

```
START.TIME <- Sys.time()
knitr::opts_chunk$set(fig.show = 'hide',
  fig.width = 8.4,
  fig.height = 5,
  out.width = "8.4in")
```

```
library(data.table)
library(bit64)

## Loading required package: bit
## Attaching package bit
## package:bit (c) 2008-2012 Jens Oehlschlaegel (GPL-2)
## creators: bit bitwhich
## coercion: as.logical as.integer as.bit as.bitwhich which
## operator: ! & | xor != ==
## querying: print length any all min max range sum summary
## bit access: length<- [ [<- [[ [[<-
## for more help type ?bit
##
## Attaching package: 'bit'
##
## The following object is masked from 'package:data.table':
##
##   setattr
##
## The following object is masked from 'package:base':
##
##   xor
##
## Attaching package bit64
## package:bit64 (c) 2011-2012 Jens Oehlschlaegel (GPL-2 with commercial restrictions)
## creators: integer64 seq :
## coercion: as.integer64 as.vector as.logical as.integer as.double as.character
as.bin
## logical operator: ! & | xor != == < <= >= >
## arithmetic operator: + - * / %/% %% ^
## math: sign abs sqrt log log2 log10
## math: floor ceiling trunc round
## querying: is.integer64 is.vector [is.atomic} [length] is.na format print
## aggregation: any all min max range sum prod
## cumulation: diff cummin cummax cumsum cumprod
## access: length<- [ [<- [[ [[<-
## combine: c rep cbind rbind as.data.frame
```

```

## for more help type ?bit64
##
## Attaching package: 'bit64'
##
## The following object is masked from 'package:bit':
##
##   still.identical
##
## The following objects are masked from 'package:base':
##
##   %in%, :, is.double, match, order, rank

library(dplyr)

##
## Attaching package: 'dplyr'
##
## The following objects are masked from 'package:data.table':
##
##   between, last
##
## The following objects are masked from 'package:stats':
##
##   filter, lag
##
## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(magrittr)
library(ggplot2)
theme_update(text = element_text(family = "",
                                  face = "plain", colour = "black", size = 20, lineheight
                                  hjust = 0.5, vjust = 0.5, angle = 0, margin = margin(),
                                  debug = FALSE))

library(nycflights13) # for airports
nycflights.airports <- airports
library(fasttime)
library(grattan)

## Loading required package: devEMF
##
## Attaching package: 'grattan'
##
## The following object is masked from 'package:datasets':

```

```
##
```

```
##      Orange
```

```
pre2008_flights <-
```

```
  rbindlist(lapply(list.files(path = "../flights/1987-2008/",
                              pattern = "csv$",
                              full.names = TRUE), fread))
```

```
pre2008.names <-
```

```
  names(pre2008_flights)
```

```
read_and_report <-
```

```
  function(filename){
    year <- gsub("^.*(2[0-9]{3}).{3,4}csv$", "\\1", filename)
    if(grepl("1.csv", filename, fixed = TRUE))
      cat(year)
    fread(filename, select = pre2008.names, showProgress = FALSE)
  }
```

```
gc(1,1)
```

```
post2008_flights <-
```

```
  rbindlist(lapply(list.files(path = "../flights", recursive = TRUE, pattern = "2[0-9]{3}*",
                              full.names = TRUE),
                    read_and_report))
```

```
flights <- rbindlist(list(pre2008_flights, post2008_flights), use.names = TRUE)
```

```
readr::write_csv(flights, path = "../1987-2015-On-Time-Performance.csv")
```

```
Sys.time()
```

```
## [1] "2016-01-07 00:49:32 AEDT"
```

```
flights <- fread("../1987-2015-On-Time-Performance.csv")
```

```
##
```

```
Read 0.0% of 165931626 rows
```

```
Read 0.5% of 165931626 rows
```

```
Read 1.0% of 165931626 rows
```

```
Read 1.5% of 165931626 rows
```

```
Read 2.1% of 165931626 rows
```

```
Read 2.6% of 165931626 rows
```

```
Read 3.1% of 165931626 rows
```

Read 3.7% of 165931626 rows
Read 4.2% of 165931626 rows
Read 4.7% of 165931626 rows
Read 5.2% of 165931626 rows
Read 5.7% of 165931626 rows
Read 6.3% of 165931626 rows
Read 6.8% of 165931626 rows
Read 7.4% of 165931626 rows
Read 7.9% of 165931626 rows
Read 8.5% of 165931626 rows
Read 9.0% of 165931626 rows
Read 9.6% of 165931626 rows
Read 10.1% of 165931626 rows
Read 10.7% of 165931626 rows
Read 11.2% of 165931626 rows
Read 11.7% of 165931626 rows
Read 12.2% of 165931626 rows
Read 12.7% of 165931626 rows
Read 13.2% of 165931626 rows
Read 13.8% of 165931626 rows
Read 14.3% of 165931626 rows
Read 14.8% of 165931626 rows
Read 15.3% of 165931626 rows
Read 15.8% of 165931626 rows
Read 16.4% of 165931626 rows
Read 16.9% of 165931626 rows
Read 17.4% of 165931626 rows
Read 17.9% of 165931626 rows
Read 18.4% of 165931626 rows
Read 18.9% of 165931626 rows
Read 19.5% of 165931626 rows
Read 20.0% of 165931626 rows
Read 20.5% of 165931626 rows
Read 21.0% of 165931626 rows
Read 21.5% of 165931626 rows
Read 22.0% of 165931626 rows
Read 22.6% of 165931626 rows
Read 23.1% of 165931626 rows
Read 23.6% of 165931626 rows
Read 24.1% of 165931626 rows
Read 24.6% of 165931626 rows
Read 25.2% of 165931626 rows
Read 25.7% of 165931626 rows

Read 26.2% of 165931626 rows
Read 26.7% of 165931626 rows
Read 27.2% of 165931626 rows
Read 27.7% of 165931626 rows
Read 28.3% of 165931626 rows
Read 28.8% of 165931626 rows
Read 29.3% of 165931626 rows
Read 29.8% of 165931626 rows
Read 30.3% of 165931626 rows
Read 30.8% of 165931626 rows
Read 31.3% of 165931626 rows
Read 31.9% of 165931626 rows
Read 32.4% of 165931626 rows
Read 32.9% of 165931626 rows
Read 33.4% of 165931626 rows
Read 33.9% of 165931626 rows
Read 34.4% of 165931626 rows
Read 35.0% of 165931626 rows
Read 35.5% of 165931626 rows
Read 36.0% of 165931626 rows
Read 36.5% of 165931626 rows
Read 37.0% of 165931626 rows
Read 37.5% of 165931626 rows
Read 38.1% of 165931626 rows
Read 38.6% of 165931626 rows
Read 39.1% of 165931626 rows
Read 39.6% of 165931626 rows
Read 40.1% of 165931626 rows
Read 40.6% of 165931626 rows
Read 41.2% of 165931626 rows
Read 41.7% of 165931626 rows
Read 42.2% of 165931626 rows
Read 42.7% of 165931626 rows
Read 43.2% of 165931626 rows
Read 43.8% of 165931626 rows
Read 44.3% of 165931626 rows
Read 44.8% of 165931626 rows
Read 45.3% of 165931626 rows
Read 45.8% of 165931626 rows
Read 46.3% of 165931626 rows
Read 46.9% of 165931626 rows
Read 47.4% of 165931626 rows
Read 47.9% of 165931626 rows

