

# MY-DATA607-Week05-Flights

Michael Y.

September 29, 2019

## Contents

<b>Week 5 - Flight Delays</b>	<b>3</b>
The assignment is as follows: . . . . .	3
Load libraries . . . . .	4
Data Loading . . . . .	6
Load the raw datafile (which I created by entering the given values into an excel spreadsheet, then saving the spreadsheet as a .csv file) :	6
Data Cleanup Work . . . . .	7
Reviewing Hadley Wickham's directive for creating "tidy" data, it is necessary to distinguish between <i>Fixed Variables</i> and <i>Measured Variables</i> . . . . .	7
<b>Fixed Variables</b> : . . . . .	7
<b>Measured Variables</b> : . . . . .	7
<b>tidyr / dplyr</b> steps: . . . . .	8
<b>GATHER</b> : . . . . .	8
<b>MUTATE(lag)</b> : . . . . .	9
<b>SPREAD</b> : . . . . .	10
<b>RENAME</b> : . . . . .	10
<b>MUTATE</b> : . . . . .	11
Chain using pipe connector "%>%" : . . . . .	12
Analyze the data . . . . .	13
Let's sort the above data by city, then by airline: . . . . .	13
Plot percent delays <b>by city</b> : . . . . .	15

Flights and delays <b>by airline</b> . . . . .	16
Simpson's Paradox: . . . . .	18
City by city: . . . . .	19
Here's a barplot: . . . . .	21
Relative market share . . . . .	24
Plot airline market share, by City: . . . . .	25
Good Weather vs. Bad Weather . . . . .	26
Conclusion . . . . .	26
Simpson's Paradox explained: . . . . .	26

---

## Week 5 - Flight Delays

The assignment is as follows:

A chart is supplied, listing flight performance (on-time vs. delayed) for two airlines (ALASKA and AM WEST) across five cities (Los Angeles, Phoenix, San Diego, San Francisco, and Seattle):

```
# pull Airline Delays JPG file from github - only usable for HTML knit -- will not work for PDF knit
airline_delays_pic_URL <- "https://raw.githubusercontent.com/myampol/MY607/master/Airline_Delays.JPG"
### Determine whether we are knitting to PDF ("latex") or HTML
whichknit <- knitr::opts_knit$get("rmarkdown.pandoc.to") # Works only if knit() is called via render()
if (!is.null(whichknit)) { # can't test the below items if whichknit returned NULL
  if (whichknit=="latex") {
    knitr::include_graphics('Airline_Delays.JPG') ### PDF knit requires that the file be local
  } else if (whichknit=="html") {
    knitr::include_graphics(airline_delays_pic_URL) ### HTML knit will accept URL to insert picture
  }
}
```

		Los Angeles	Phoenix	San Diego	San Francisco	Seattle
ALASKA	on time	497	221	212	503	1,841
	delayed	62	12	20	102	305
AM WEST	on time	694	4,840	383	320	201
	delayed	117	415	65	129	61

Source: *Numbersense*, Kaiser Fung, McGraw Hill, 2013

- (1) Create a .CSV file (or optionally, a MySQL database!) that includes all of the information above. You're encouraged to use a "wide" structure similar to how the information appears above, so that you can practice tidying and transformations as described below.
- (2) Read the information from your .CSV file into R, and use `tidyr` and `dplyr` as needed to tidy and transform your data.
- (3) Perform analysis to compare the arrival delays for the two airlines.
- (4) Your code should be in an R Markdown file, posted to [rpubs.com](https://rpubs.com), and should include narrative descriptions of your data cleanup work, analysis, and conclusions.

## Load libraries

```
#library(readr)  
#library(stringr)  
library(tidyr)  
library(dplyr)
```

```
##  
## Attaching package: 'dplyr'  
  
## The following objects are masked from 'package:stats':  
##  
##     filter, lag  
  
## The following objects are masked from 'package:base':  
##  
##     intersect, setdiff, setequal, union
```

```
library(ggplot2)  
library(kableExtra)
```

```
##  
## Attaching package: 'kableExtra'
```

```
## The following object is masked from 'package:dplyr':  
##  
##      group_rows
```

## Data Loading

Load the raw datafile (which I created by entering the given values into an excel spreadsheet, then saving the spreadsheet as a .csv file) :

```
#setwd("C:/Users/Michael/Dropbox/priv/CUNY/MSDS/201909-Fall/DATA607_Tati_Andy/20190929_Week05")
### Don't pull file from local drive
#inputfile <- "InputFlightData.csv"
### Pull the input data file from github, rather than local drive
inputfile <- "https://raw.githubusercontent.com/myampol/MY607/master/InputFlightData.csv"
rawflights <- read.csv(inputfile,stringsAsFactors = F)
rawflights %>% kable() %>% kable_styling(c("striped", "bordered"))
```

X	X.1	Los.Angeles	Phoenix	San.Diego	San.Francisco	Seattle
ALASKA	on time	497	221	212	503	1841
	delayed	62	12	20	102	305
		NA	NA	NA	NA	NA
AM WEST	on time	694	4840	383	320	201
	delayed	117	415	65	129	61

Modify the column headings to clean up the naming:

```
rf1 <- rename(.data = rawflights, Airline=X, Status=`X.1`,
              LosAngeles=`Los.Angeles`, SanDiego=`San.Diego`, SanFrancisco=`San.Francisco`)
rf1 %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	Status	LosAngeles	Phoenix	SanDiego	SanFrancisco	Seattle
ALASKA	on time	497	221	212	503	1841
	delayed	62	12	20	102	305
		NA	NA	NA	NA	NA
AM WEST	on time	694	4840	383	320	201
	delayed	117	415	65	129	61

## Data Cleanup Work

Reviewing Hadley Wickham's directive for creating "tidy" data, it is necessary to distinguish between *Fixed Variables* and *Measured Variables* .

### Fixed Variables :

- (a) the airline (AM West or Alaska), and
- (b) the city (Los Angeles, Phoenix, San Diego, San Francisco, and Seattle)

### Measured Variables :

- (c) the count of **ON TIME** flights and
- (d) the count of **DELAYED** flights

for each (Airline, City) pair.

tidyr / dplyr steps:

**GATHER :**

First, use gather to put all the Cities into one column, and also drop the row containing NAs:

