## all 5 cities: aggregated pollutant measurements (2016)
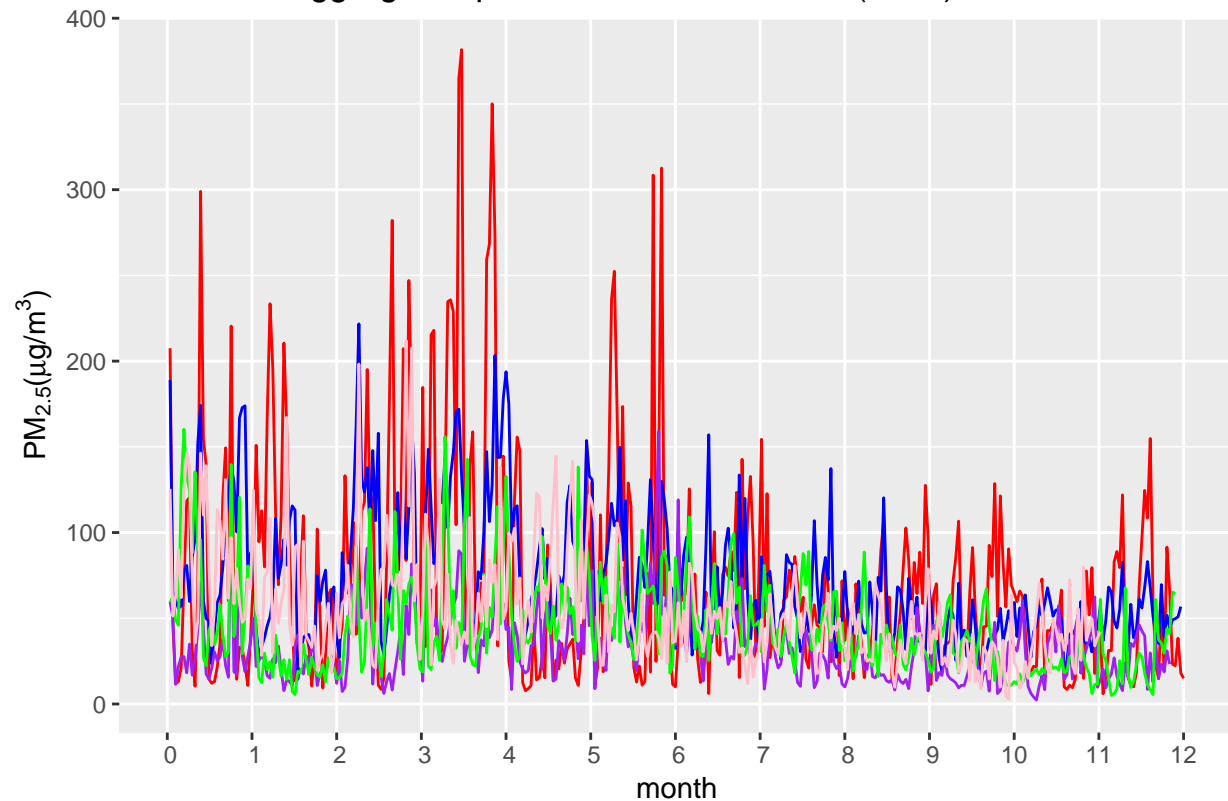


```
five_table <- rbind(beijing_ag, chengdu_ag, guangzhou_ag, shanghai_ag, shenyang_ag)

five_table$city <- rep(0, nrow(five_table))

length(chengdu_ag$Date)
```

```
## [1] 365
```

```
five_table$city[1:365] <- "beijing"
five_table$city[366:731] <- "chengdu"
five_table$city[732:1092] <- "guangzhou"
five_table$city[1093:1455] <- "shanghai"
five_table$city[1456:nrow(five_table)] <- "shenyang"

ggplot() + geom_point(aes(x = c(1:nrow(first_six))/30.3, y = first_six$mean_pm25, col = first_six$month)
```
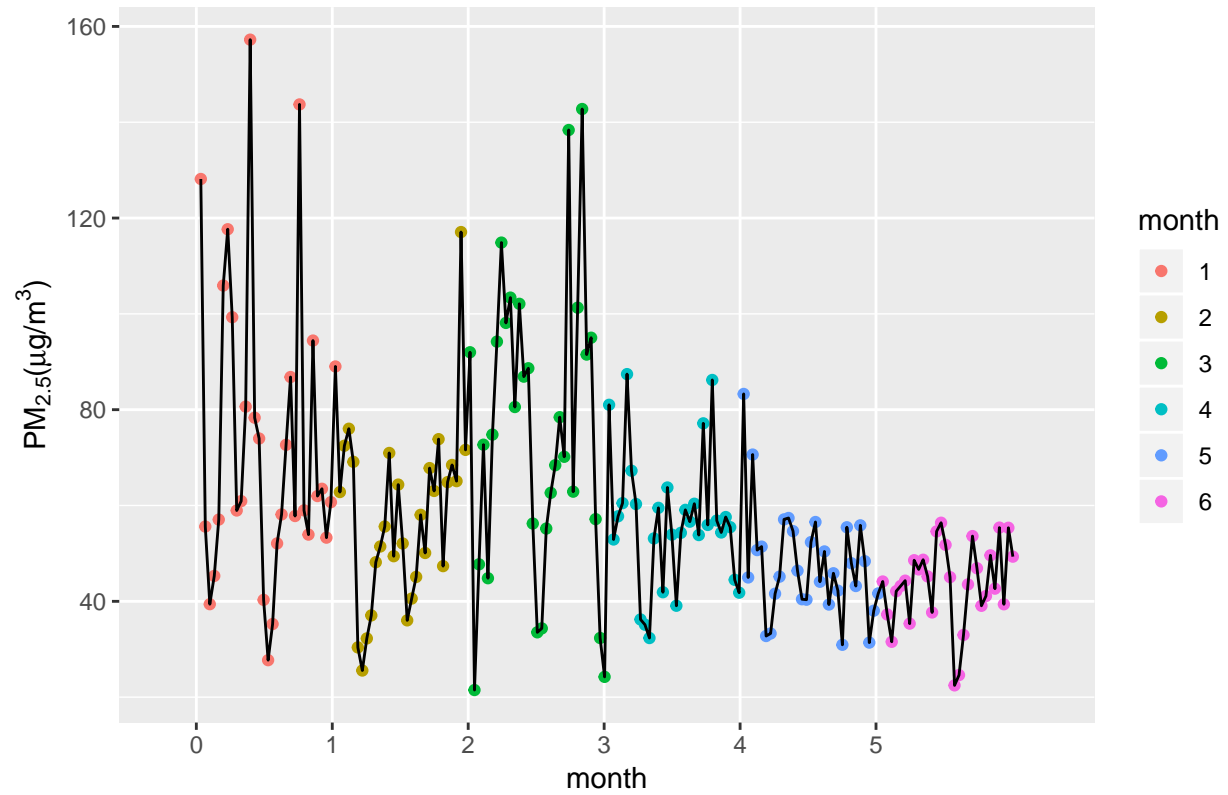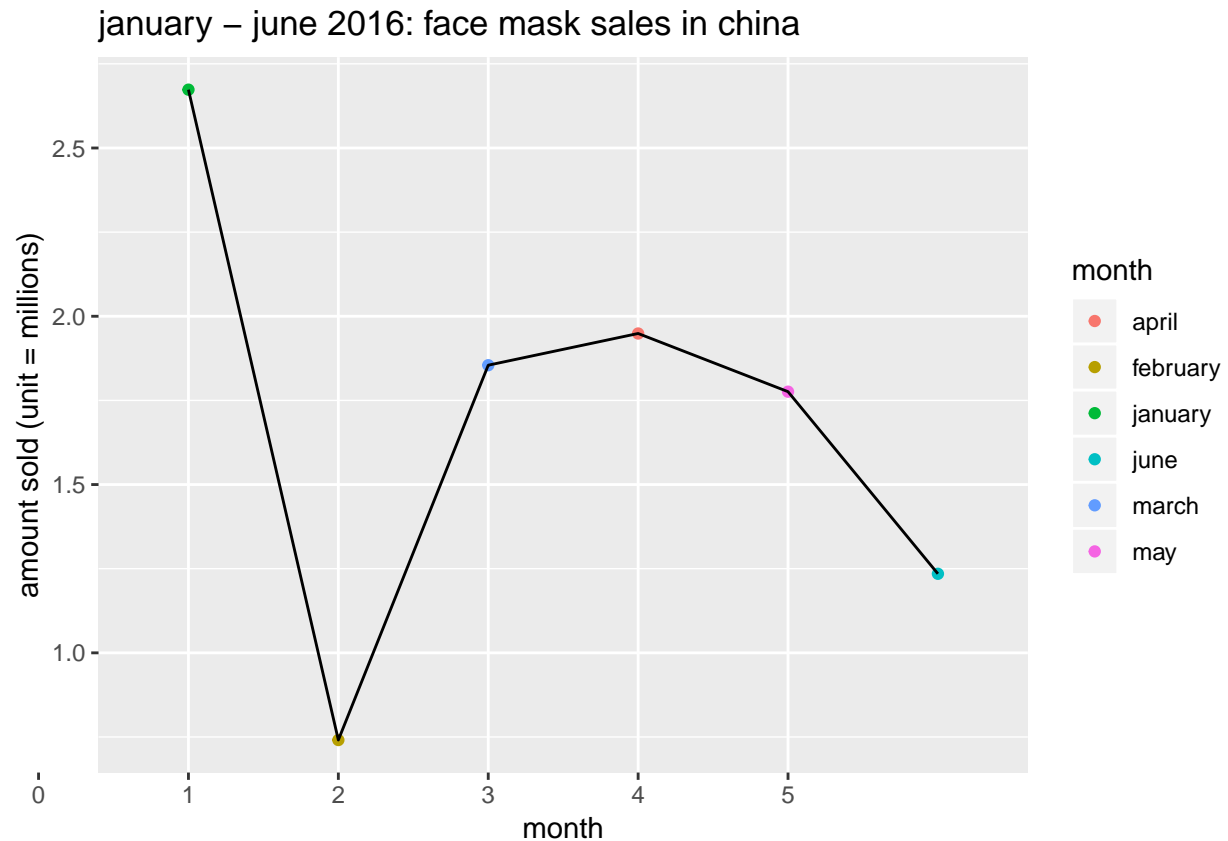
january – june 2016: pollution in china

```
ggplot() + geom_point(aes(x =c(1:6) , y = masks$volume/1000000, col = month)) + geom_line(aes(x = c(1:6
```

## january – june 2016: face mask sales in china



```
?merge
```

```
## Help on topic 'merge' was found in the following packages:
##
##   Package           Library
##   data.table        /Library/Frameworks/R.framework/Versions/3.5/Resources/library
##   raster            /Library/Frameworks/R.framework/Versions/3.5/Resources/library
##   sp                /Library/Frameworks/R.framework/Versions/3.5/Resources/library
##   base              /Library/Frameworks/R.framework/Resources/library
##   lava              /Library/Frameworks/R.framework/Versions/3.5/Resources/library
##
##
## Using the first match ...
```

```
merged <- merge(filtered, counted_frame, by.x = "device_id", by.y = "device_id.Var1")
beijing_ag
```

```
## # A tibble: 366 x 3
##    Date     mean_aqi   moe
##    <chr>       <dbl> <dbl>
## 1 1/1/16      208.   88.8
## 2 1/10/16      47.5  22.0
## 3 1/11/16      11.8   2.91
## 4 1/12/16      13.1   5.81
## 5 1/13/16      21.0  13.3
```

```
##  6 1/14/16     71.6 33.8
##  7 1/15/16     118.  49.1
##  8 1/16/16     120.  24.8
##  9 1/17/16     31.1 34.3
## 10 1/18/16     10.4  5.66
## # ... with 356 more rows
```