Read 48.4% of 165931626 rows
Read 48.9% of 165931626 rows
Read 49.4% of 165931626 rows
Read 50.0% of 165931626 rows
Read 50.5% of 165931626 rows
Read 51.0% of 165931626 rows
Read 51.5% of 165931626 rows
Read 52.1% of 165931626 rows
Read 52.6% of 165931626 rows
Read 53.1% of 165931626 rows
Read 53.6% of 165931626 rows
Read 54.2% of 165931626 rows
Read 54.7% of 165931626 rows
Read 55.2% of 165931626 rows
Read 55.7% of 165931626 rows
Read 56.2% of 165931626 rows
Read 56.8% of 165931626 rows
Read 57.3% of 165931626 rows
Read 57.8% of 165931626 rows
Read 58.3% of 165931626 rows
Read 58.8% of 165931626 rows
Read 59.3% of 165931626 rows
Read 59.9% of 165931626 rows
Read 60.4% of 165931626 rows
Read 60.9% of 165931626 rows
Read 61.4% of 165931626 rows
Read 61.9% of 165931626 rows
Read 62.5% of 165931626 rows
Read 63.0% of 165931626 rows
Read 63.5% of 165931626 rows
Read 64.0% of 165931626 rows
Read 64.6% of 165931626 rows
Read 65.1% of 165931626 rows
Read 65.6% of 165931626 rows
Read 66.1% of 165931626 rows
Read 66.6% of 165931626 rows
Read 67.1% of 165931626 rows
Read 67.7% of 165931626 rows
Read 68.2% of 165931626 rows
Read 68.7% of 165931626 rows
Read 69.2% of 165931626 rows
Read 69.7% of 165931626 rows
Read 70.3% of 165931626 rows

Read 70.8% of 165931626 rows
Read 71.3% of 165931626 rows
Read 71.8% of 165931626 rows
Read 72.3% of 165931626 rows
Read 72.8% of 165931626 rows
Read 73.3% of 165931626 rows
Read 73.8% of 165931626 rows
Read 74.4% of 165931626 rows
Read 74.9% of 165931626 rows
Read 75.4% of 165931626 rows
Read 75.9% of 165931626 rows
Read 76.4% of 165931626 rows
Read 76.9% of 165931626 rows
Read 77.4% of 165931626 rows
Read 77.9% of 165931626 rows
Read 78.4% of 165931626 rows
Read 78.9% of 165931626 rows
Read 79.4% of 165931626 rows
Read 79.9% of 165931626 rows
Read 80.4% of 165931626 rows
Read 80.9% of 165931626 rows
Read 81.4% of 165931626 rows
Read 81.9% of 165931626 rows
Read 82.4% of 165931626 rows
Read 82.9% of 165931626 rows
Read 83.5% of 165931626 rows
Read 84.0% of 165931626 rows
Read 84.5% of 165931626 rows
Read 85.0% of 165931626 rows
Read 85.5% of 165931626 rows
Read 86.0% of 165931626 rows
Read 86.5% of 165931626 rows
Read 87.0% of 165931626 rows
Read 87.5% of 165931626 rows
Read 88.0% of 165931626 rows
Read 88.5% of 165931626 rows
Read 89.0% of 165931626 rows
Read 89.5% of 165931626 rows
Read 90.0% of 165931626 rows
Read 90.5% of 165931626 rows
Read 91.0% of 165931626 rows
Read 91.5% of 165931626 rows
Read 92.0% of 165931626 rows

```

Read 92.5% of 165931626 rows
Read 93.0% of 165931626 rows
Read 93.5% of 165931626 rows
Read 94.0% of 165931626 rows
Read 94.5% of 165931626 rows
Read 95.0% of 165931626 rows
Read 95.5% of 165931626 rows
Read 96.0% of 165931626 rows
Read 96.5% of 165931626 rows
Read 97.0% of 165931626 rows
Read 97.5% of 165931626 rows
Read 98.0% of 165931626 rows
Read 98.5% of 165931626 rows
Read 99.0% of 165931626 rows
Read 99.5% of 165931626 rows
Read 165931626 rows and 29 (of 29) columns from 15.111 GB file in 00:03:46

```

```

# flights <- readRDS("../1987-2015-On-Time-Performance.rds")

```

```

flightsSanFran <- flights[Origin %in% c("SFO", "OAK") | Dest %in% c("SFO", "OAK")]
sample.frac = 0.2
sample.weight.int = as.integer(round(1/sample.frac))
flights <- flights[sample(.N, .N * sample.frac)]

```

```

# First we want a time for each flight. This is more difficult than it might seem.
# We need to concatenate the Year, Month, and DayofMonth fields, but we also need
# to take into account the various time zones of the airports in the database.
integer.cols <- grep("Time$", names(flights))

```

```

Sys.time()

```

```

## [1] "2016-01-07 00:54:09 AEDT"

```

```

for (j in integer.cols){
  set(flights, j = j, value = as.integer(flights[[j]]))
}

```

```

Sys.time()

```

```

## [1] "2016-01-07 00:54:09 AEDT"

```

```

# See stackoverflow: links and comments under my question

```

```

create_DepDateTime <- function(DT){

```



```

setkey(DT, Year, Month, DayofMonth, DepTime)
unique_dates <- unique(DT[,list(Year, Month, DayofMonth, DepTime)])
unique_dates[,DepDateTime := fastPOSIXct(sprintf("%d-%02d-%02d %s", Year, Month, DayofMonth,
sub("([0-9]{2})([0-9]{2})", "\\1:\\2:00", DepTime),
perl = TRUE)),
tz = "GMT")]

DT[unique_dates]
}

create_ArrDateTime <- function(DT){
  setkey(DT, Year, Month, DayofMonth, ArrTime)
  unique_dates <- unique(DT[,list(Year, Month, DayofMonth, ArrTime)])
  unique_dates[,ArrDateTime := fastPOSIXct(sprintf("%d-%02d-%02d %s", Year, Month, DayofMonth,
sub("([0-9]{2})([0-9]{2})", "\\1:\\2:00", ArrTime),
perl = TRUE)),
tz = "GMT")]

  DT[unique_dates]
}

flights <- create_DepDateTime(flights)
flights <- create_ArrDateTime(flights)
#flights[,`:=`(Year = NULL, Month = NULL, DayofMonth = NULL, DepTime = NULL, ArrTime = NULL),
Sys.time()

## [1] "2016-01-07 00:56:00 AEDT"

```

```

# Now we join it to the airports dataset from nycflights13 to obtain time zone information
Sys.time()

## [1] "2016-01-07 00:56:00 AEDT"

airports <- as.data.table(airports)
airports <- airports[,list(faa, tz)]
gc(1,1)

##           used      (Mb) gc trigger      (Mb) max used      (Mb)
## Ncells   551113    29.5  11554252    617.1   551113    29.5
## Vcells 819139108 6249.6 2325204027 17739.9 819139108 6249.6

setnames(airports, old = c("faa", "tz"), new = c("Origin", "tzOrigin"))
setkey(airports, Origin)
setkey(flights, Origin)
flights <- flights[airports]
setnames(airports, old = c("Origin", "tzOrigin"), new = c("Dest", "tzDest"))

```

```

setkey(flights, Dest)
flights <- flights[airports]
rm(airports)
gc(1,1)

##           used      (Mb) gc trigger      (Mb)  max used      (Mb)
## Ncells    551171    29.5   9243401    493.7    551171    29.5
## Vcells 879030196 6706.5 2325204027 17739.9 879030196 6706.5

# The joins produce NAs when the airports table isn't present in the flights table.
flights <- flights[!is.na(Origin)]
gc(1,1)

##           used      (Mb) gc trigger      (Mb)  max used      (Mb)
## Ncells    551145    29.5   7394720    395.0    551145    29.5
## Vcells 879007270 6706.3 2325204027 17739.9 879007270 6706.3

Sys.time()

## [1] "2016-01-07 00:56:36 AEDT"

```

```

Sys.time()

## [1] "2016-01-07 00:56:36 AEDT"

# setting keys doesn't improve timing
flights[,`:=`(DepDateTimeZulu = DepDateTime - lubridate::hours(tzOrigin))]
flights[,`:=`(ArrDateTimeZulu = ArrDateTime - lubridate::hours(tzDest))]
Sys.time()

## [1] "2016-01-07 01:00:20 AEDT"

```

```

# Flights typically follow a weekly cycle, so we should obtain the week in the dataset.
# Pretty quick!
Sys.time()

## [1] "2016-01-07 01:00:20 AEDT"

setkey(flights, Year, Month, DayofMonth)
unique_dates <- unique(flights)
unique_dates <- unique_dates[,list(Year, Month, DayofMonth)]
unique_dates[,Week := (Year - 1987L) * 52 + data.table::yday(sprintf("%d-%02d-%02d", Year,
unique_dates[,Week := Week - min(Week)]

```

```
flights <- flights[unique_dates]  
Sys.time()  
## [1] "2016-01-07 01:00:37 AEDT"
```

Flights 1987-2015

Hugh P

January 7, 2016

1

There were 164 million flights from 1987-09-30 23:40:00 to 2015-11-01 09:04:00.