```
rf2 <- gather(data = rf1 , key = City, value = NumFlights, ... = LosAngeles:Seattle, na.rm = T)
rf2 %>% kable() %>% kable_styling(c("striped", "bordered"))
```

	Airline	Status	City	NumFlights
1	ALASKA	on time	LosAngeles	497
2		delayed	LosAngeles	62
4	AM WEST	on time	LosAngeles	694
5		delayed	LosAngeles	117
6	ALASKA	on time	Phoenix	221
7		delayed	Phoenix	12
9	AM WEST	on time	Phoenix	4840
10		delayed	Phoenix	415
11	ALASKA	on time	SanDiego	212
12		delayed	SanDiego	20
14	AM WEST	on time	SanDiego	383
15		delayed	SanDiego	65
16	ALASKA	on time	SanFrancisco	503
17		delayed	SanFrancisco	102
19	AM WEST	on time	SanFrancisco	320
20		delayed	SanFrancisco	129
21	ALASKA	on time	Seattle	1841
22		delayed	Seattle	305
24	AM WEST	on time	Seattle	201
25		delayed	Seattle	61



MUTATE(lag) :

Next, use mutate(lag) to propagate the Airline names downward (from each odd-numbered row) to fill the missing airline name (on the subsequent even-numbered row):

```
rf3 <- mutate(.data = rf2, Airline= ifelse(Airline=="", lag(Airline), Airline))
rf3 %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	Status	City	NumFlights
ALASKA	on time	LosAngeles	497
ALASKA	delayed	LosAngeles	62
AM WEST	on time	LosAngeles	694
AM WEST	delayed	LosAngeles	117
ALASKA	on time	Phoenix	221
ALASKA	delayed	Phoenix	12
AM WEST	on time	Phoenix	4840
AM WEST	delayed	Phoenix	415
ALASKA	on time	SanDiego	212
ALASKA	delayed	SanDiego	20
AM WEST	on time	SanDiego	383
AM WEST	delayed	SanDiego	65
ALASKA	on time	SanFrancisco	503
ALASKA	delayed	SanFrancisco	102
AM WEST	on time	SanFrancisco	320
AM WEST	delayed	SanFrancisco	129
ALASKA	on time	Seattle	1841
ALASKA	delayed	Seattle	305
AM WEST	on time	Seattle	201
AM WEST	delayed	Seattle	61

## SPREAD :

Now, use spread to put the “on time” and “delayed” counts into separate columns:

```
rf4 <- spread(data = rf3, key = Status, value = NumFlights)
rf4 %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	City	delayed	on time
ALASKA	LosAngeles	62	497
ALASKA	Phoenix	12	221
ALASKA	SanDiego	20	212
ALASKA	SanFrancisco	102	503
ALASKA	Seattle	305	1841
AM WEST	LosAngeles	117	694
AM WEST	Phoenix	415	4840
AM WEST	SanDiego	65	383
AM WEST	SanFrancisco	129	320
AM WEST	Seattle	61	201

## RENAME :

Use rename to improve the names of the “on time” and “delayed” columns:

```
rf5 <- rename(.data = rf4, NumFlightsDelayed=delayed, NumFlightsOnTime=`on time`)
rf5 %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	City	NumFlightsDelayed	NumFlightsOnTime
ALASKA	LosAngeles	62	497
ALASKA	Phoenix	12	221
ALASKA	SanDiego	20	212
ALASKA	SanFrancisco	102	503
ALASKA	Seattle	305	1841
AM WEST	LosAngeles	117	694
AM WEST	Phoenix	415	4840
AM WEST	SanDiego	65	383
AM WEST	SanFrancisco	129	320
AM WEST	Seattle	61	201

## MUTATE :

Use mutate to compute and append NumFlightsTotal:

```
rf6 <- mutate(.data = rf5, NumFlightsTotal = NumFlightsDelayed + NumFlightsOnTime)
rf6 %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	City	NumFlightsDelayed	NumFlightsOnTime	NumFlightsTotal
ALASKA	LosAngeles	62	497	559
ALASKA	Phoenix	12	221	233
ALASKA	SanDiego	20	212	232
ALASKA	SanFrancisco	102	503	605
ALASKA	Seattle	305	1841	2146
AM WEST	LosAngeles	117	694	811
AM WEST	Phoenix	415	4840	5255
AM WEST	SanDiego	65	383	448
AM WEST	SanFrancisco	129	320	449
AM WEST	Seattle	61	201	262

Use mutate to compute and append the *percentage* of delayed and ontime flights at each city:

```
rf7 <- mutate(.data = rf6, PctFlightsDelayed=NumFlightsDelayed/NumFlightsTotal,
                PctFlightsOnTime=NumFlightsOnTime/NumFlightsTotal)
rf7 %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	City	NumFlightsDelayed	NumFlightsOnTime	NumFlightsTotal	PctFlightsDelayed	PctFlightsOnTime
ALASKA	LosAngeles	62	497	559	0.110912	0.889088
ALASKA	Phoenix	12	221	233	0.051502	0.948498
ALASKA	SanDiego	20	212	232	0.086207	0.913793
ALASKA	SanFrancisco	102	503	605	0.168595	0.831405
ALASKA	Seattle	305	1841	2146	0.142125	0.857875
AM WEST	LosAngeles	117	694	811	0.144266	0.855734
AM WEST	Phoenix	415	4840	5255	0.078972	0.921028
AM WEST	SanDiego	65	383	448	0.145089	0.854911
AM WEST	SanFrancisco	129	320	449	0.287305	0.712695
AM WEST	Seattle	61	201	262	0.232824	0.767176

The above shows the data tidying and manipulation step-by-step.