```
cities <- merge(beijing_ag, chengdu_ag, by.x = "Date", by.y = "Date")
cities <- merge(cities, guangzhou_ag, by.x = "Date", by.y = "Date")

cities <- data.frame(
  Date = cities$Date,
  beijing = cities$mean_aqi.x,
  chengdu = cities$mean_aqi.y,
  guangzhou_ag = cities$mean_aqi
)

cities <- merge(cities, shanghai_ag, by.x = "Date", by.y = "Date")
cities <- merge(cities, shenyang_ag, by.x = "Date", by.y = "Date")

cities[1:10, -9]
```

```
##        Date     beijing    chengdu guangzhou_ag mean_aqi.x      moe.x
## 1   1/1/16 207.50000 188.95833     59.58333    59.25000  8.935761
## 2  1/10/16  47.50000  66.00000     50.75000    62.54167 10.496290
## 3  1/11/16  11.83333  60.54167     11.47826    48.91667 12.870312
## 4  1/12/16  13.12500  55.66667     21.70833    45.95833  9.493610
## 5  1/13/16  20.95833  57.25000     27.58333   126.09524 30.541619
## 6  1/14/16  71.58333  78.12500     27.66667   160.25000 19.460663
## 7  1/15/16 118.45833  80.81250     18.20000   133.08333 30.371492
## 8  1/16/16 120.41667  59.84615     34.91667   124.37500 13.454634
## 9  1/17/16  31.08333  85.58333     16.79167    88.37500 42.661267
## 10 1/18/16  10.37500  96.87500     26.29167   135.41667 25.286045
##    mean_aqi.y     moe.y
## 1   125.45833 26.97661
## 2    51.25000 24.49534
## 3    63.08333 41.03118
## 4    90.04167 51.72038
## 5    62.04167 36.66592
## 6   123.45833 65.50073
## 7   146.75000 24.10890
## 8   138.87500 16.18725
## 9    72.95833 33.26505
## 10   35.62500 11.26677
```

```
cities <- data.frame(
  Date = cities$Date,
  beijing = cities$beijing,
  chengdu = cities$chengdu,
  guangzhou = cities$guangzhou_ag,
  shanghai = cities$mean_aqi.x,
  shenyang = cities$mean_aqi.y
)
```

```r
#write.csv(cities, "cities.csv")


week <- c()
mult_7s <- c(0:46)

for ( i in 1: 46 ){
  week <- append(week, rep(mult_7s[i], 7))
}
week <- append(week, rep(47, 5))

cities$week <- week

bj <- summarise(
  group_by(cities[,c(1,2,7)], week),