2 San Francisco

```
Sys.time()

## [1] "2016-01-07 01:00:37 AEDT"

setkey(flightsSanFran, Year, Month, DayofMonth)
unique_dates <- unique(flightsSanFran)
unique_dates <- unique_dates[,list(Year, Month, DayofMonth)]
unique_dates[,Week := (Year - 1987L) * 52 + data.table::yday(sprintf("%d-%02d-%02d", Year,
unique_dates[,Week := Week - min(Week)]
flightsSanFran <- flightsSanFran[unique_dates]
Sys.time()

## [1] "2016-01-07 01:00:41 AEDT"
```

```
maxN <- function(x, N=2){
  len <- length(x)
  if(N>len){
    warning('N greater than length(x). Setting N=length(x)')
    N <- length(x)
  }
  sort(x,partial=len-N+1)[len-N+1]
}

setkey(unique_dates, Week)
flightsSanFran %>%
  filter(!(Origin %in% c("SFO", "OAK") & Dest %in% c("SFO", "OAK"))) %>%
  mutate(SF_airport = ifelse(Origin %in% c("SFO", "OAK"),
                             Origin,
                             Dest)) %>%
  count(Week, SF_airport) %>%
  group_by(SF_airport) %>%
  mutate(label.text = ifelse(n == maxN(n), paste(" ", SF_airport), NA_character_)) %>%
  setkey(Week) %>%
  data.table::merge.data.table(unique(unique_dates)) %>%
  mutate(Date = fastPOSIXct(sprintf("%d-%02d-%02d", Year, Month, DayofMonth), tz = "GMT")
```

```

    n = n) %>% # not a sample
  ggplot(aes(x = Date, y = n, color = SF_airport, group = SF_airport)) +
  geom_point() +
  geom_text(aes(label = label.text),
            nudge_y = 0.5,
            nudge_x = 1,
            hjust = 0,
            fontface = "bold",
            size = 5) +
  theme(legend.position = "none") +
  geom_line(size = 0.5) +
  #
  geom_vline(xintercept = as.numeric(as.POSIXct("2001-09-11"))) +
  scale_x_datetime(date_breaks = "5 years",
                  date_labels = "%Y",
                  minor_breaks = seq(as.POSIXct("1987-12-31"), as.POSIXct("2014-12-31")))

## Warning: Removed 2920 rows containing missing values (geom_text).

```

```

carriers <- as.data.table(airlines)
if("carrier" %in% names(carriers))
  setnames(carriers, old = "carrier", new = "UniqueCarrier")

setkey(carriers, UniqueCarrier)
set(carriers, j = 1L, value = as.character(carriers[[1L]]))
set(carriers, j = 2L, value = gsub("^([A-Za-z]+)\\s.*$", "\\1", carriers[[2L]]))

flightsSanFran %>%
  filter(Origin %in% c("SFO", "OAK")) %>%
  count(Year, Month, Origin, UniqueCarrier) %>%
  group_by(UniqueCarrier) %>%
  filter(sum(n) > (2015 - 1987) * 12 * 30) %>%
  mutate(Date = Year + (Month - 1)/12) %>%
  setkey(UniqueCarrier) %>%
  merge(carriers) %>%
  ggplot(aes(x = Date, y = n * sample.weight.int, color = name, group = interaction(name,
  geom_smooth(span = 0.25, se = FALSE) +
  geom_text(aes(label = ifelse(Date == max(Date),
                             name,
                             NA_character_),
             vjust = ifelse(name == "Southwest" & Origin == "SFO",
                             -0.5,
                             0.5)),