Chain using pipe connector “%>%” :

Using the pipe connector “%>%”, all the above steps can be specified in a single chain:

```
tidy_flights <- rawflights %>%
  rename(.data = ., Airline=X, Status=`X.1`, LosAngeles=`Los.Angeles`,
        SanDiego=`San.Diego`,
        SanFrancisco=`San.Francisco`) %>%
  gather( data = ., key = City, value = NumFlights, ... = LosAngeles:Seattle, na.rm = T) %>%
  mutate(.data = ., Airline= ifelse(Airline=="", lag(Airline), Airline)) %>%
  spread( data = ., key = Status, value = NumFlights) %>%
  rename(.data = ., NumFlightsDelayed=delayed, NumFlightsOnTime=`on time`) %>%
  mutate(.data = ., NumFlightsTotal = NumFlightsDelayed + NumFlightsOnTime) %>%
  mutate(.data = ., PctFlightsDelayed=NumFlightsDelayed/NumFlightsTotal,
        PctFlightsOnTime=NumFlightsOnTime/NumFlightsTotal)
tidy_flights %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	City	NumFlightsDelayed	NumFlightsOnTime	NumFlightsTotal	PctFlightsDelayed	PctFlightsOnTime
ALASKA	LosAngeles	62	497	559	0.110912	0.889088
ALASKA	Phoenix	12	221	233	0.051502	0.948498
ALASKA	SanDiego	20	212	232	0.086207	0.913793
ALASKA	SanFrancisco	102	503	605	0.168595	0.831405
ALASKA	Seattle	305	1841	2146	0.142125	0.857875
AM WEST	LosAngeles	117	694	811	0.144266	0.855734
AM WEST	Phoenix	415	4840	5255	0.078972	0.921028
AM WEST	SanDiego	65	383	448	0.145089	0.854911
AM WEST	SanFrancisco	129	320	449	0.287305	0.712695
AM WEST	Seattle	61	201	262	0.232824	0.767176

The above result matches that from the step-by-step process.

## Analyze the data

Let's sort the above data by city, then by airline:

```
arrange(tidy_flights, City, Airline) %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	City	NumFlightsDelayed	NumFlightsOnTime	NumFlightsTotal	PctFlightsDelayed	PctFlightsOnTime
ALASKA	LosAngeles	62	497	559	0.110912	0.889088
AM WEST	LosAngeles	117	694	811	0.144266	0.855734
ALASKA	Phoenix	12	221	233	0.051502	0.948498
AM WEST	Phoenix	415	4840	5255	0.078972	0.921028
ALASKA	SanDiego	20	212	232	0.086207	0.913793
AM WEST	SanDiego	65	383	448	0.145089	0.854911
ALASKA	SanFrancisco	102	503	605	0.168595	0.831405
AM WEST	SanFrancisco	129	320	449	0.287305	0.712695
ALASKA	Seattle	305	1841	2146	0.142125	0.857875
AM WEST	Seattle	61	201	262	0.232824	0.767176

```
ALASKA_Phoenix_delays <- filter(.data=tidy_flights, Airline=="ALASKA" & City=="Phoenix") %>%  
  select(PctFlightsDelayed)  
AMWEST_Phoenix_delays <- filter(.data=tidy_flights, Airline=="AM WEST" & City=="Phoenix") %>%  
  select(PctFlightsDelayed)  
ALASKA_SanFrancisco_delays <- filter(.data=tidy_flights, Airline=="ALASKA" & City=="SanFrancisco") %>%  
  select(PctFlightsDelayed)  
AMWEST_SanFrancisco_delays <- filter(.data=tidy_flights, Airline=="AM WEST" & City=="SanFrancisco") %>%  
  select(PctFlightsDelayed)
```

What's noticeable is that looking *city-by-city*, the Percentage of Flights Delayed is smaller for ALASKA than it is for AMWEST:

Percentage of flights delayed for each airline, by city:

```
Pct_Delays_by_City <- tidy_flights %>%  
  select(.data = ., -NumFlightsDelayed, -NumFlightsOnTime,  
          -NumFlightsTotal, -PctFlightsOnTime) %>%  
  spread(data = ., key = Airline, value = PctFlightsDelayed)  
Pct_Delays_by_City %>% kable() %>% kable_styling(c("striped", "bordered"))
```

City	ALASKA	AM WEST
LosAngeles	0.110912	0.144266
Phoenix	0.051502	0.078972
SanDiego	0.086207	0.145089
SanFrancisco	0.168595	0.287305
Seattle	0.142125	0.232824