  beijing = mean(beijing)
)

cd <- summarise(
  group_by(cities[,c(1,3,7)], week),
  chengdu = mean(chengdu)
)

gz <- summarise(
  group_by(cities[,c(1,4,7)], week),
  guangzhou = mean(guangzhou)
)

sha <- summarise(
  group_by(cities[,c(1,5,7)], week),
  shanghai = mean(shanghai)
)

she <- summarise(
  group_by(cities[,c(1,6,7)], week),
  shenyang = mean(shenyang)
)

cities_week <- data.frame(
  week = she$week,
  beijing = bj$beijing,
  chengdu = cd$chengdu,
  guangzhou = gz$guangzhou,
  shanghai = sha$shanghai,
  shenyang = she$shenyang
)

#write.csv(cities_week, "cities_week.csv")

cities_week$week <- cities_week$week/3.9


#write.csv(cities_week, "cities_week_mo.csv")
```

```r
event_days <- str_extract(events$timestamp, pattern = "2016-[0-9][0-9]-[0-9][0-9]")

unique_dates <- names(table(event_days))
unique_dates <- unique_dates[2:8]


event_dates <- c("4/30/16", "5/1/16", "5/2/16", "5/3/16", "5/4/16", "5/5/16", "5/6/16", "5/7/16", "5/8/
event_dates <- event_dates[2:8]

date_indices_b <- NULL
date_indices_c <- NULL
date_indices_g <- NULL
date_indices_sha <- NULL
date_indices_she <- NULL

for (i in 1:nrow(beijing_ag)){
  if (beijing_ag$Date[i] %in% event_dates){
    date_indices_b[i] <- i
  } else{
    date_indices_b[i] <- 0
  }
}

for (i in 1:nrow(chengdu_ag)){
  if (chengdu_ag$Date[i] %in% event_dates){
    date_indices_c[i] <- i
  } else{
    date_indices_c[i] <- 0
  }
}

for (i in 1:nrow(guangzhou_ag)){
  if (guangzhou_ag$Date[i] %in% event_dates){
    date_indices_g[i] <- i
  } else{
    date_indices_g[i] <- 0
  }
}

for (i in 1:nrow(shanghai_ag)){
  if (shanghai_ag$Date[i] %in% event_dates){
    date_indices_sha[i] <- i
  } else{
    date_indices_sha[i] <- 0
  }
}

for (i in 1:nrow(shenyang_ag)){
  if (shenyang_ag$Date[i] %in% event_dates){
    date_indices_she[i] <- i
  } else{
    date_indices_she[i] <- 0
  }
```

```r
}
extract_b <- date_indices_b[date_indices_b != 0]
events_beijing <- beijing_ag[extract_b, ]
