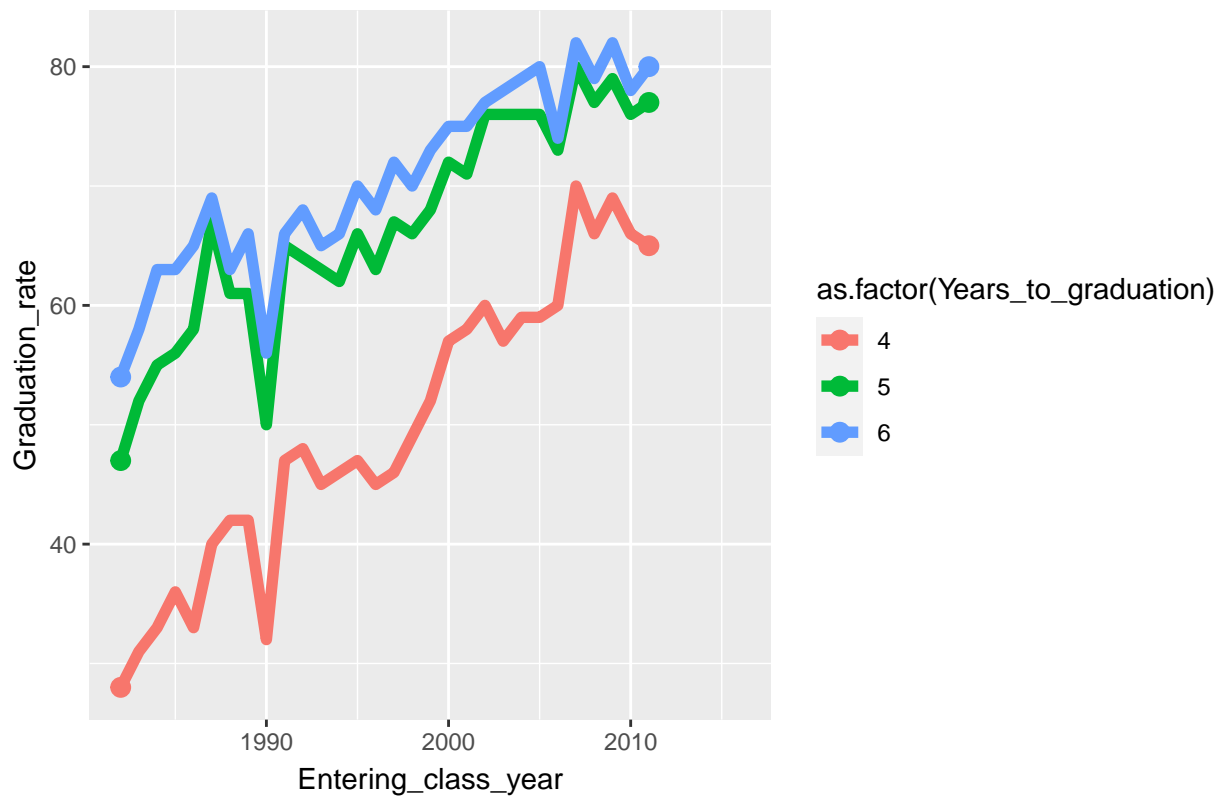
The Year is 2010.



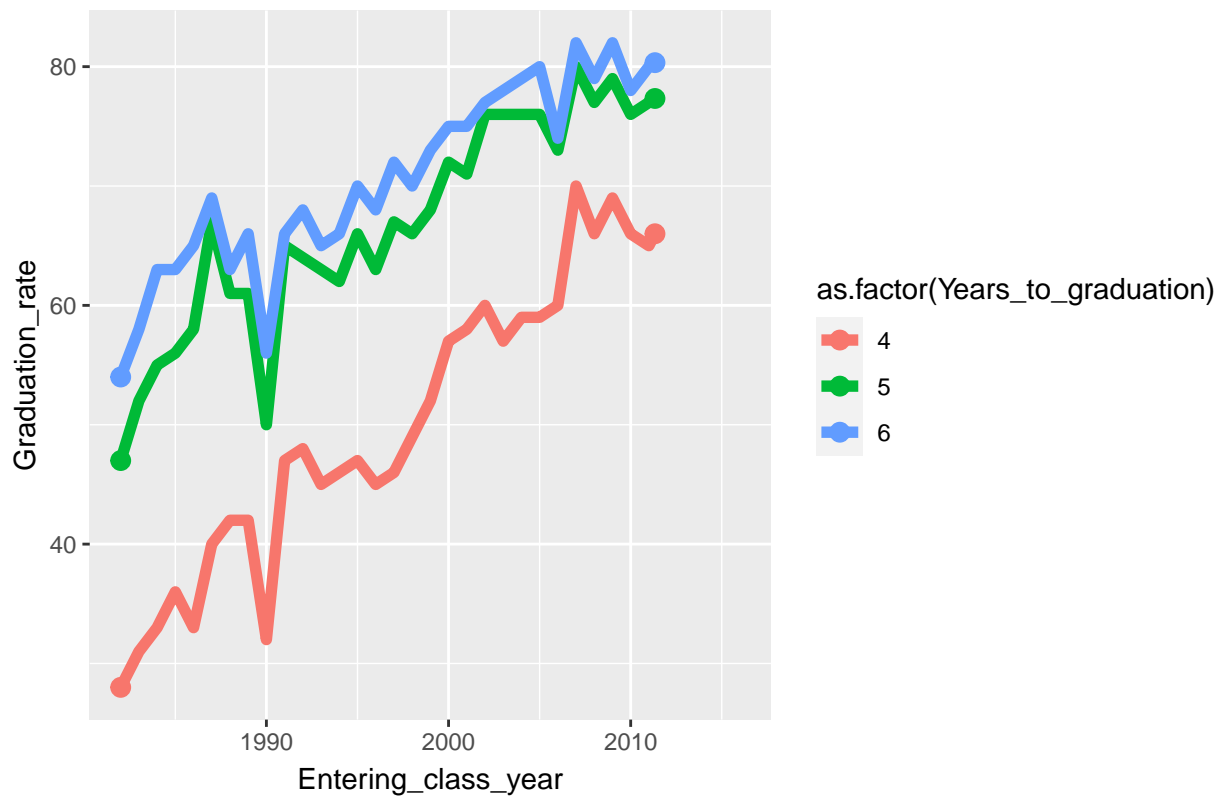
The Year is 2011.