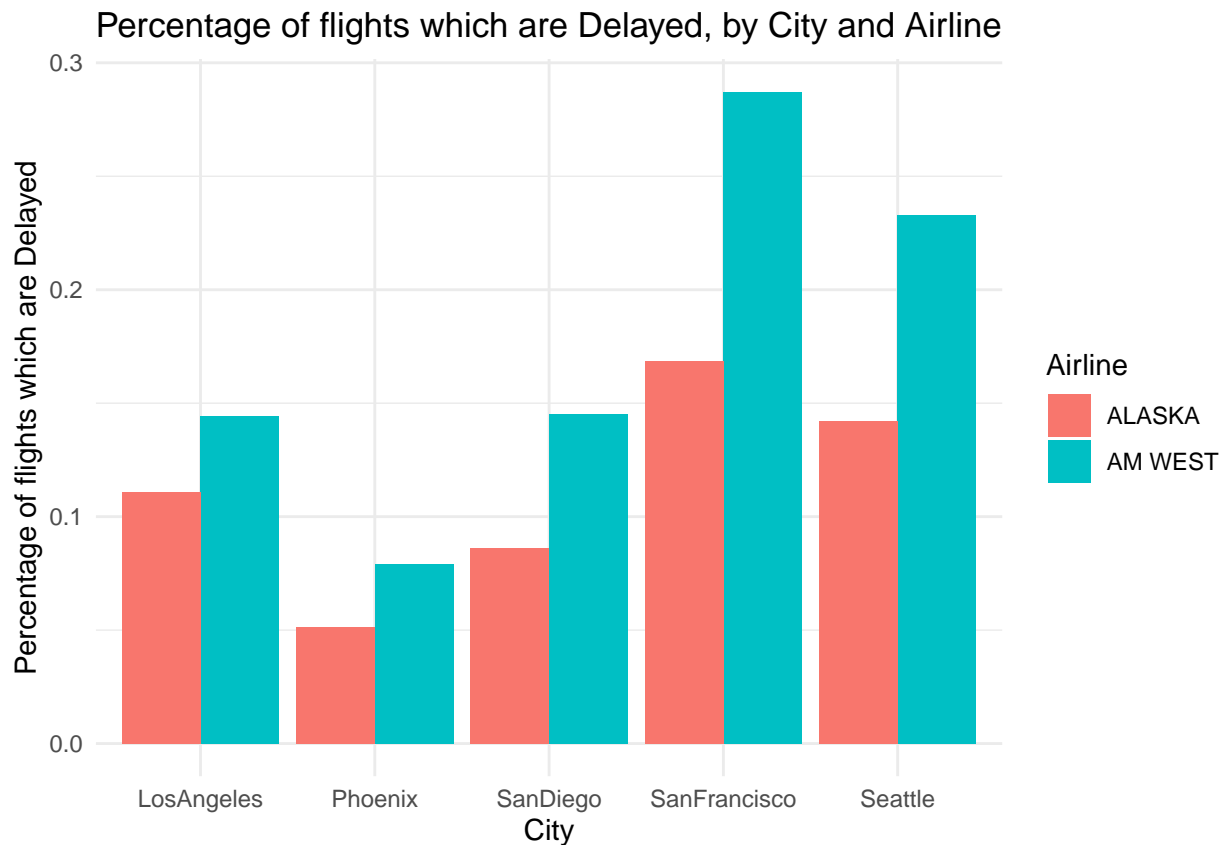
For example, at Phoenix, ALASKA's delays are 5.2% while AM WEST's delays are 7.9% .

At SanFrancisco, ALASKA's delays are 16.9% while AM WEST's delays are 28.7% .

The above relationship holds *for each city*.

Plot percent delays by city:

```
Pct_Delays_by_City %>% gather(data = ., key = Airline,  
                               value = Pct_Delays_by_City,...=ALASKA:`AM WEST`) %>%  
  ggplot(data = ., aes(factor(City), Pct_Delays_by_City, fill = Airline)) +  
  geom_bar(stat="identity", position = "dodge") +  
  theme_minimal()+  
  labs( x="City", y="Percentage of flights which are Delayed") +  
  ggtitle("Percentage of flights which are Delayed, by City and Airline")
```



## Flights and delays by airline

Question: How many TOTAL flights does each airline have, and what PERCENT are delayed?

```
ResultsByAirline <- group_by(tidy_flights,Airline) %>% summarize(  
  TotalDelays=sum(NumFlightsDelayed),  
  TotalOnTime=sum(NumFlightsOnTime),  
  TotalFlights=sum(NumFlightsTotal),  
  PctDelayed=TotalDelays/TotalFlights,  
  PctOnTime=TotalOnTime/TotalFlights)  
ResultsByAirline %>% kable() %>% kable_styling(c("striped", "bordered"))
```

Airline	TotalDelays	TotalOnTime	TotalFlights	PctDelayed	PctOnTime
ALASKA	501	3274	3775	0.132715	0.867285
AM WEST	787	6438	7225	0.108927	0.891073

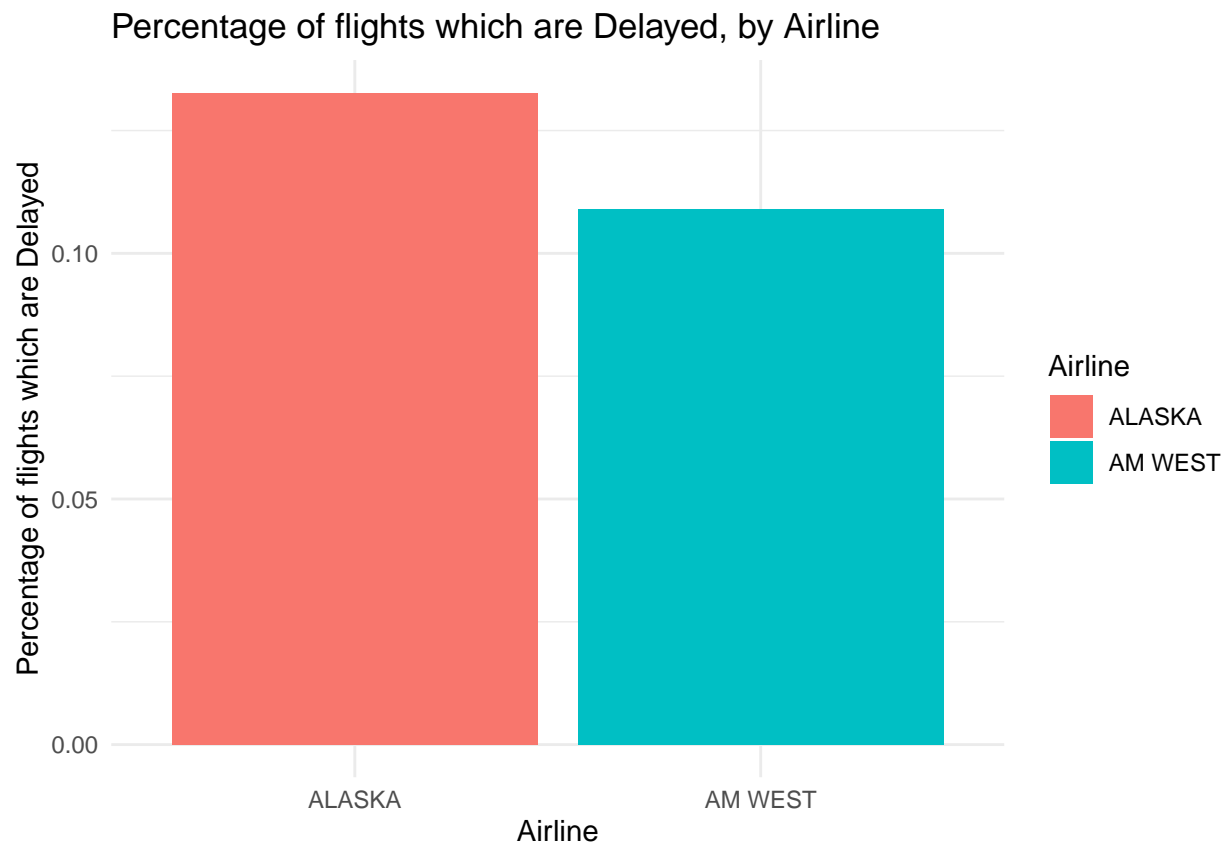
```
ALASKAdelays <- filter(.data = ResultsByAirline, Airline=="ALASKA") %>% select(PctDelayed)  
AMWESTdelays <- filter(.data = ResultsByAirline, Airline=="AM WEST") %>% select(PctDelayed)  
ALASKAtotals <- filter(.data = ResultsByAirline, Airline=="ALASKA") %>% select(TotalFlights)  
AMWESTtotals <- filter(.data = ResultsByAirline, Airline=="AM WEST") %>% select(TotalFlights)
```