```

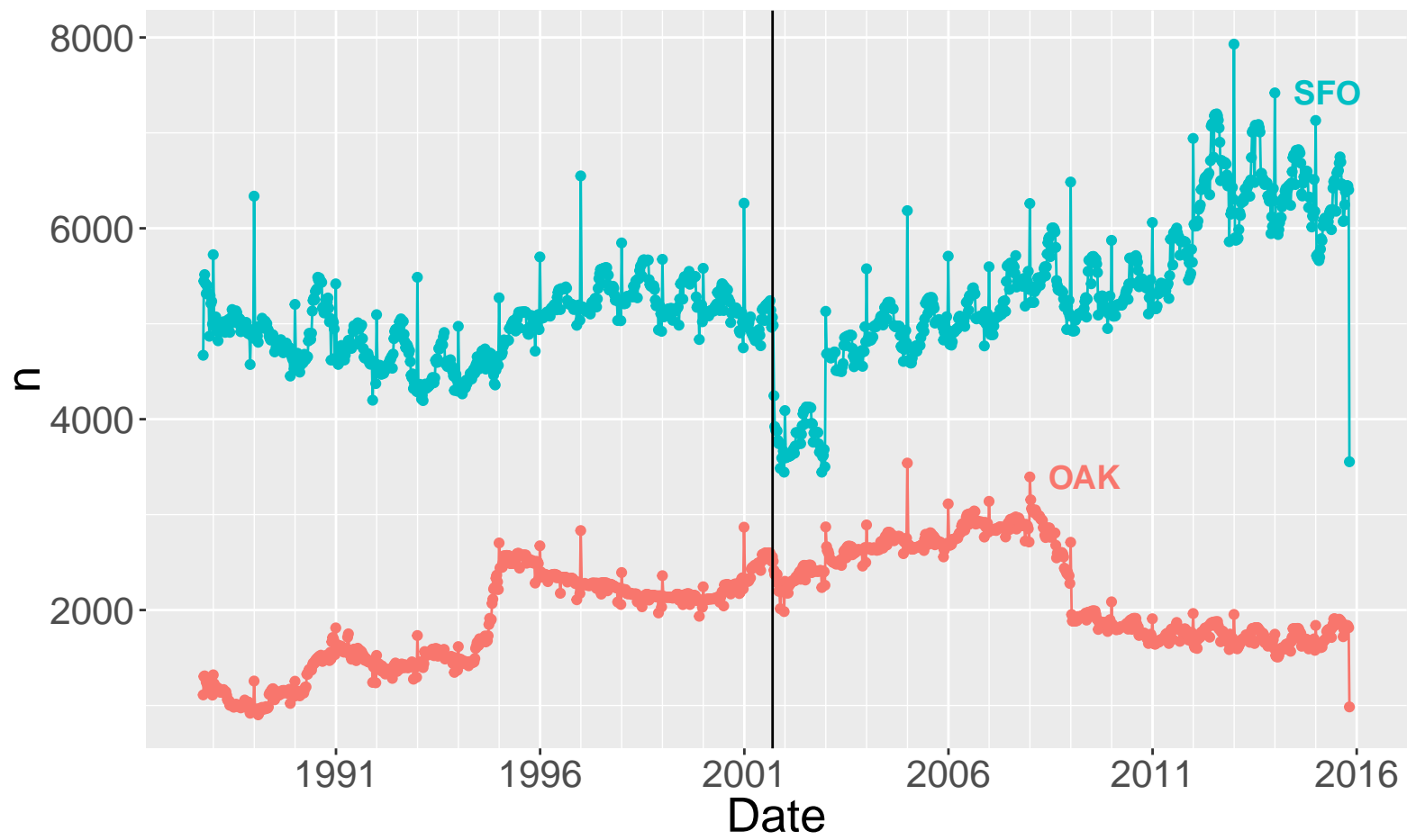


Figure 2.1: Number of depatures over time from Oakland and San Francisco Intl.


```

      nudge_x = 0.75,
      size = 5) + theme(legend.position = "none") +
  annotate("blank", x = 2019, y = 0) +
  facet_grid(Origin ~ .) +
  theme(text = element_text(size = 16))

## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: span too small. fewer data values than degrees of freedom.
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: pseudoinverse used at 2002.1
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: neighborhood radius 0.17125
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: reciprocal condition number 0
## Warning in simpleLoess(y, x, w, span, degree = degree, parametric = parametric,
: There are other near singularities as well. 0.029327
## Warning: Removed 4579 rows containing missing values (geom.text).

```

After September 11, flights from SFO fell, whereas OAK's volume did not. Flights fell more in SFO than they did in OAK because most of OAK's flights are from Southwest, which did not change its flight patterns. Furthermore, United was affected more than most airlines from the aftermath of the attacks.

```

top_5_carriers <-
  flights %>%
    count(UniqueCarrier) %>%
    arrange(desc(n)) %>%
    mutate(TopN = 1:n() <= 5) %>%
    mutate(Carrier_other = ifelse(TopN, UniqueCarrier, "Other")) %>%
    select(-n) %>%
    setkey(UniqueCarrier)

flights %>%
  setkey(UniqueCarrier) %>%
  merge(top_5_carriers) %>%
  count(Carrier_other, Year) %>%
  ggplot(aes(x = Year, y = n * sample.weight.int, color = Carrier_other, group = Carrier_other)) +
  geom_line() +
  scale_colour_brewer(palette = "Accent") +
  scale_y_continuous(label = scales::comma)

```

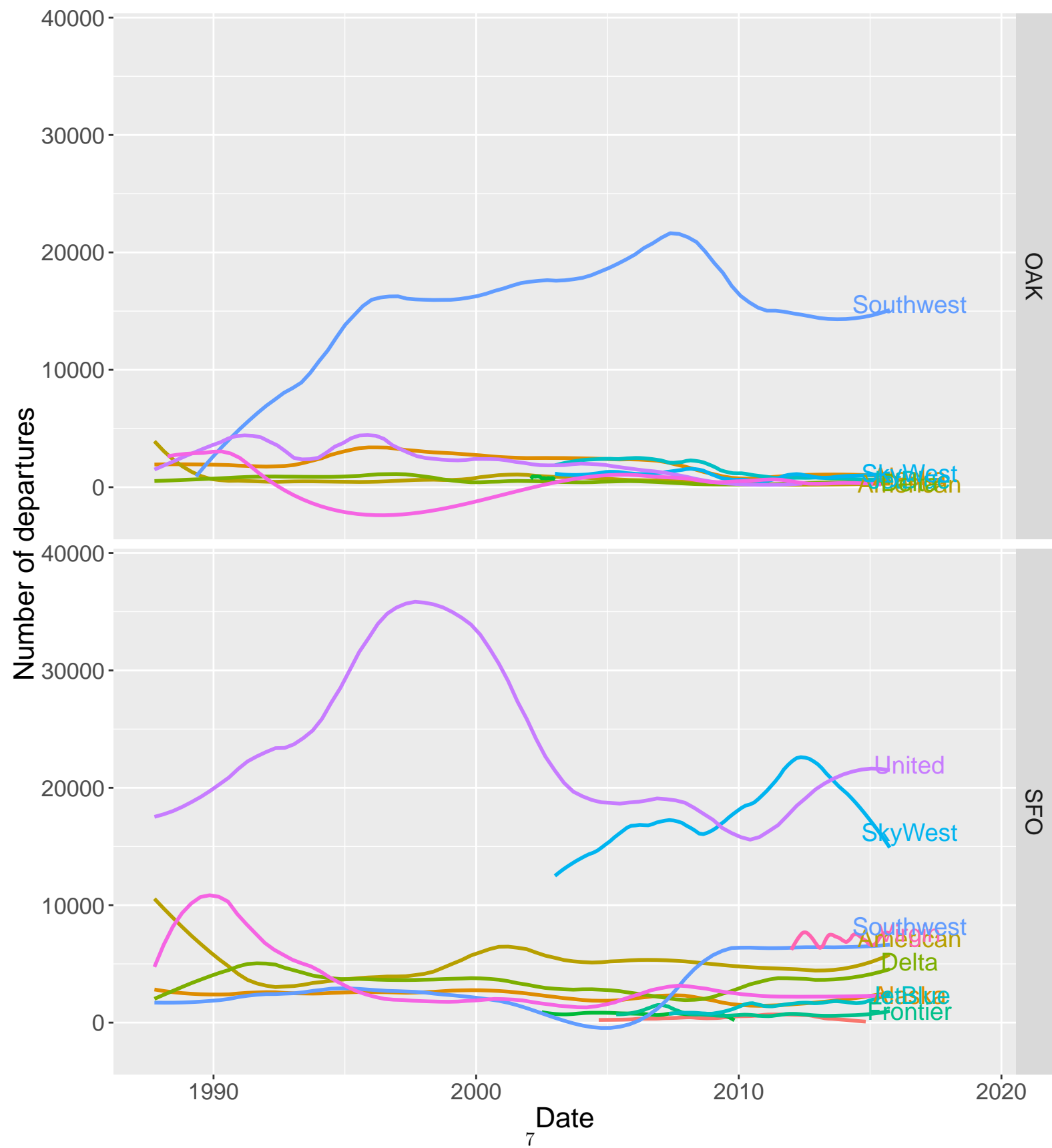
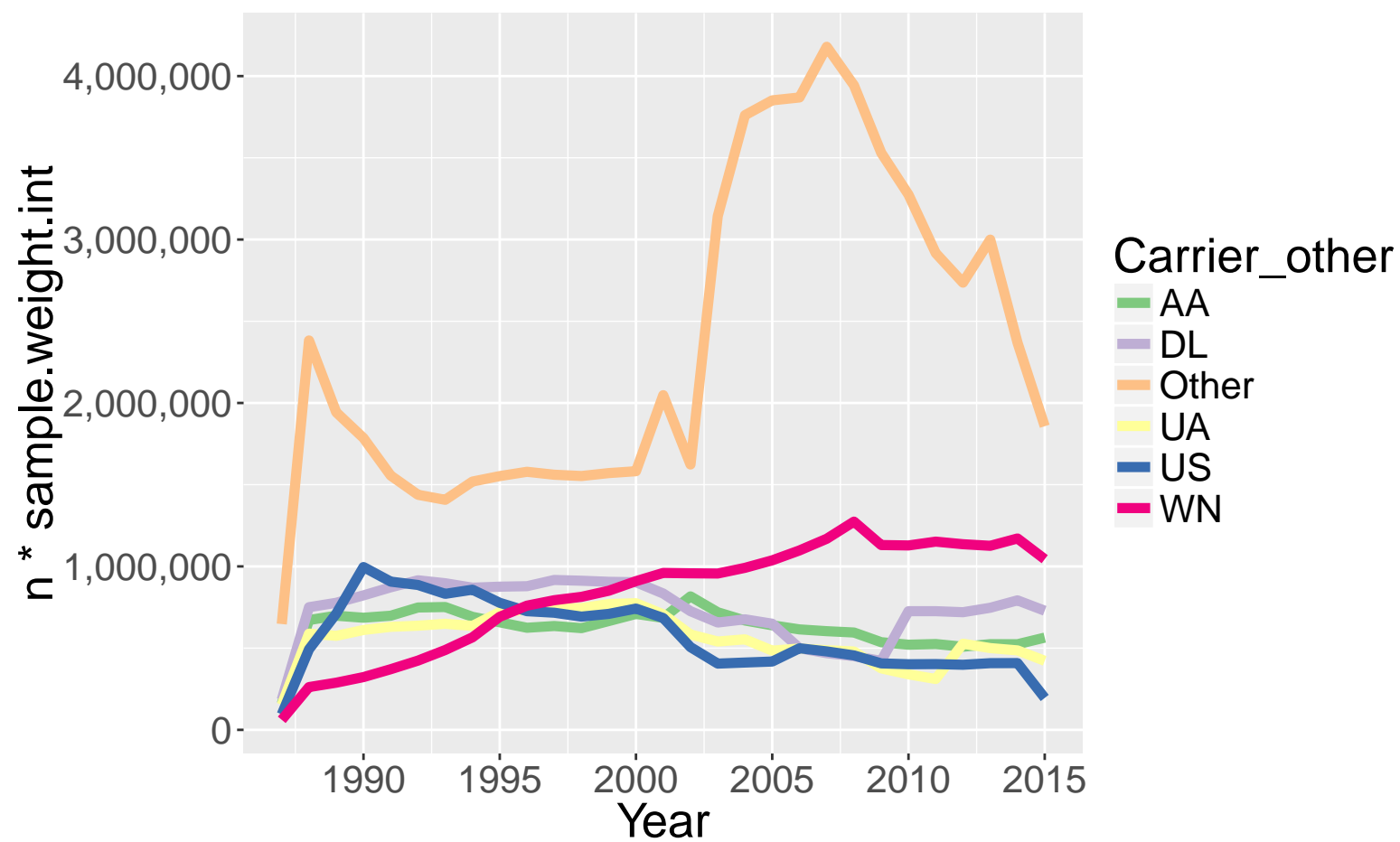