events_beijing$day <- unique_dates
events_beijing <- events_beijing[c(2:8), ]


extract_c <- date_indices_c[date_indices_c != 0]
events_chengdu <- chengdu_ag[extract_c, ]
events_chengdu$day <- unique_dates
events_chengdu <- events_chengdu[c(2:8), ]

extract_g <- date_indices_g[date_indices_g != 0]
events_guangzhou <- guangzhou_ag[extract_g, ]
events_guangzhou$day <- unique_dates
events_guangzhou <- events_guangzhou[c(2:8), ]

extract_sha <- date_indices_sha[date_indices_sha != 0]
events_shanghai <- shanghai_ag[extract_sha, ]
events_shanghai$day <- unique_dates
events_shanghai <- events_shanghai[c(2:8), ]

extract_she <- date_indices_sha[date_indices_sha != 0]
events_shenyang <- shenyang_ag[extract_she, ]
events_shenyang$day <- unique_dates
events_shenyang <- events_shenyang[c(2:8), ]


events$day <- event_days

time1 <- Sys.time()



# uhhhhh
for (i in 1:100){
  if (events$day[i] %in% events_beijing$day){
    events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day == events$day[i]]
  } else{
    events$beijing[i] <- 0
  }
}
```

```
## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
```

```
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
```

```
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
```

```
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
```

```
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
```

```
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length

## Warning in events$beijing[i] <- events_beijing$mean_aqi[events_beijing$day
## == : number of items to replace is not a multiple of replacement length
```

```r
Sys.time() - time1
```

```
## Time difference of 8.463037 secs
```

```r
# do this differently, assign all to zero
# index out matching date values
# only go thru those


events_beijing$mean_aqi[events_beijing$day == events$day[1]]
```

```
## [1] 61.95833       NA
```

```r
t <- c(1:nrow(events))[events$day == events_beijing$day[2]]
unique_dates
```

```