These results show that while ALASKA runs about *half as many flights* (3775) as its competitor AM WEST (7225), a larger percentage (13.3 %) of ALASKA's flights are delayed vs. AM WEST, which suffered delays on only 10.9 % of its flights.



Plot:

```
ResultsByAirline %>% select(.data = ., Airline, PctDelayed) %>%  
  ggplot(data = ., aes(factor(Airline), PctDelayed, fill = Airline)) +  
  geom_bar(stat="identity", position = "dodge") +  
  theme_minimal() +  
  labs(x="Airline", y="Percentage of flights which are Delayed") +  
  ggtitle("Percentage of flights which are Delayed, by Airline")
```



Simpson's Paradox:

So, this is curious!

On a city-by-city basis, ALASKA "beat" AM WEST by having better on-time performance at every city.

But on an overall basis, AM WEST had the best overall on-time results!

This is a manifestation of a curiosity known as "Simpson's Paradox".

[https://en.wikipedia.org/wiki/Simpson%27s\\_paradox](https://en.wikipedia.org/wiki/Simpson%27s_paradox)

How could this paradox be explained?

Although the data shows that each airline serves the same 5 cities, they seem to focus on different markets.

Perhaps looking more closely at the different cities served by each airline can help explain?

City by city:

First, let's determine how many TOTAL flights go to each city, and what percent are delayed?

```
ResultsByCity <- group_by(tidy_flights, City) %>% summarize(  
  TotalDelays=sum(NumFlightsDelayed),  
  TotalOnTime=sum(NumFlightsOnTime),  
  TotalFlights=sum(NumFlightsTotal),  
  PctDelayed=TotalDelays/TotalFlights,  
  PctOnTime=TotalOnTime/TotalFlights)  
ResultsByCity %>% kable() %>% kable_styling(c("striped", "bordered"))
```

City	TotalDelays	TotalOnTime	TotalFlights	PctDelayed	PctOnTime
LosAngeles	179	1191	1370	0.130657	0.869343
Phoenix	427	5061	5488	0.077806	0.922194
SanDiego	85	595	680	0.125000	0.875000
SanFrancisco	231	823	1054	0.219165	0.780835
Seattle	366	2042	2408	0.151993	0.848007

```
PhoenixDelays <- filter(.data = ResultsByCity, City=="Phoenix") %>% select(PctDelayed)  
SanFranciscoDelays <- filter(.data = ResultsByCity, City=="SanFrancisco") %>% select(PctDelayed)
```

This shows that the *smallest* percentage (7.8%) of flights are delayed at *Phoenix*, while the *largest* percentage (21.9%) of flights are delayed at *SanFrancisco*:

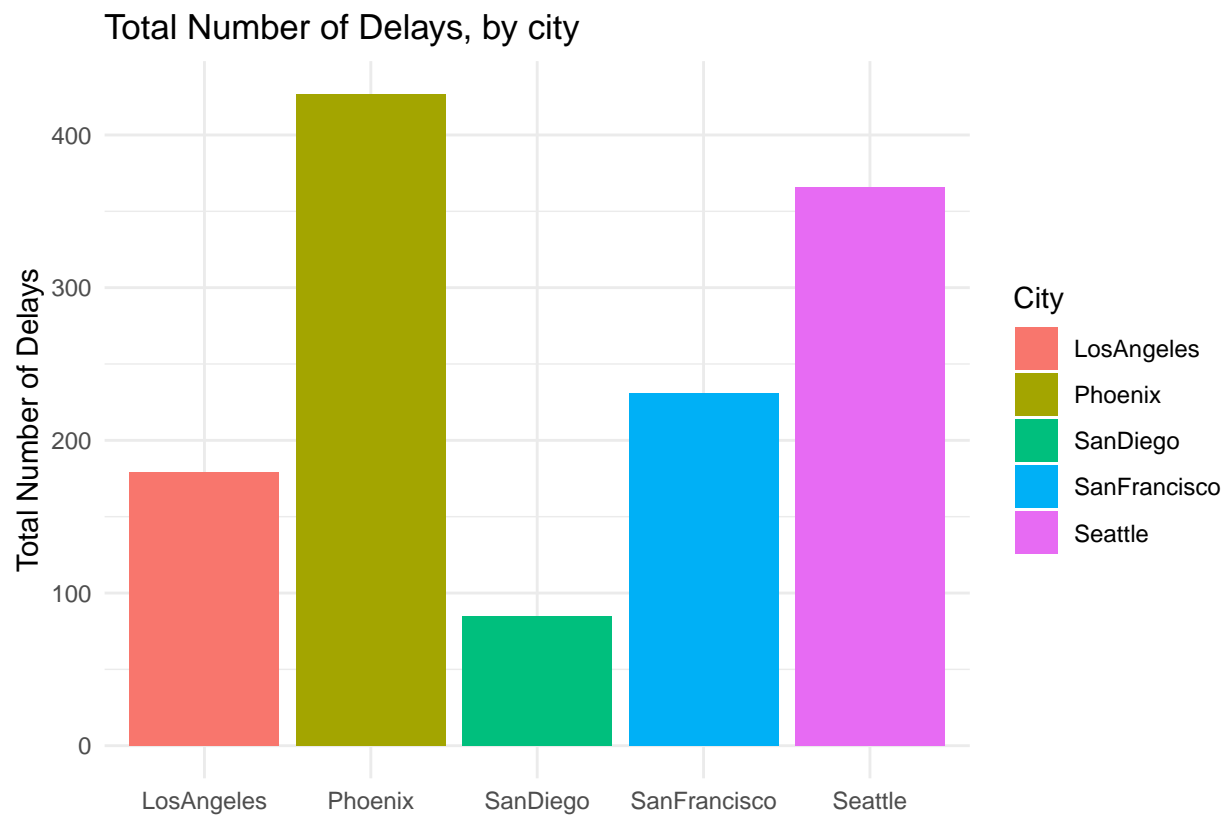
ResultsByCity, sorted by PctDelayed:

```
arrange(ResultsByCity, PctDelayed) %>% kable() %>% kable_styling(c("striped", "bordered"))
```

City	TotalDelays	TotalOnTime	TotalFlights	PctDelayed	PctOnTime
Phoenix	427	5061	5488	0.077806	0.922194
SanDiego	85	595	680	0.125000	0.875000
LosAngeles	179	1191	1370	0.130657	0.869343
Seattle	366	2042	2408	0.151993	0.848007
SanFrancisco	231	823	1054	0.219165	0.780835

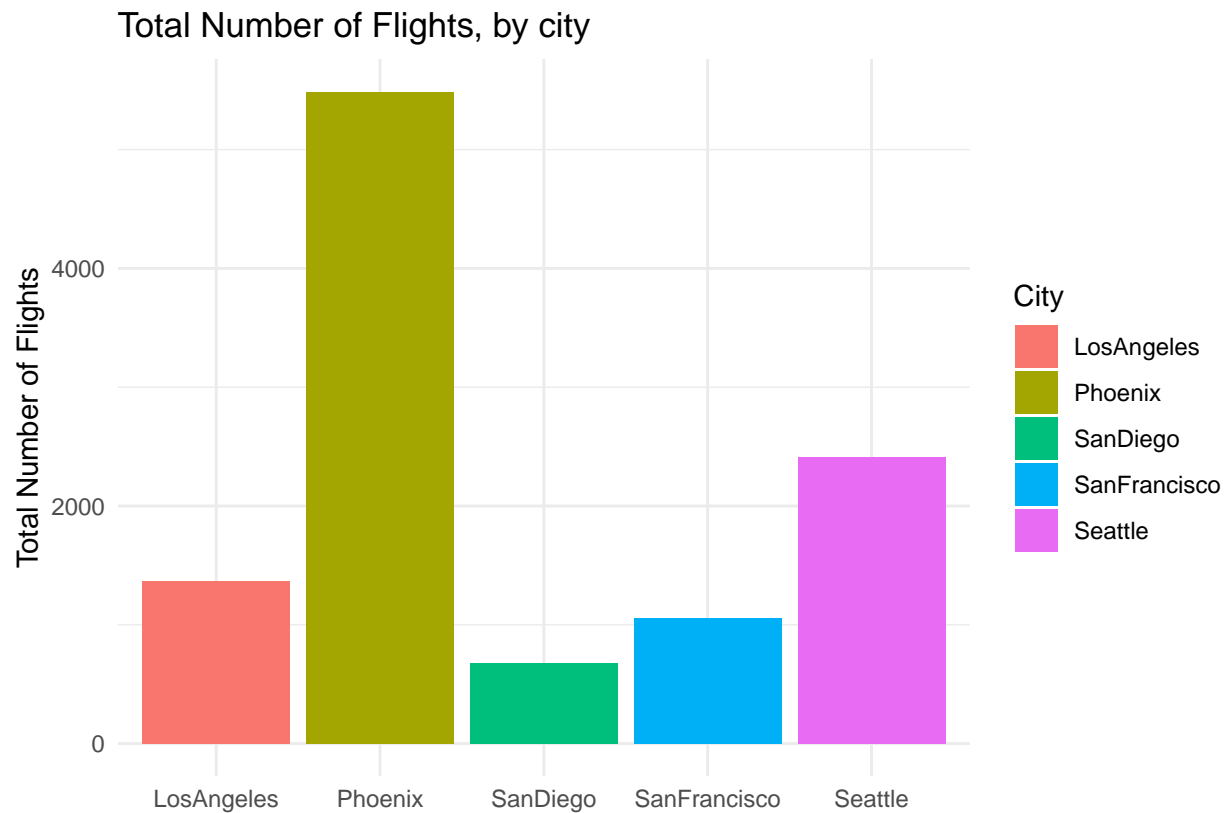
Here's a barplot:

```
select(.data = ResultsByCity, City, TotalDelays) %>%  
  ggplot(data = ., aes(factor(City), TotalDelays, fill = City)) +  
  geom_bar(stat="identity", position = "dodge") +  
  theme_minimal() +  
  labs(x="", y="Total Number of Delays") +  
  ggtitle("Total Number of Delays, by city")
```