Figure 2.2: Number of depatures over time from Oakland and San Francisco Intl.



```

majorAirportThreshold = 10

major_airports <-
  flights[,.(n = .N), by = Dest][order(-n)] %>% # flights %>% count(Dest) %>% arrange(d
  mutate(TopN = 1:n() <= majorAirportThreshold) %>%
  mutate(AirportOther = ifelse(TopN, Dest, "Other_airport")) %>%
  select(-n) %>%
  setkey(Dest)

airports_by_volume_by_year <- flights[major_airports][,.(n = .N * sample.weight.int), by

airports_by_volume_by_2014 <-
  airports_by_volume_by_year %>%
  filter(Year == 2014) %>%
  filter(AirportOther != "AirportOther") %>%
  merge(select(nycflights.airports, faa, name), by.x = "AirportOther", by.y = "faa") %>%
  arrange(desc(n))
gc(0,1)

##           used      (Mb) gc trigger      (Mb) max used      (Mb)
## Ncells    707571    37.8   2423100    129.5   707571    37.8
## Vcells 988365629 7540.7 2678775838 20437.5 988365629 7540.7

setkey(flights, Dest)
gc(0,1)

##           used      (Mb) gc trigger      (Mb) max used      (Mb)
## Ncells    707518    37.8   2423100    129.5   707518    37.8
## Vcells 988362573 7540.7 2678775838 20437.5 988362573 7540.7

airports_by_volume_by_year %>%
  filter(AirportOther != "Other_airport", Year > 1987L, Year < 2015L) %>%
  merge(select(nycflights.airports, faa, name), by.x = "AirportOther", by.y = "faa") %>%
  mutate(name = factor(name, levels = airports_by_volume_by_2014$name)) %>%
  ggplot(aes(x = Year, y = n, group = name, color = name)) +
  geom_line()
gc(0,1)

##           used      (Mb) gc trigger      (Mb) max used      (Mb)
## Ncells    714769    38.2   2423100    129.5   714769    38.2
## Vcells 988383059 7540.8 2678775838 20437.5 988383059 7540.8

```

```

rel_vol_major_airports <-
  flights[major_airports][ ,.(n = .N * sample.weight.int), by = list(Year, AirportOther)]
  filter(AirportOther != "Other_airport", Year > 1987L, Year < 2015L) %>%
  arrange(Year) %>%
  group_by(AirportOther) %>%
  mutate(rel = n/first(n)) %>%
  merge(select(nycflights.airports, faa, name), by.x = "AirportOther", by.y = "faa")

last_values <-
  rel_vol_major_airports %>%
  filter(Year == max(Year)) %>%
  arrange(rel)

rel_vol_major_airports %>%
  mutate(name = factor(name, levels = rev(last_values$name))) %>%
  ggplot(aes(x = Year, y = rel, group = name, color = name)) +
  geom_line()

```

```

otp201510 <-
  fread("../dep_delay/On_Time_On_Time_Performance_2015_10.csv")

##
Read 43.2% of 486165 rows
Read 80.2% of 486165 rows
Read 486165 rows and 110 (of 110) columns from 0.204 GB file in 00:00:04

city_decoder <-
  otp201510 %>%
  select(contains("Origin")) %>%
  unique

setkey(city_decoder, OriginCityMarketID)

gc(T,T)

##           used   (Mb) gc trigger   (Mb)    max used   (Mb)
## Ncells    716669   38.3   2423100   129.5     716669   38.3
## Vcells 1033347922 7883.9 2678775838 20437.5 1033347922 7883.9

city_market_decoder <-
  fread("../metadata/L_CITY_MARKET_ID.csv") %>%
  setnames(old = c("Code", "Description"),
           new = c("OriginCityMarketID", "OriginCityMarketDescription")) %>%