## [1] "2016-05-01" "2016-05-02" "2016-05-03" "2016-05-04" "2016-05-05"
## [6] "2016-05-06" "2016-05-07"
```

```r
events$date <- str_extract(events$timestamp, pattern = "[0-9]+-[0-9]+-[0-9]+")

event1 <- events[1:500000,]

event2 <- events[500001:1000000,]

event3 <- events[1000001:1500000,]

event4 <- events[1500001:2000000,]

event5 <- events[2000001:2500000,]

event5 <- events[2500001:3032372,]


dates_aqis <- data.frame(
  day = unique_dates,
  aqi_beijing = events_beijing$mean_aqi,
  aqi_chengdu = events_chengdu$mean_aqi,
  aqi_guangzhou = events_guangzhou$mean_aqi,
  aqi_shanghai = events_shanghai$mean_aqi,
  aqi_shenyang = events_shenyang$mean_aqi
)

#write.csv(dates_aqis, "dates_aqis.csv")


randomized_index <- sample(1:nrow(events), 100000, replace = FALSE)

randomized_events <- events[randomized_index, ]

unique_dates
```

```
## [1] "2016-05-01" "2016-05-02" "2016-05-03" "2016-05-04" "2016-05-05"
## [6] "2016-05-06" "2016-05-07"
```

```r
events_short <- events[ ,c(1, 3,4,5,9)]

apr_30 <- events_short[events$day == unique_dates[1], ]
may_1 <- events_short[events$day == unique_dates[2], ]
may_2 <- events_short[events$day == unique_dates[3], ]
may_3 <- events_short[events$day == unique_dates[4], ]
may_4 <- events_short[events$day == unique_dates[5], ]
may_5 <- events_short[events$day == unique_dates[6], ]
may_6 <- events_short[events$day == unique_dates[7], ]
may_7 <- events_short[events$day == unique_dates[8], ]
```

```r
may_8 <- events_short[events$day == unique_dates[9], ]



#write.csv(may_1, "may_1.csv")
#write.csv(may_2, "may_2.csv")
#write.csv(may_3, "may_3.csv")
#write.csv(may_4, "may_4.csv")
#write.csv(may_5, "may_5.csv")
#write.csv(may_6, "may_6.csv")
#write.csv(may_7, "may_7.csv")



#write.csv(dates_aqis, "dates_aqis.csv")



#write.csv(randomized_events, "randomized_events.csv")
ggplot() + geom_point(aes(x = c(1:nrow(events_beijing)), y = events_beijing$mean_aqi)) + geom_line(aes(
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