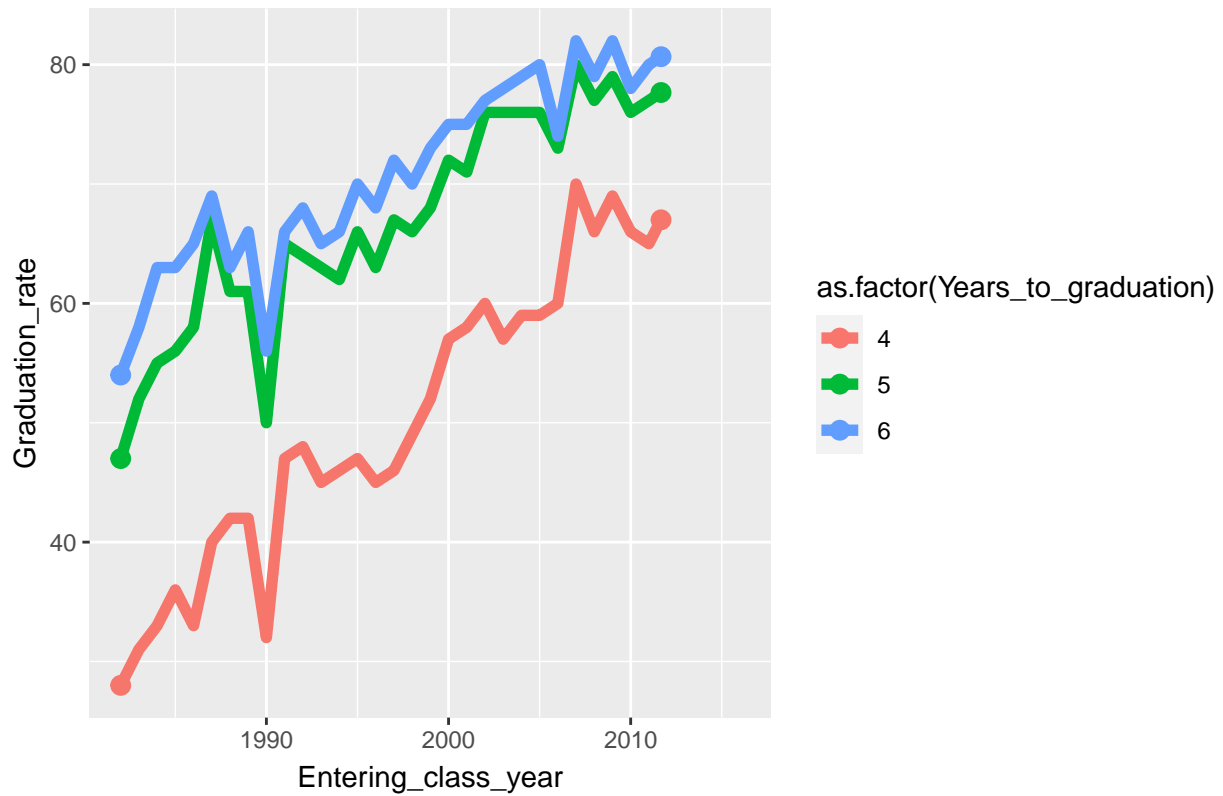
The Year is 2011.



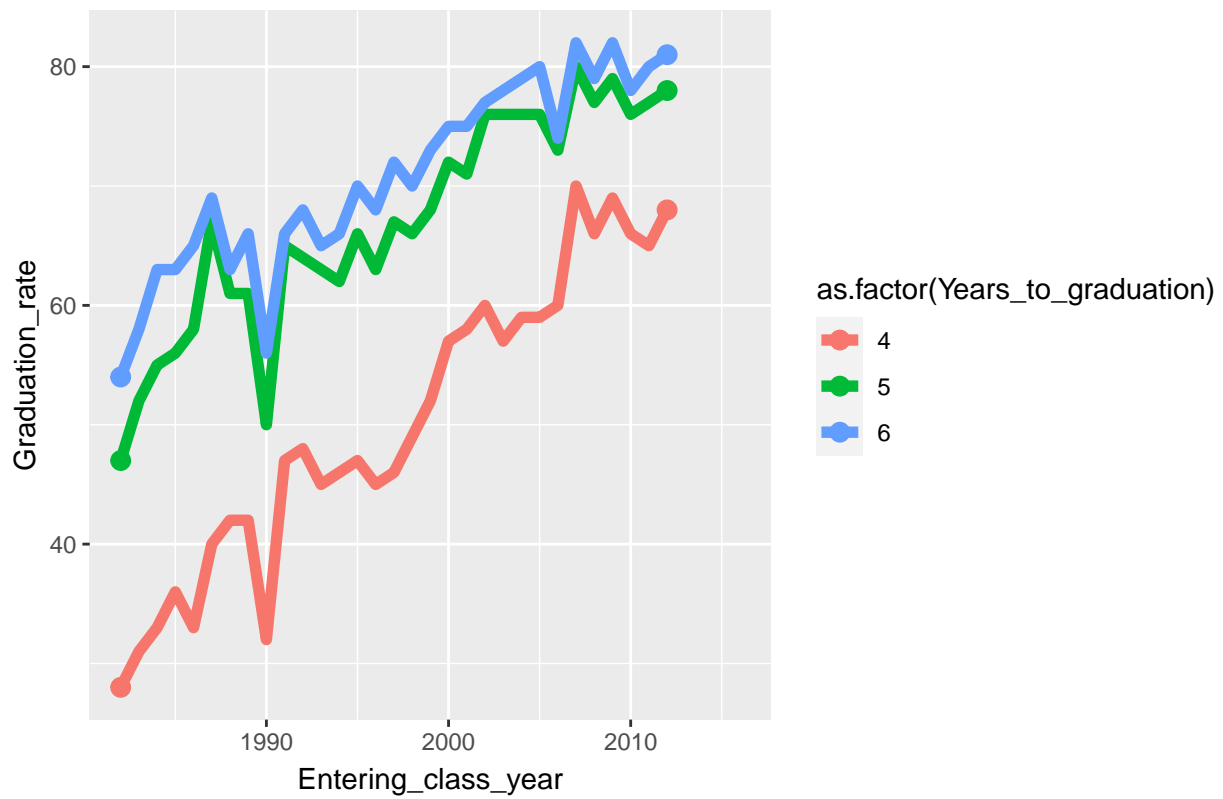
The Year is 2011.