```

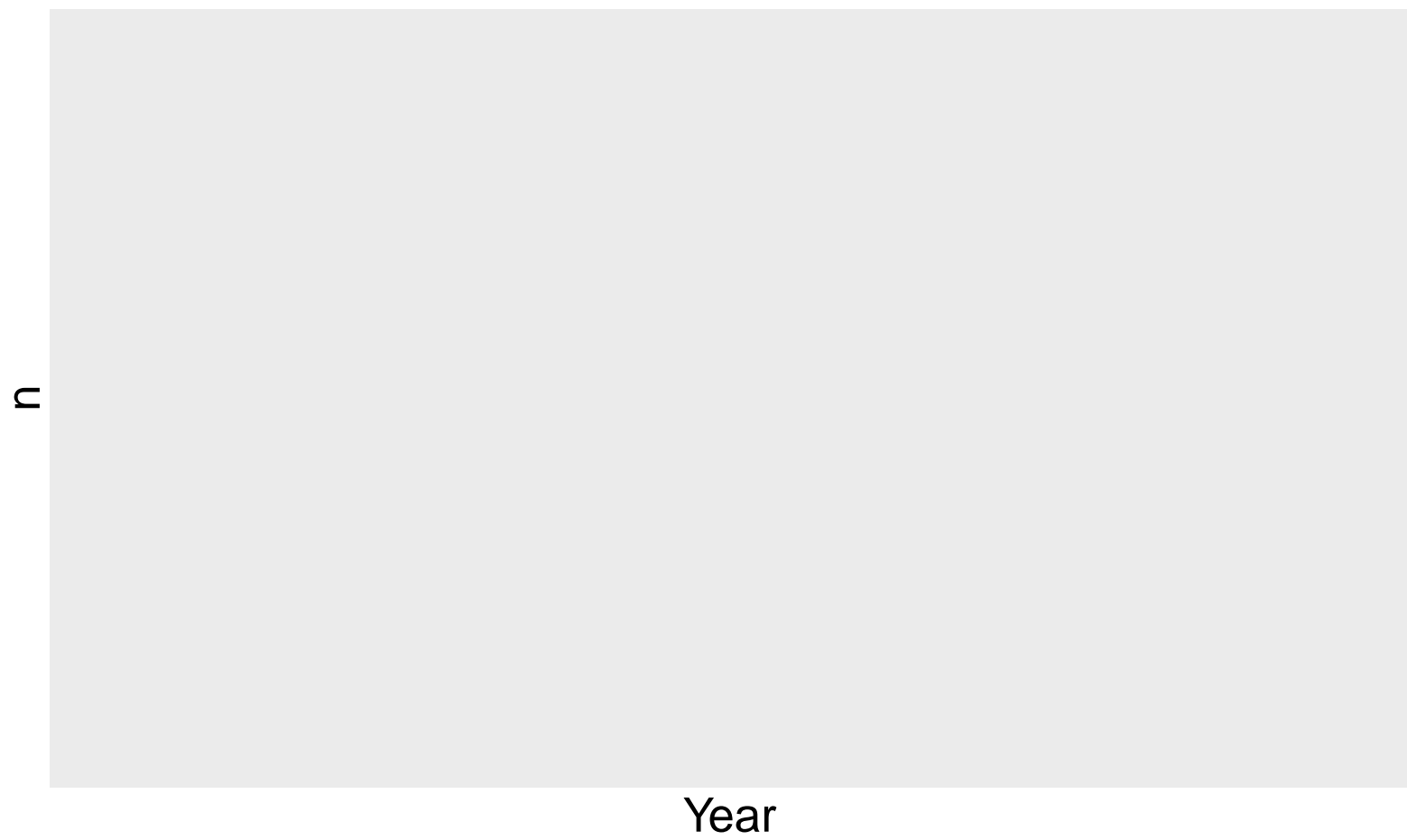


Figure 2.3: Annual flights by the top 10 airports by total volume.

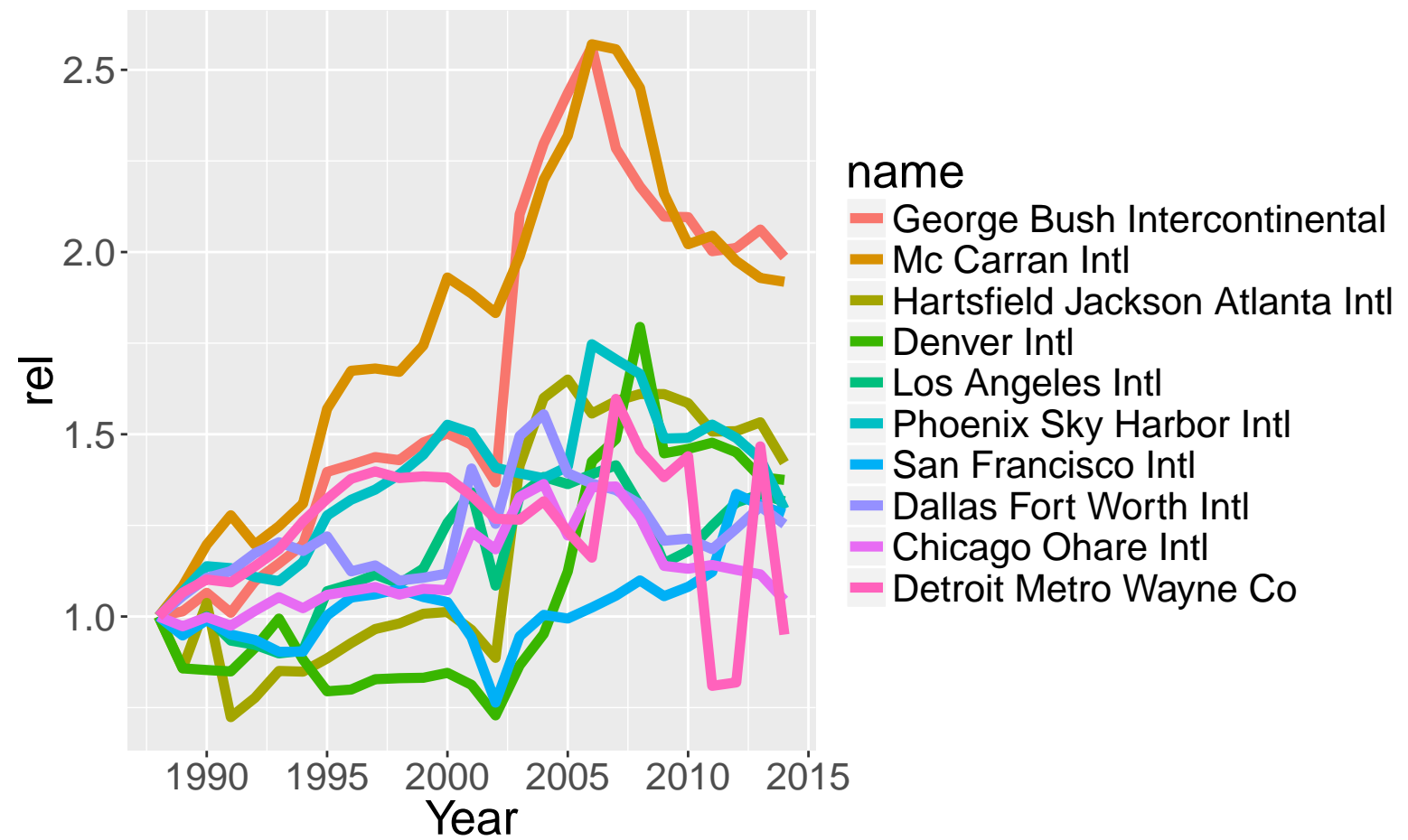


Figure 2.4: Annual flights by airport, 1988 = 1.

```

  setkey(OriginCityMarketID)
city_market_decoder[,OriginCityMarketID := as.integer(OriginCityMarketID)]
city_decoder <- merge(city_decoder, city_market_decoder, by = "OriginCityMarketID", all.=T)
gc(T,T)

```

```

##          used   (Mb) gc trigger   (Mb)    max used   (Mb)
## Ncells   722202  38.6   2423100  129.5     722202   38.6
## Vcells 1033374194 7884.1 2678775838 20437.5 1033374194 7884.1

```

```

market_volume_by_year <-
  flightsSanFran %>%
  filter(Dest %in% c("SFO", "OAK")) %>%
  merge(city_decoder, by = "Origin") %>%
  count(Year, OriginCityMarketDescription) %>%
  mutate(State = gsub("^.*([A-Z]{2}).*$", "\\1", OriginCityMarketDescription)) %>%
  filter(n > 3650) %>%
  mutate(Label = ifelse(Year == max(Year), OriginCityMarketDescription, NA_character_)) %>%
  arrange(Year, desc(n))

mkt.vol.by.yr <- function(year, colname){
  magrittr::extract2(dplyr::filter(market_volume_by_year, Year == year), colname)
}

market_volume_by_year %>%
  mutate(OriginCityMarketDescription = factor(OriginCityMarketDescription, levels = mkt.vol.by.yr$OriginCityMarketDescription)) %>%
  ggplot(aes(x = Year, y = n, color = OriginCityMarketDescription, group = OriginCityMarketDescription)) +
  #facet_grid(State ~ .) +
  geom_line() +
  #geom_text(aes(label = Label)) +
  #geom_dl(method = list("top.points", dl.trans(y = y+0.25), fontfamily = "bold"), aes(label = Label)) +
  theme(legend.position = "none") -> p
direct.label(p, list("top.points", dl.trans(y = y+0.25), fontface="bold"))

## Error in eval(expr, envir, enclos): could not find function "direct.label"

```

```

FAA_aircraft <-
  fread("../metadata/planes.csv") %>%
  setnames(old = c("tailnum", "year"), new = c("TailNum", "YearOfReg")) %>%
  setkey(TailNum)

```



```

flights %>%
  group_by(Origin, Dest) %>%
  filter(n() > 50000) %>%
  mutate(Route = paste0(Origin, "-", Dest),
         RevRoute = paste0(Dest, "-", Origin),
         maxRoute = pmax(Route, RevRoute)) %>%
  ggplot(aes(x = ActualElapsedTime)) +
  geom_density(aes(fill = maxRoute), alpha = 0.5) + xlim(0,300)

## Warning: Removed 32765 rows containing non-finite values (stat.density).