```
## Warning: Removed 1 rows containing missing values (geom_path).
```
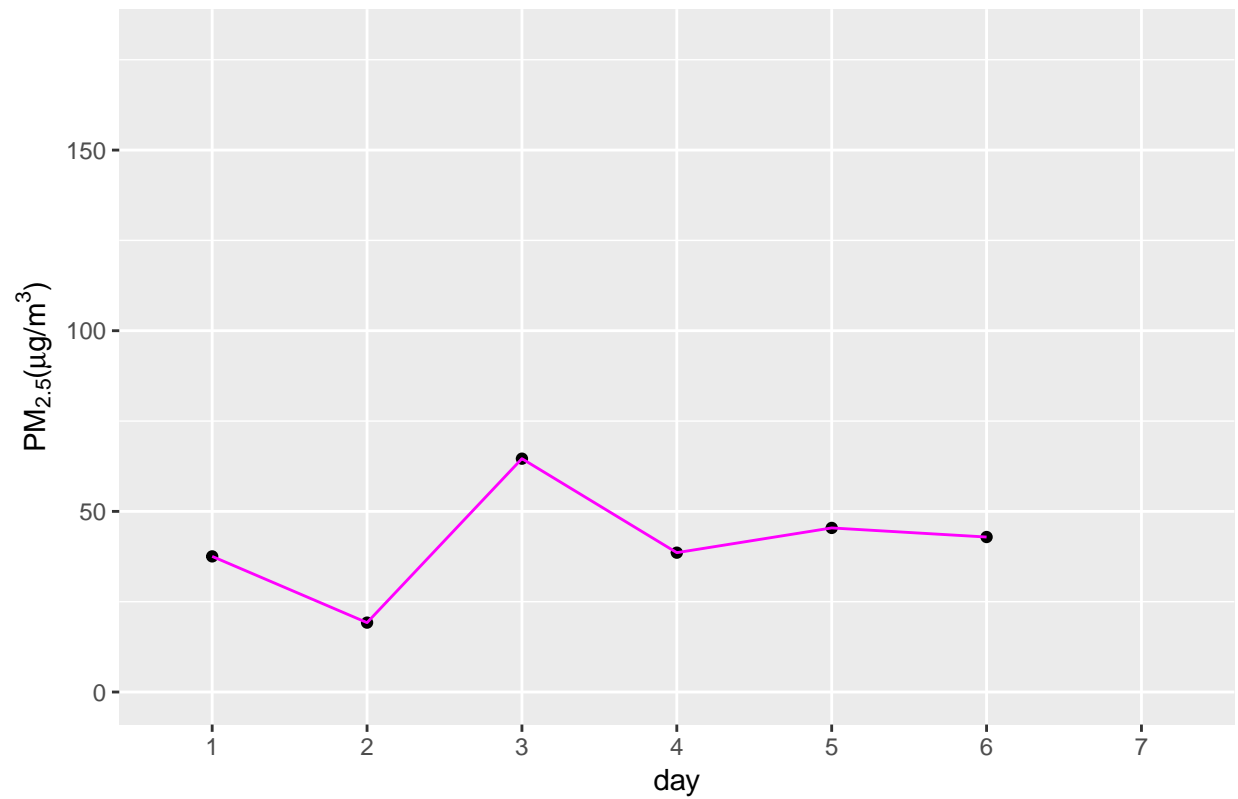


beijing AQ: 5/1–5/7

```r
ggplot() + geom_point(aes(x = c(1:nrow(events_chengdu)), y = events_chengdu$mean_aqi)) + geom_line(aes(
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