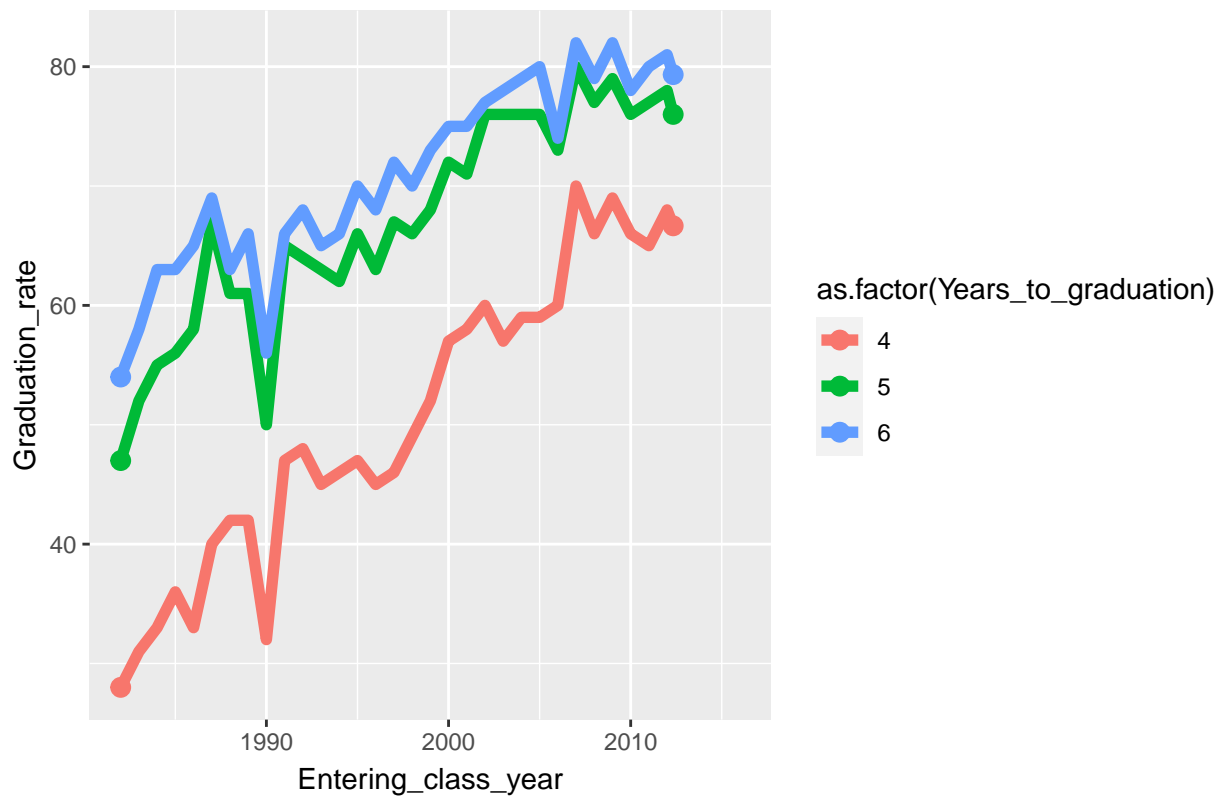
The Year is 2012.



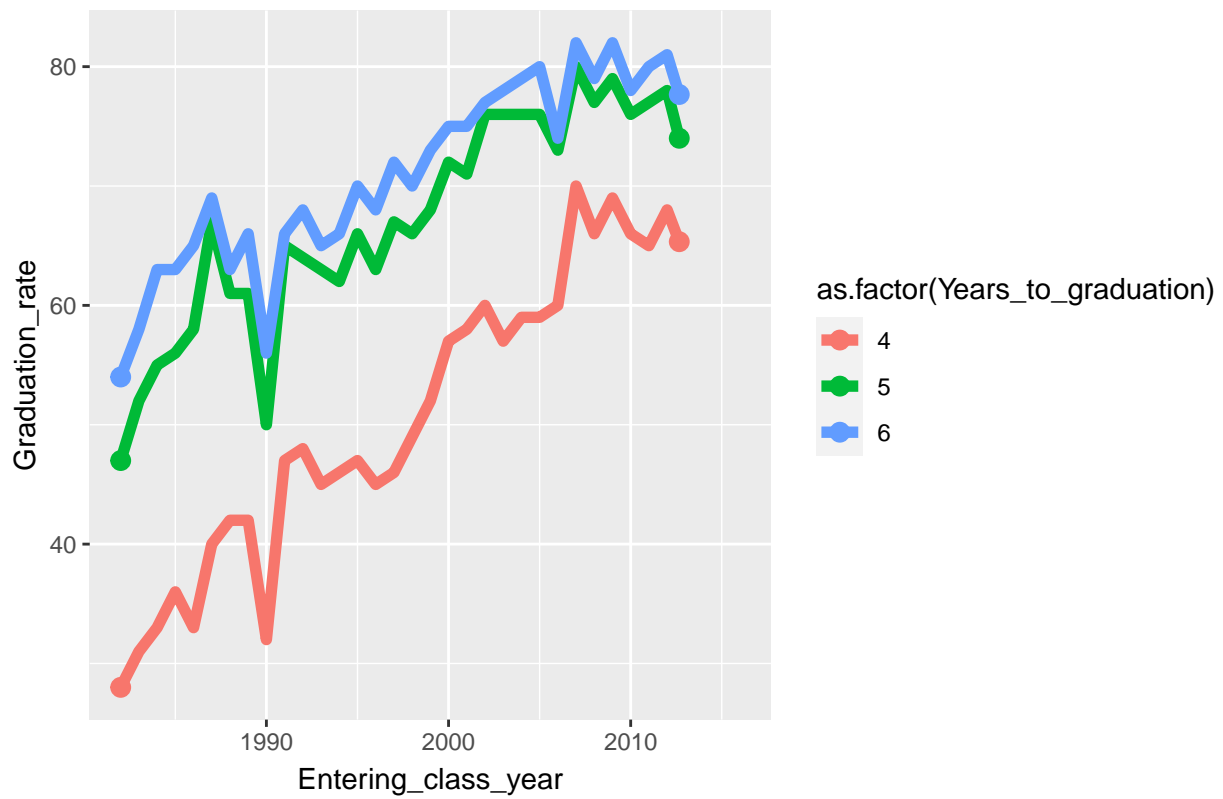
The Year is 2012.