```

```

flights %>%
  select(Origin, Dest, ActualElapsedTime) %>%
  group_by(Origin, Dest) %>%
  summarise(average_time = mean(ActualElapsedTime, na.rm = TRUE),
            sd_time = sd(ActualElapsedTime, na.rm = TRUE),
            n = n()) %>%
  mutate(avg_less_sd = (sd_time - average_time) / average_time) %>%
  arrange(avg_less_sd) %>%
  mutate(Route = paste0(Origin, "-", Dest),
         Label = ifelse(Route %in% c('ROC-JFK', 'SLC-PHX', 'DCA-LGA', 'ORD-EWR'), Route,
                        hasLabel = !is.na(Label)) %>%
  ggplot(aes(x = average_time, y = sd_time)) +
  #geom_point(aes(alpha = n/max(n))) + scale_alpha_identity() +
  geom_point(aes(size = n, fill = hasLabel, alpha = hasLabel), color = "black", stroke = 0) +
  scale_fill_manual(values = c(Orange, "red")) +
  scale_alpha_manual(values = c(0.5, 1)) +
  geom_text(aes(label = Label), color = "red", fontface = "bold", hjust = 1.1, vjust = 0) +
  coord_cartesian(xlim = c(0,480), ylim = c(0,50)) +
  scale_x_continuous("Average elapsed time", expand = c(0,0)) +
  scale_y_continuous("SD of time", expand = c(0,0))

## Warning: Removed 486 rows containing missing values (geom.point).
## Warning: Removed 8275 rows containing missing values (geom.text).

```

2.1 Effect of 9-11

```

flights %>%
  group_by(Year, Month, DayofMonth) %>%
  summarise(prop_cancelled = mean(Cancelled)) %>%

```

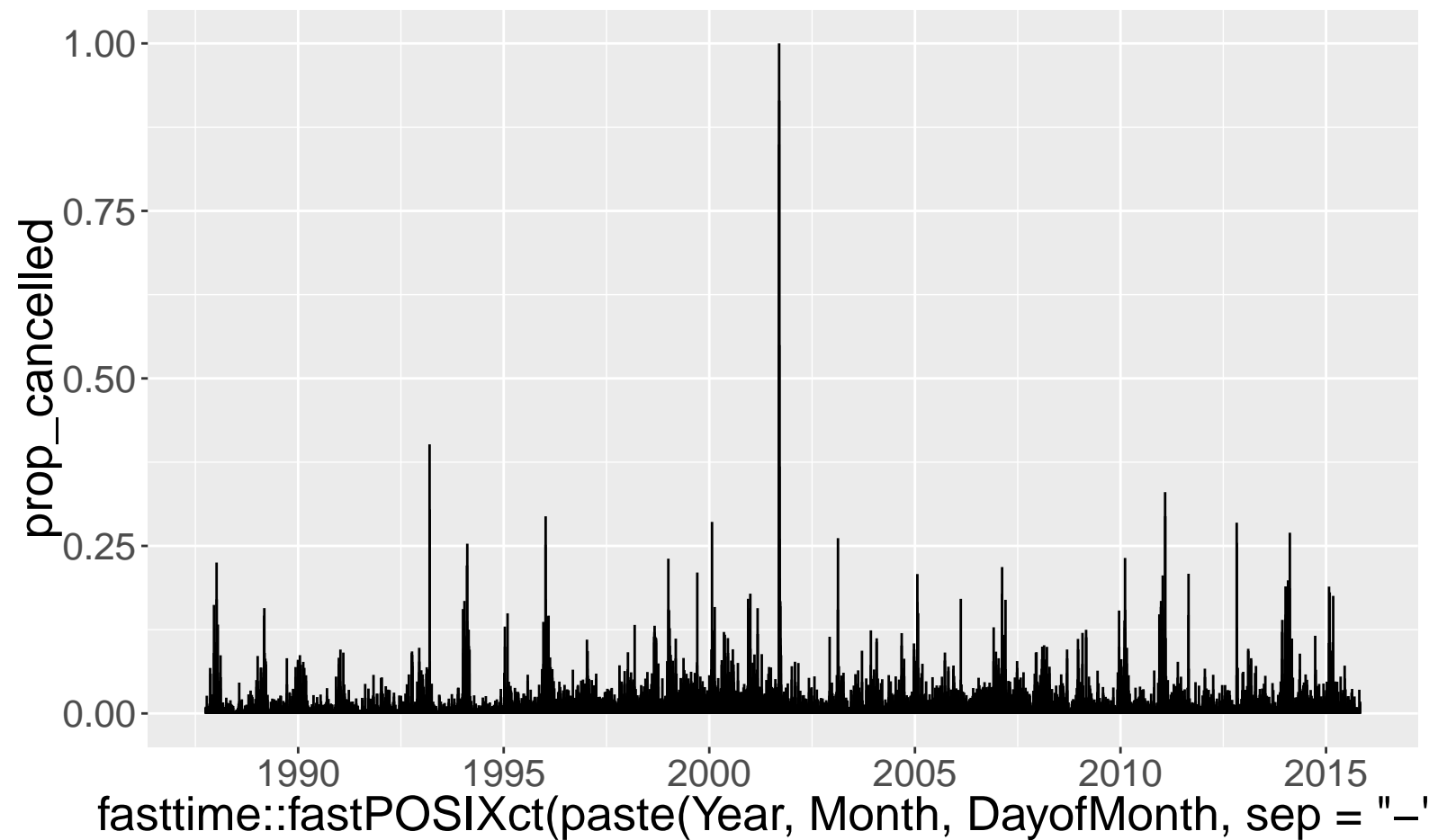


Figure 2.5: Proportion of flights cancelled by date.

```
ggplot(aes(x = fasttime::fastPOSIXct(paste(Year, Month, DayofMonth, sep = "-")), y = prop_cancelled)) +  
  geom_bar(stat = "identity", width=1)
```

```
flights %>%  
  group_by(Year, Month, DayofMonth) %>%  
  summarise(prop_cancelled = mean(Cancelled)) %>%  
  ungroup %>%  
  mutate(rank = dense_rank(prop_cancelled)) %>%  
  ggplot(aes(x = jitter(rank, amount = 0.1), y = prop_cancelled)) + geom_bar(stat = "identity")  
  
## Warning: position_stack requires non-overlapping x intervals
```