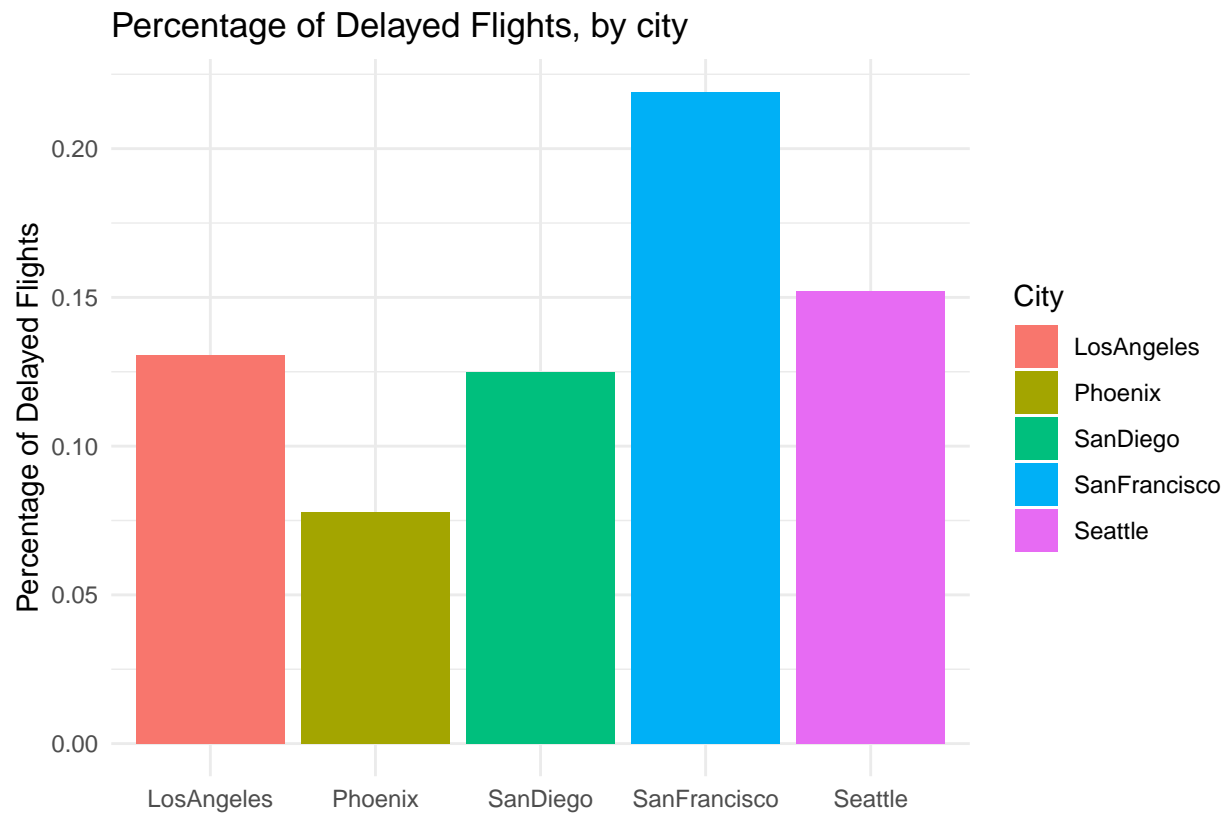
Although the largest absolute number of *delays* occurs in Phoenix, it is the city with the largest number of *overall* flights.

```
select(.data = ResultsByCity, City, TotalFlights) %>%  
  ggplot(data = ., aes(factor(City), TotalFlights, fill = City)) +  
  geom_bar(stat="identity", position = "dodge") +  
  theme_minimal() +  
  labs(x="", y="Total Number of Flights") +  
  ggtitle("Total Number of Flights, by city")
```



Indeed, the percentage of delays at Phoenix is the lowest across all 5 cities.

```
select(.data = ResultsByCity, City, PctDelayed) %>%  
  ggplot(data = ., aes(factor(City), PctDelayed, fill = City)) +  
  geom_bar(stat="identity", position = "dodge") +  
  theme_minimal() +  
  labs(x="", y="Percentage of Delayed Flights") +  
  ggtitle("Percentage of Delayed Flights, by city")
```



Relative market share

Let's determine the relative market share of each airline at each city, to see if that helps explain the delays:

Use merge to join the *Totals(by city)* onto the *tidy\_flights* dataframe.