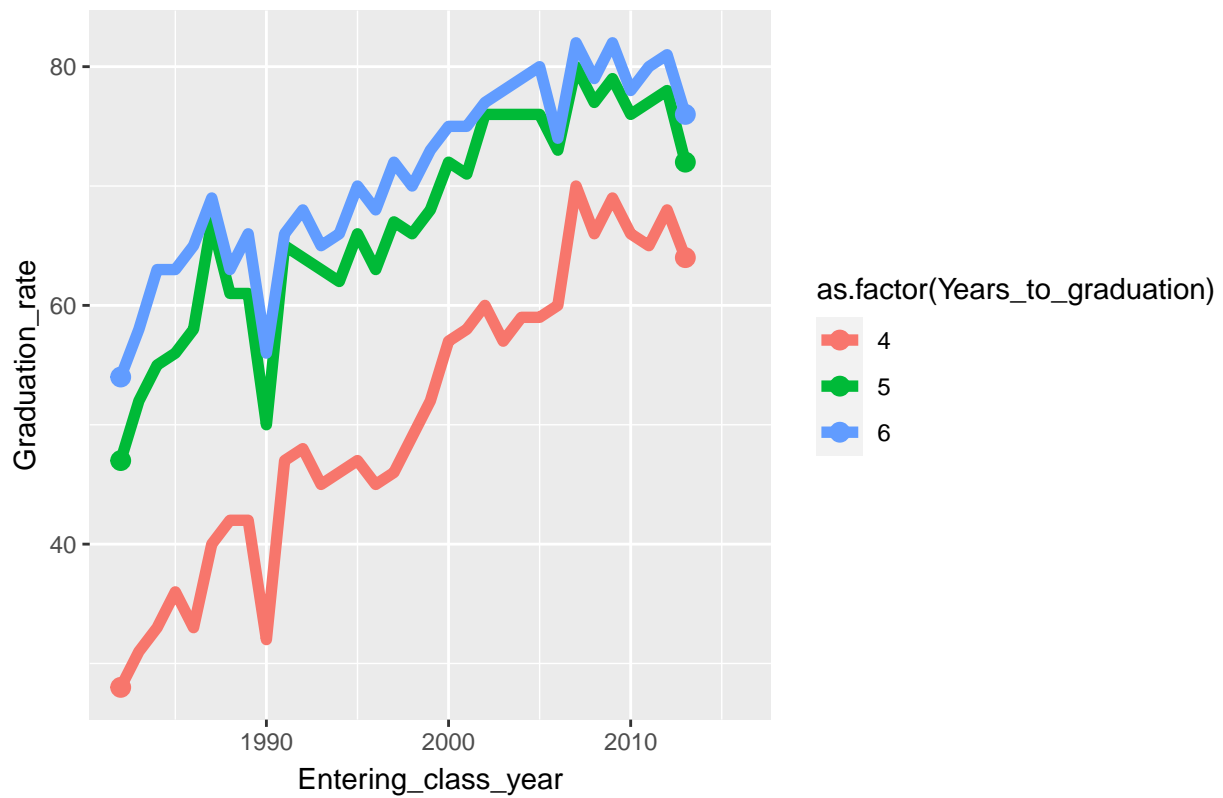
The Year is 2012.



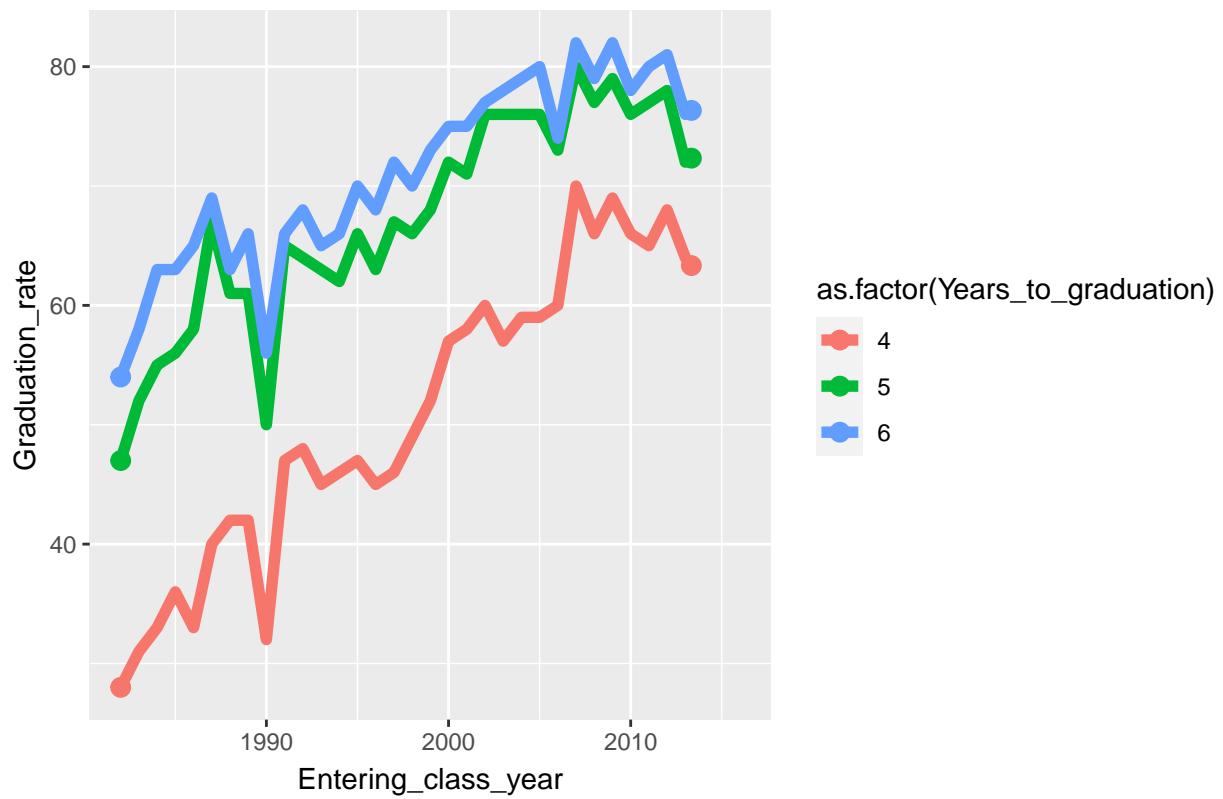
The Year is 2013.