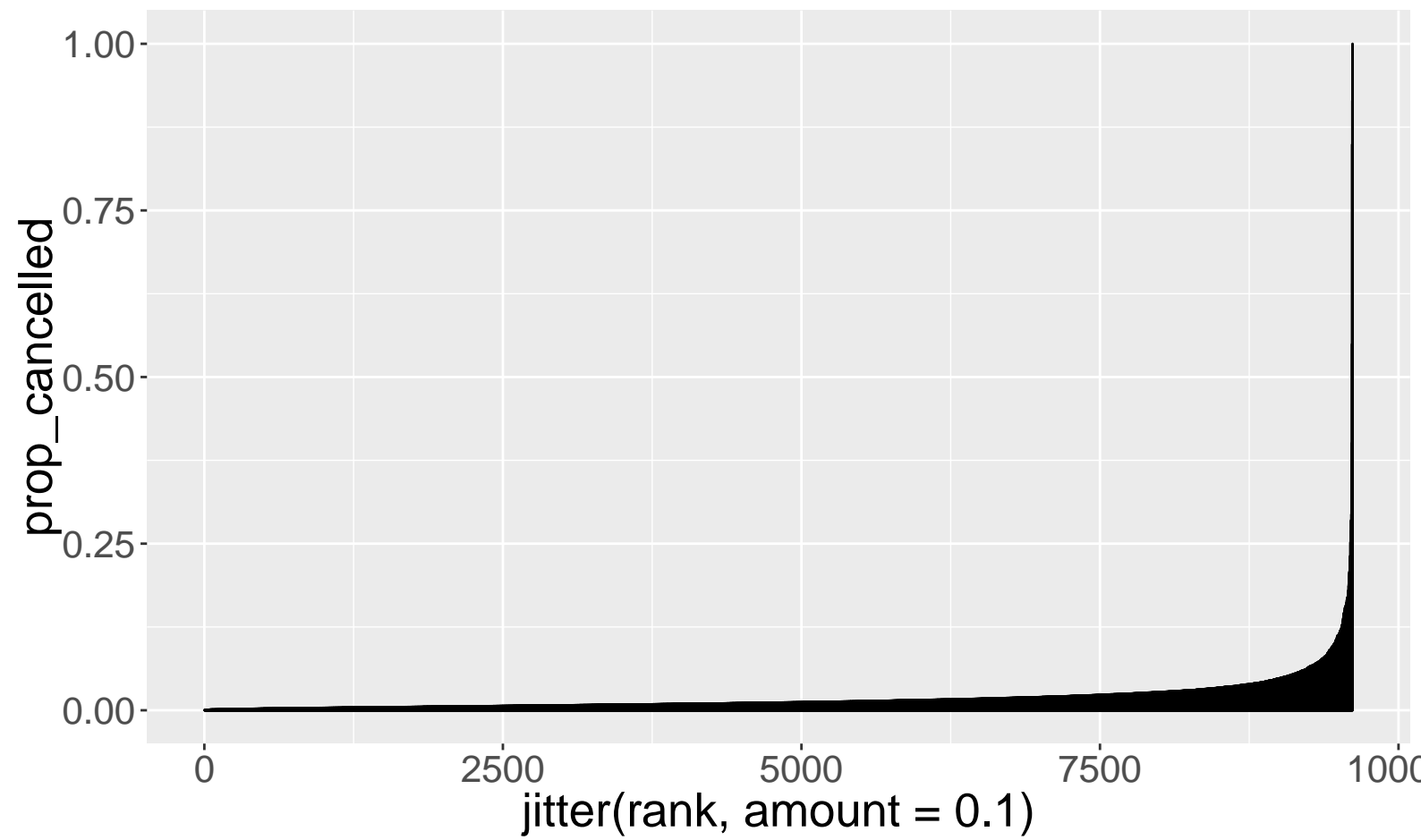
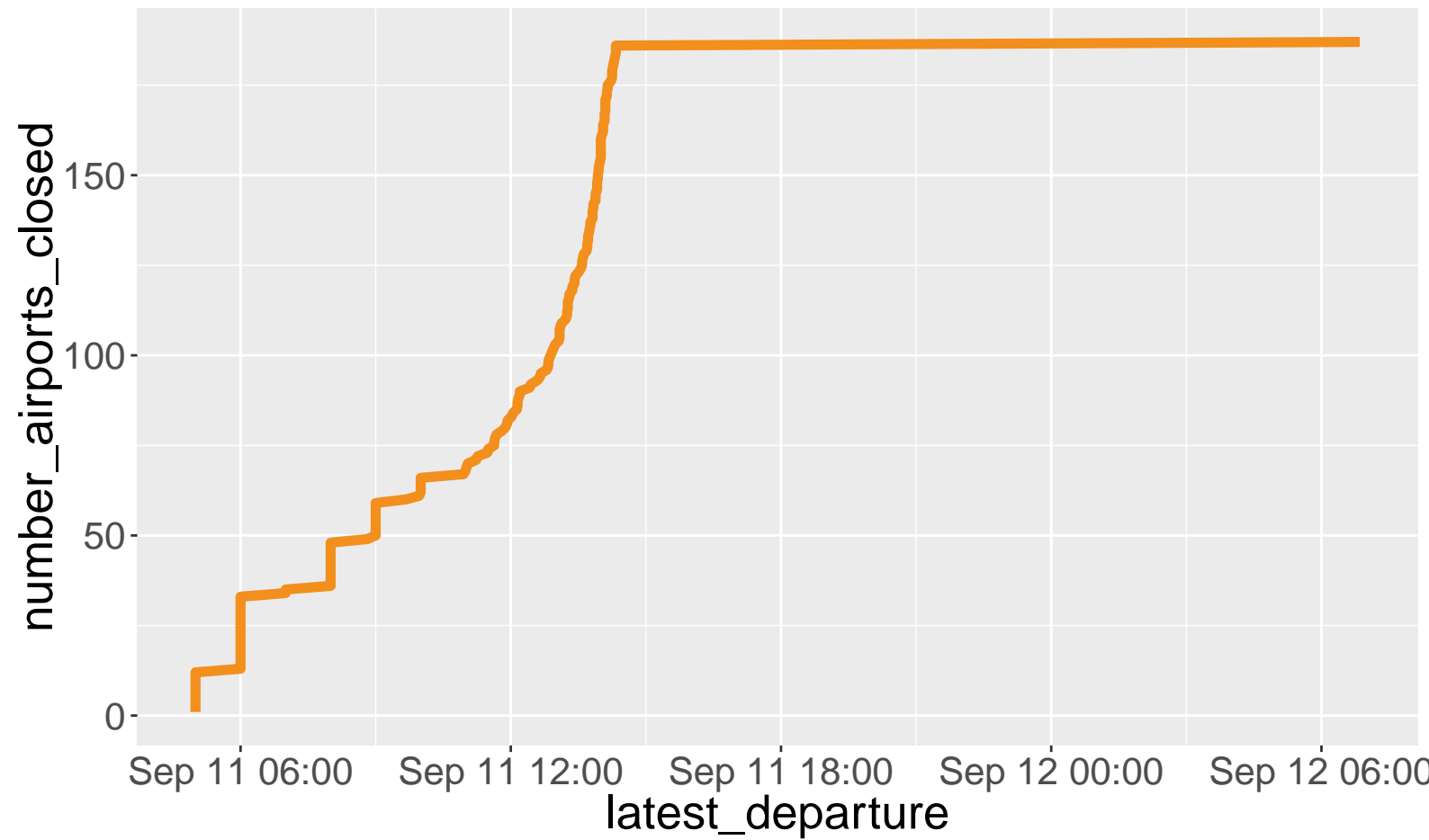


Figure 2.6: Pareto distribution of cancellation proportions.

Figure 2.7: Number of airports closed by UTC (determined by date of last departure)



```
flights %>%  
  filter(Year == 2001, Month == 9, DayofMonth == 11) %>%  
  group_by(Origin) %>%  
  summarise(latest_departure = max(DepDateTimeZulu)) %>%  
  ungroup %>%  
  arrange(latest_departure) %>%  
  mutate(number_airports_closed = 1:n()) %>%  
  ggplot(aes(x = latest_departure, y = number_airports_closed)) +  
  geom_line(group = 1) +  
  geom_vline(xintercept = as.numeric(as.POSIXct("2001-09-11 09:17:00")))
```

```
FINISH.TIME <- Sys.time()
```

Compiled in 12.4807684818904