```
## Warning: Removed 1 rows containing missing values (geom_path).
```



chengdu AQ: 5/1–5/7

```r
ggplot() + geom_point(aes(x = c(1:nrow(events_guangzhou)), y = events_guangzhou$mean_aqi)) + geom_line(a
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```
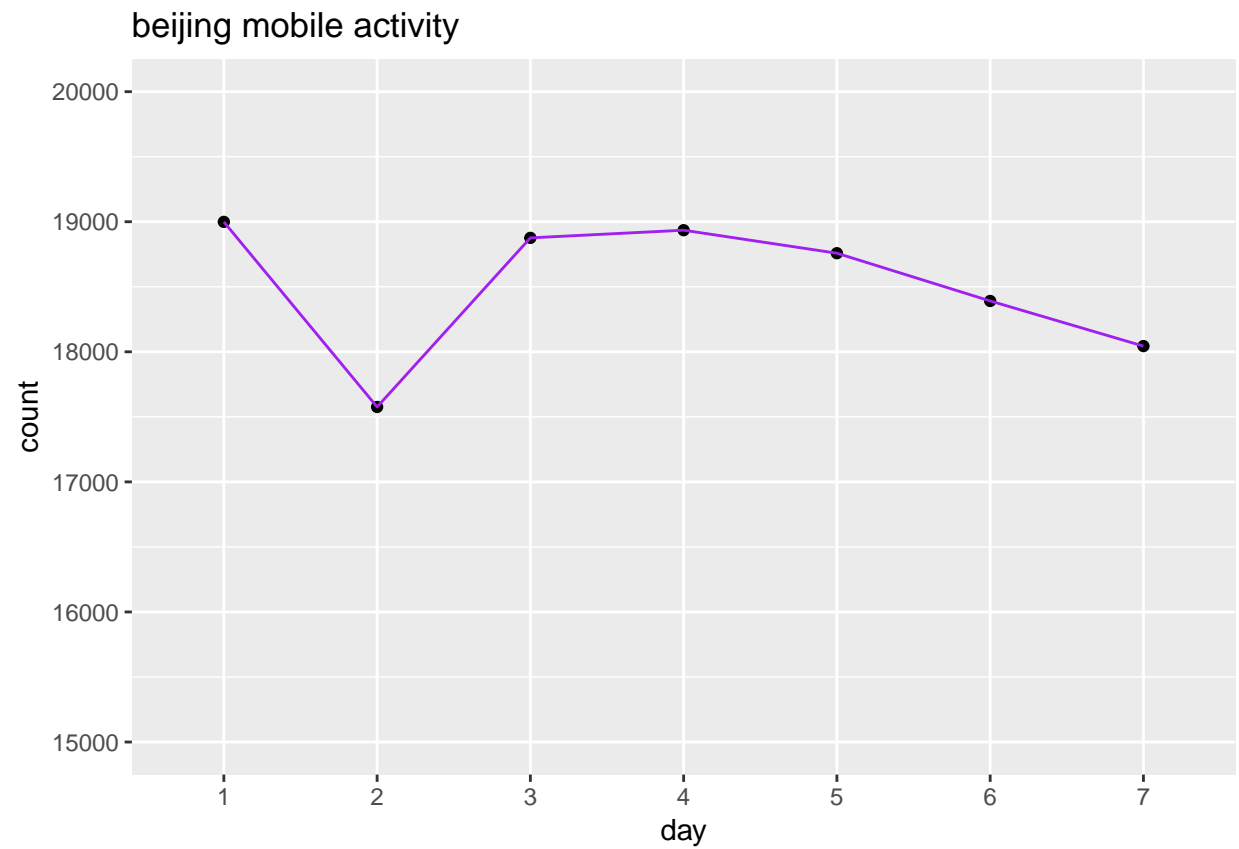
```
## Warning: Removed 1 rows containing missing values (geom_path).
```

guangzhou AQ: 5/1–5/7



```
ggplot() + geom_point(aes(x = c(1:nrow(events_shanghai)), y = events_shanghai$mean_aqi)) + geom_line(aes
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

```
## Warning: Removed 1 rows containing missing values (geom_path).
```

## shanghai AQ: 5/1–5/7



```
ggplot() + geom_point(aes(x = c(1:nrow(events_shenyang)), y = events_shenyang$mean_aqi)) + geom_line(aes
```