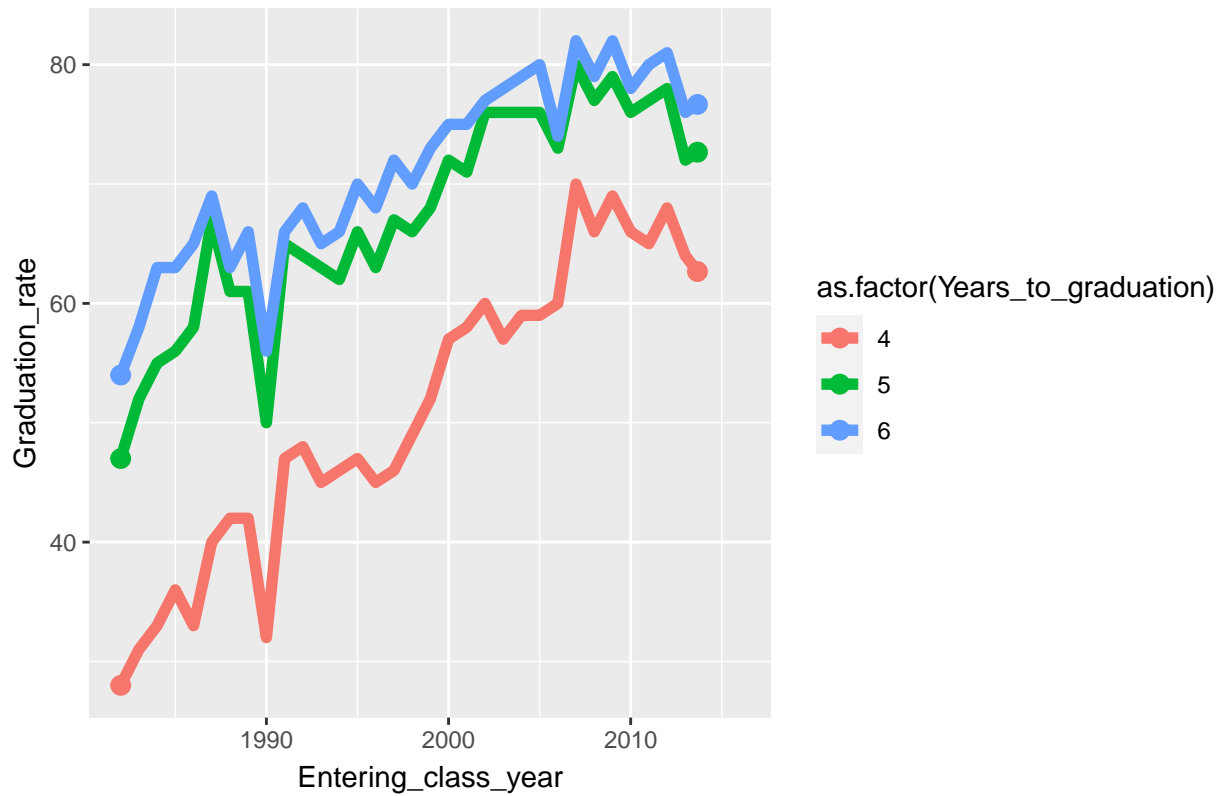
The Year is 2013.



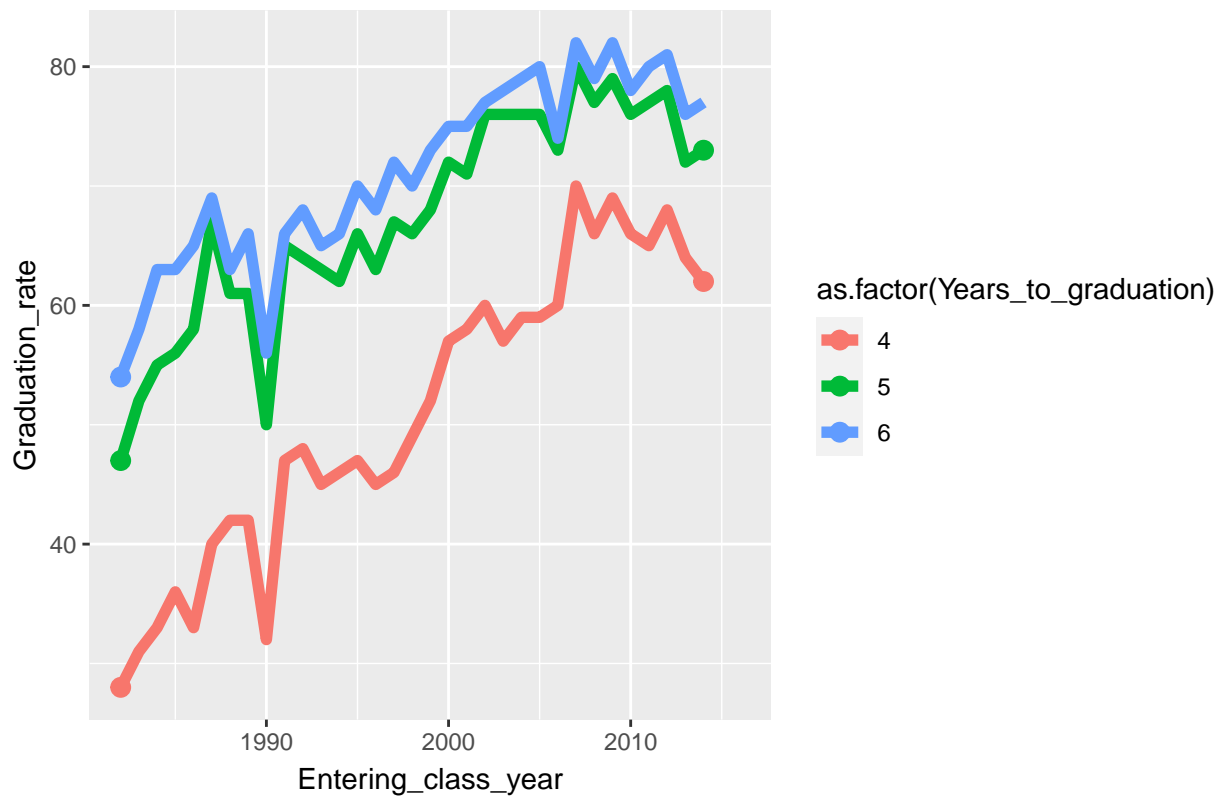
The Year is 2013.