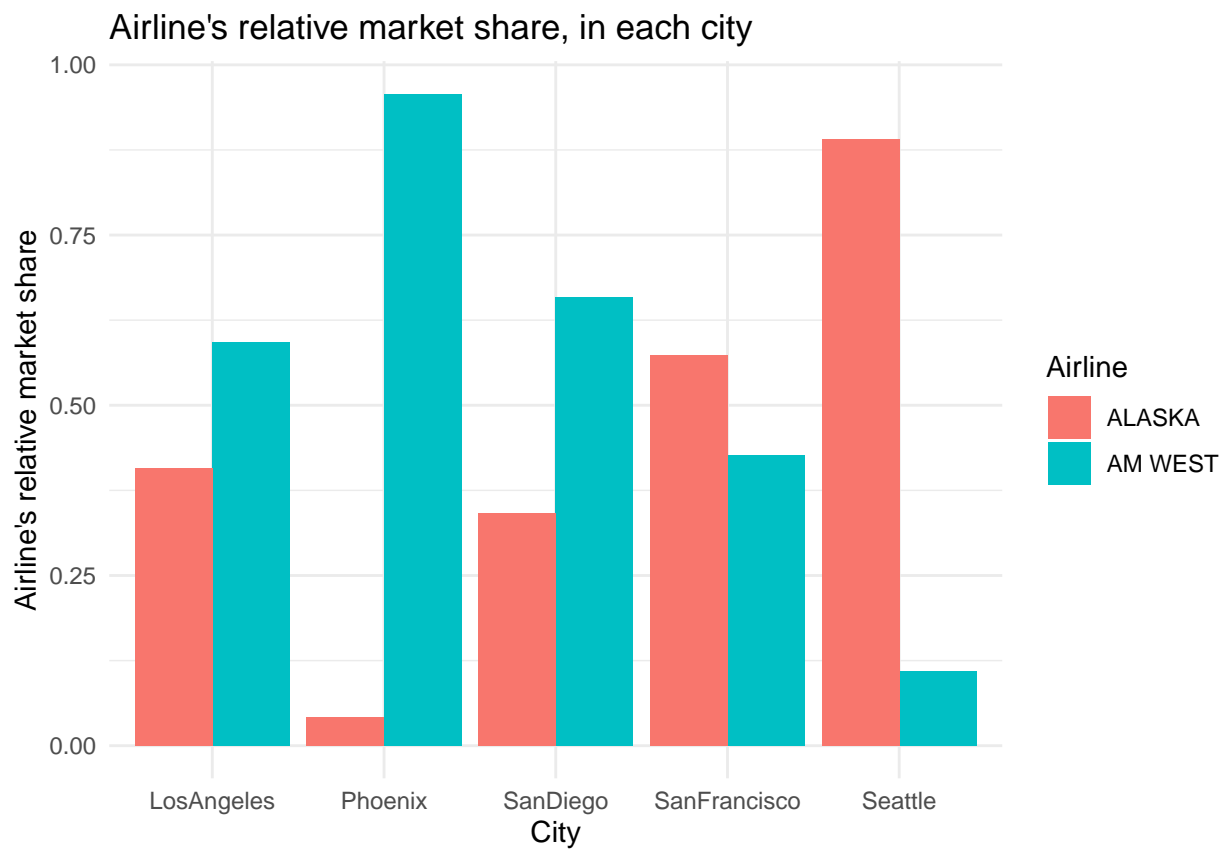
```
big_flights <- tidy_flights %>%  
  merge(ResultsByCity) %>%  
  arrange(.data = ., City, Airline) %>%  
  mutate(.data = ., ShareOfDelays=NumFlightsDelayed/TotalDelays,  
          ShareOfOnTime=NumFlightsOnTime/TotalOnTime,  
          ShareOfFlights=NumFlightsTotal/TotalFlights)  
big_flights %>% kable() %>% kable_styling(c("striped", "bordered"))
```

City	Airline	NumFlightsDelayed	NumFlightsOnTime	NumFlightsTotal	PctFlightsDelayed	PctFlightsOnTime	TotalDelays	TotalOnTime
LosAngeles	ALASKA	62	497	559	0.110912	0.889088	179	11
LosAngeles	AM WEST	117	694	811	0.144266	0.855734	179	11
Phoenix	ALASKA	12	221	233	0.051502	0.948498	427	50
Phoenix	AM WEST	415	4840	5255	0.078972	0.921028	427	50
SanDiego	ALASKA	20	212	232	0.086207	0.913793	85	5
SanDiego	AM WEST	65	383	448	0.145089	0.854911	85	5
SanFrancisco	ALASKA	102	503	605	0.168595	0.831405	231	8
SanFrancisco	AM WEST	129	320	449	0.287305	0.712695	231	8
Seattle	ALASKA	305	1841	2146	0.142125	0.857875	366	20
Seattle	AM WEST	61	201	262	0.232824	0.767176	366	20



Plot airline market share, by City:

```
select(.data = big_flights, City=City,Airline=Airline,ShareOfFlights=ShareOfFlights) %>%  
  ggplot(data = ., aes(factor(City), ShareOfFlights, fill = Airline)) +  
  geom_bar(stat="identity", position = "dodge") +  
  theme_minimal()+  
  labs( x="City", y="Airline's relative market share") +  
  ggtitle("Airline's relative market share, in each city")
```



The above helps clarify the picture.

### Good Weather vs. Bad Weather

*Phoenix* is known to be a city with “good weather” - it is a desert location where it seldom rains.

AM WEST (which no longer exists as an independent entity due to its 2005 merger with US Air, which then merged in 2015 with American Airlines) was based in Phoenix, which is where it dominated the market. It also had larger market share (than ALASKA) in both *San Diego* and *Los Angeles*, both comparatively “good weather” cities.

On the other hand, ALASKA Airlines is based in Seattle, *which is cloudy/rainy for more than 300 days per year* (the exception being July and August.) Alaska Airlines flies mainly up and down the west coast, including flights to Alaska (hence its name) plus San Francisco, where it had larger market share than AM West. *San Francisco is known for being foggy much of the time*, which results in a large percentage of flight delays there.

### Conclusion

Simpson’s Paradox explained:

The explanation for the paradox, “*How could one airline (ALASKA) have better on-time performance at each city, while the other airline (AM WEST) has better on-time performance overall?*” is found in the nature of cities in which each airline chooses to predominate, and the respective propensity for delays in such cities.

An airline which flies mainly to “bad weather” locations like *Seattle* and *San Francisco*, where a larger percentage of flights experience delays, is likely to have worse *overall* on-time performance when compared against an airline which flies mostly to “good weather” cities like *Phoenix*, *Los Angeles*, and *San Diego*.

Even if an airline boasts better on-time performance at each city, its overall performance can suffer because of its route map.

In this regard, ALASKA has won each of the “*battles*” (based on *within-city* comparison), but AM WEST has won the “*war*.”