```
## Warning: Removed 1 rows containing missing values (geom_point).
```

```
## Warning: Removed 1 rows containing missing values (geom_path).
```
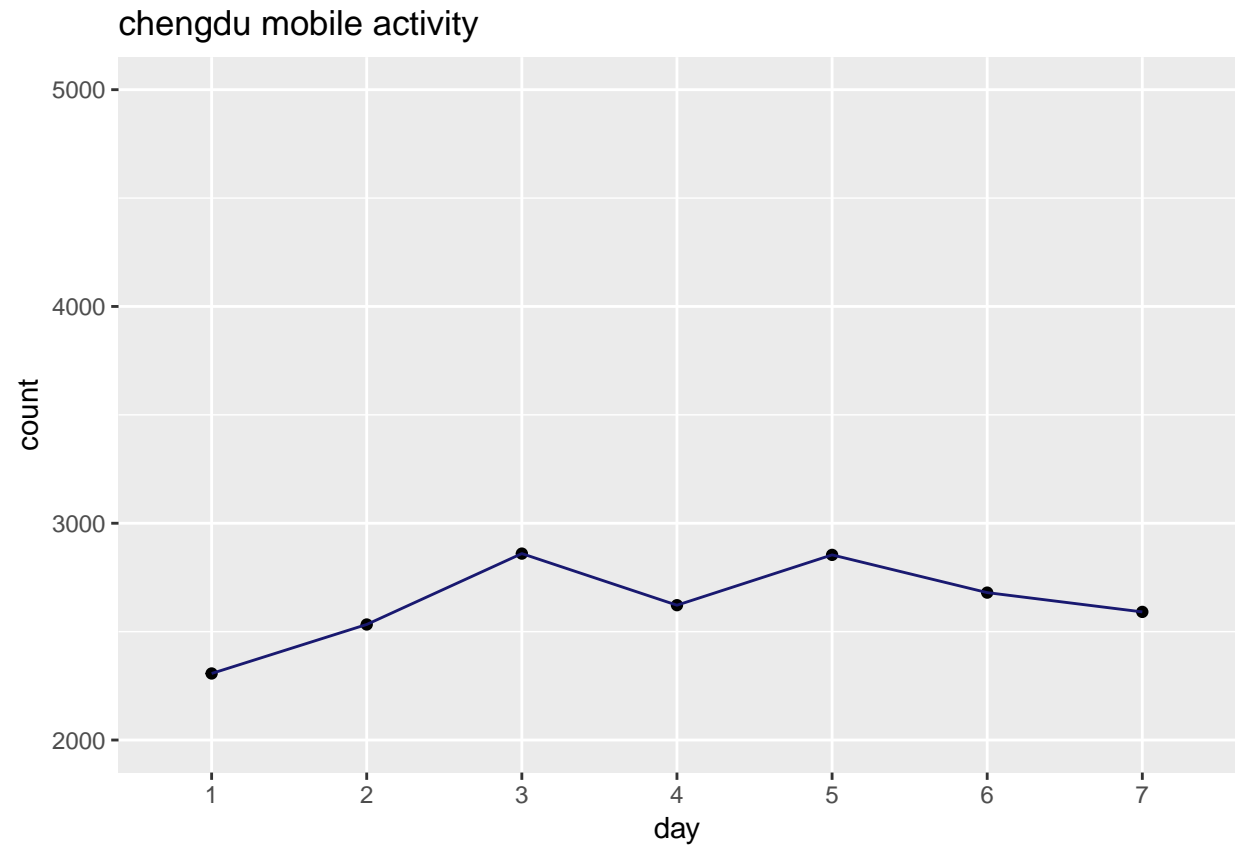
shenyang AQ: 5/1–5/7



## ACTIVITY PLOTS BY CITY

```
activity <- read.csv("~/Documents/caL/2019/cyplan101/projects/assignment3/activity_counts.csv")

ggplot() + geom_point(aes(x = c(1:7), y = activity$beijing)) + geom_line(aes(x = c(1:7), y = activity$be
```

## beijing mobile activity



```
ggplot() + geom_point(aes(x = c(1:7), y = activity$chengdu)) + geom_line(aes(x = c(1:7), y = activity$ch
```

## chengdu mobile activity



```
ggplot() + geom_point(aes(x = c(1:7), y = activity$guangzhou)) + geom_line(aes(x = c(1:7), y = activity$
```