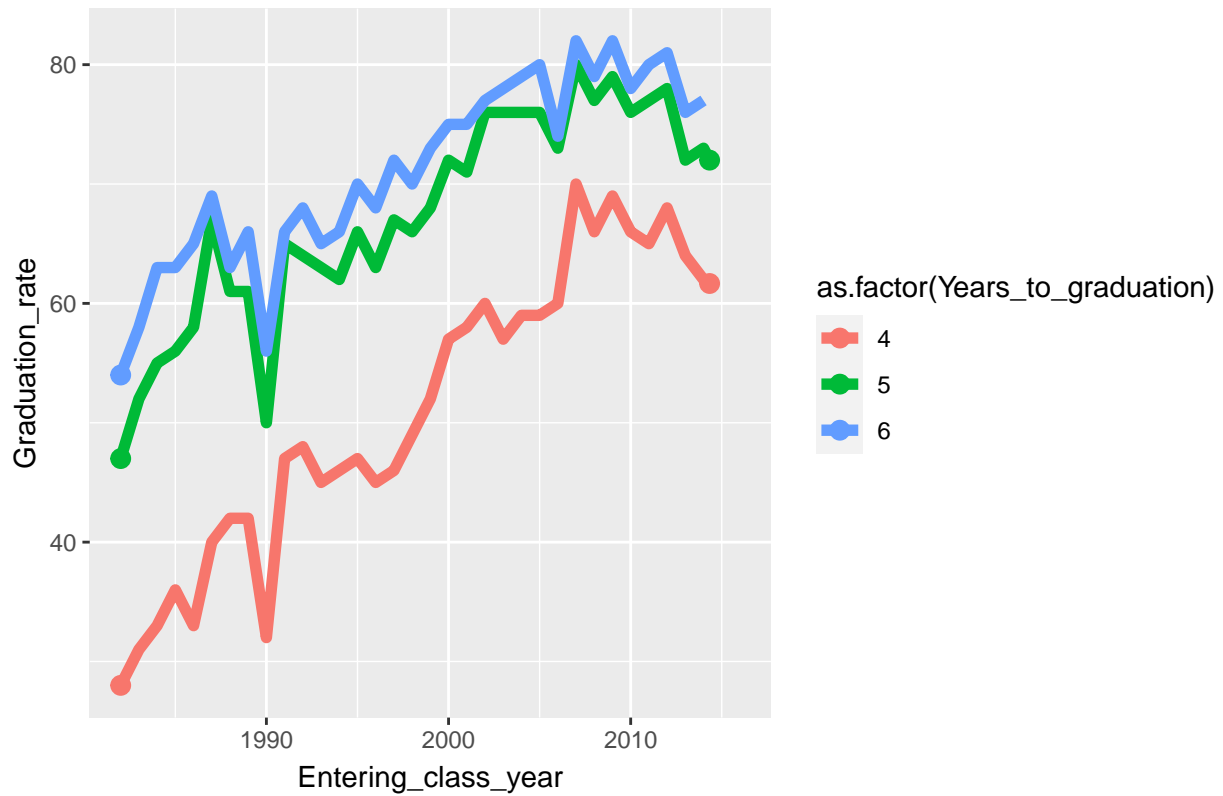
The Year is 2014.



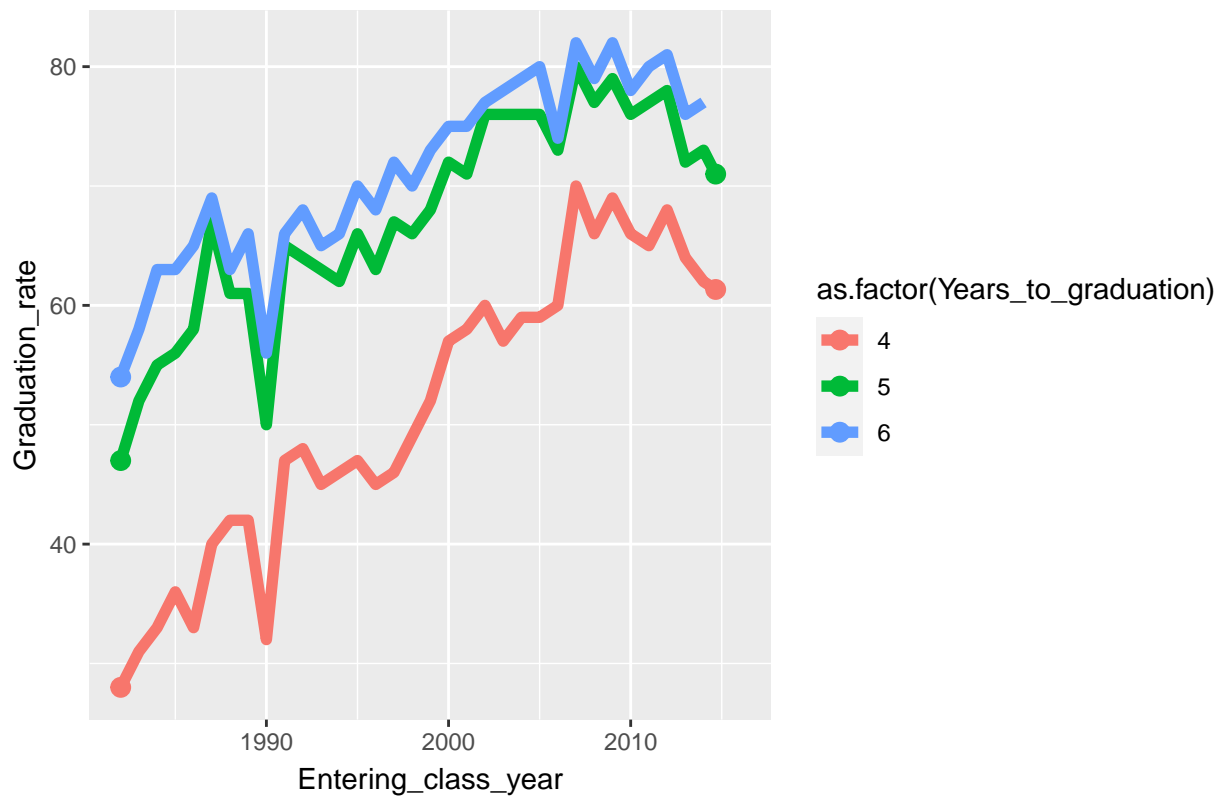
The Year is 2014.



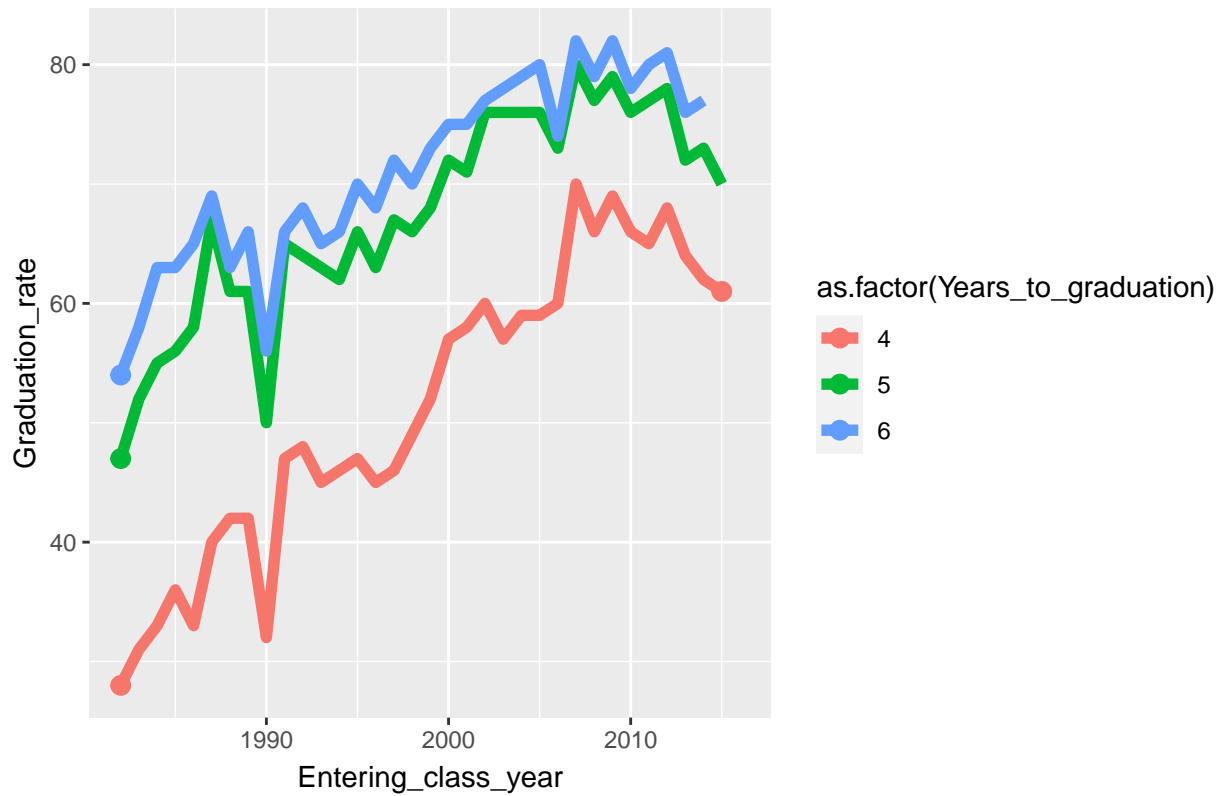
The Year is 2014.



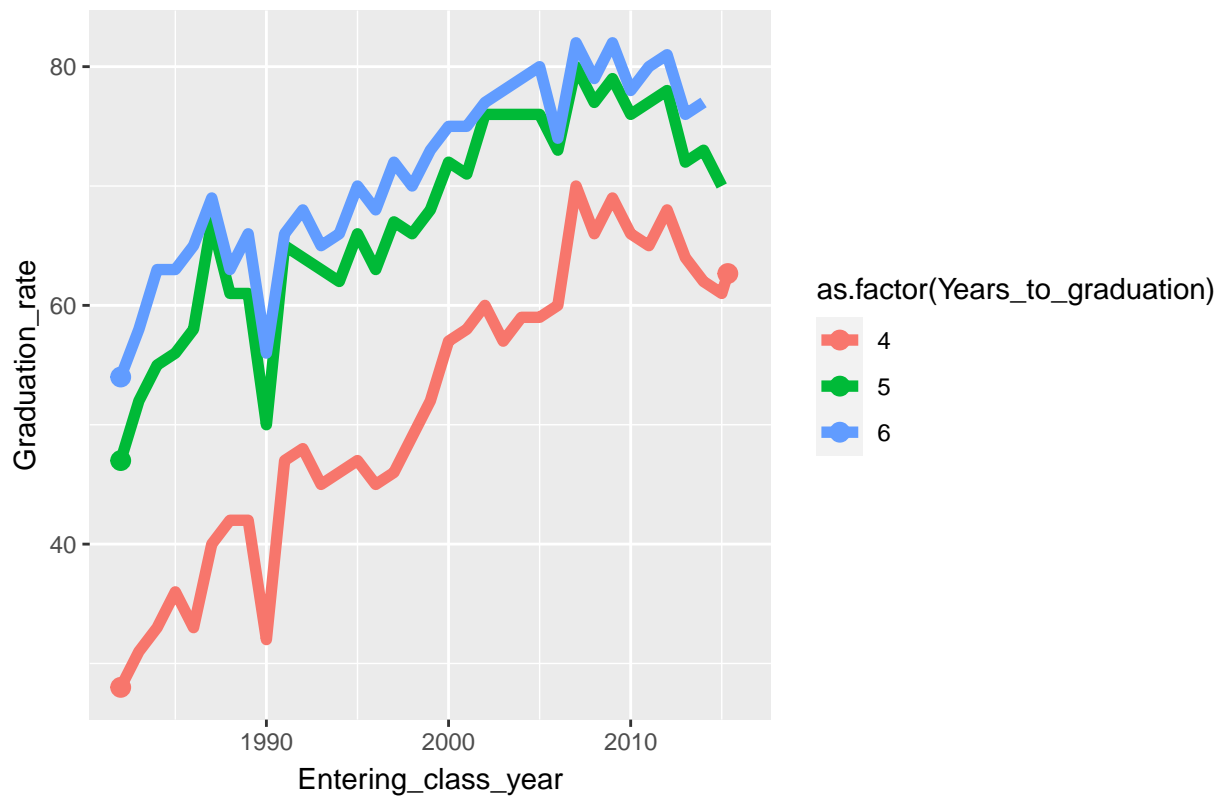
The Year is 2015.

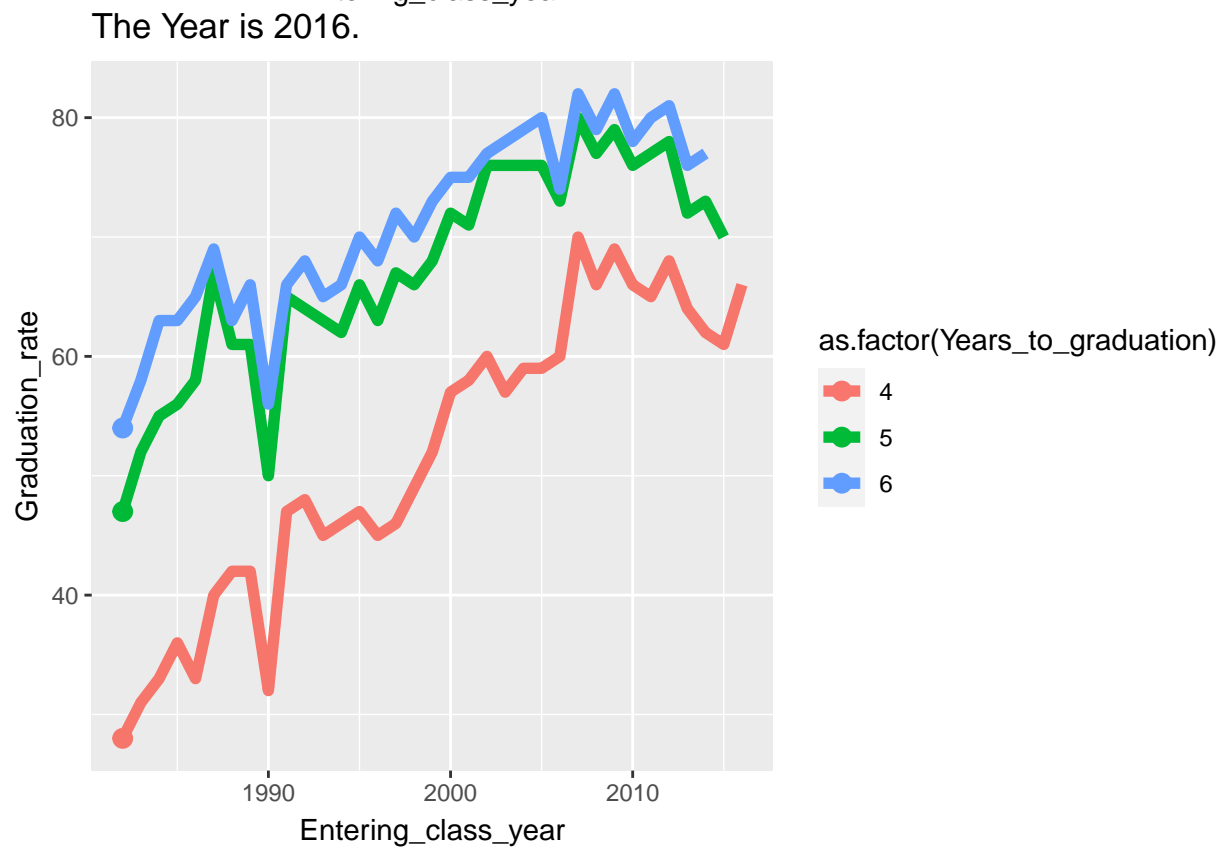
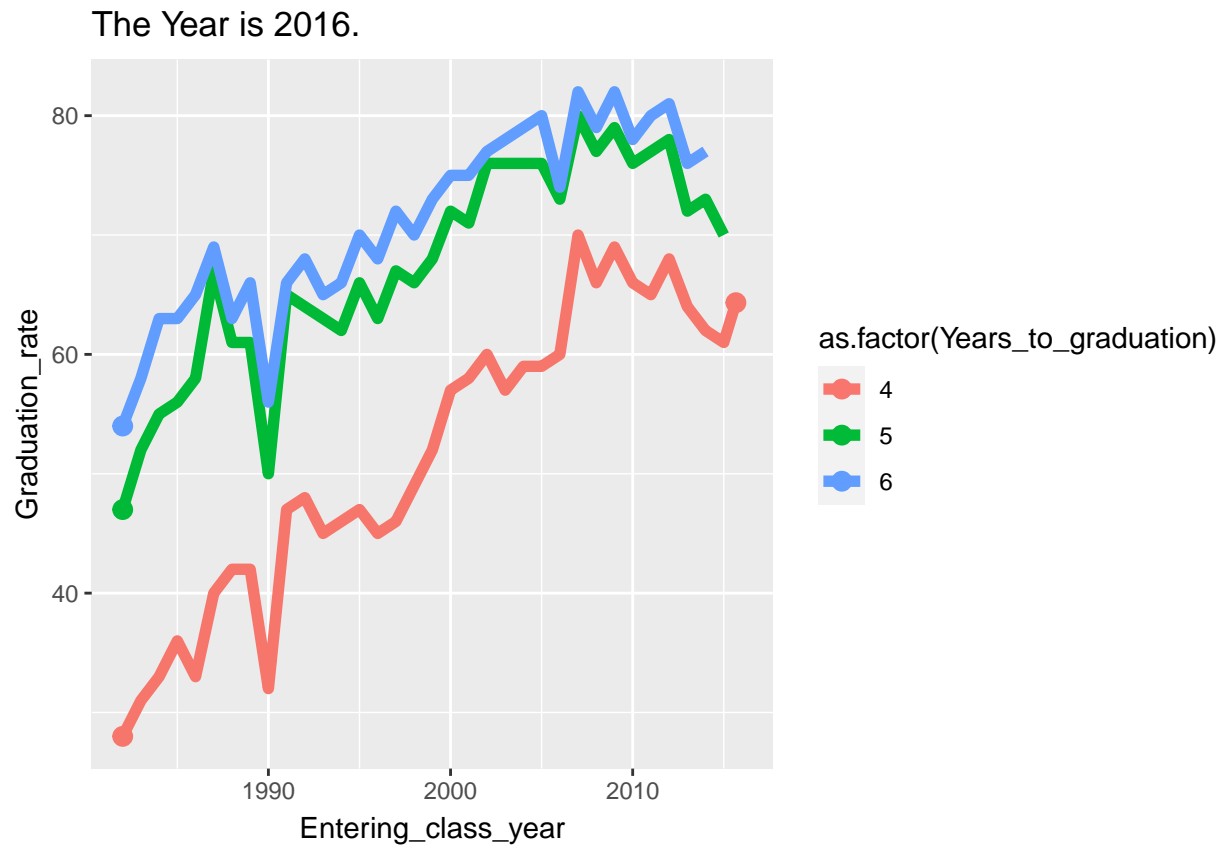


The Year is 2015.



The Year is 2015.





c. The animation improves the plot because:

- It makes the plot more engaging for viewers.
- It accentuates the story that the graduation rate is going up over years for all years to graduation. The animation worsens the plot because:
- It requires a higher level of attention from viewers.
- It might obscure the story, especially for viewers without statistic background.