

# DATA 607 Assignment 5

Adam Gersowitz

2/23/2020

## Introduction

This assignment is focused on **Tidying and Transforming** data for analysis. Prior to being transformed this data would be difficult to analyze and work with due to its format and inconsistencies.

## Importing the Data

I start by bringing in the .csv file from a github repository and making sure “air” is a dataframe. I make sure to convert all blank cells to null or “NA” values. I do this because functions such as fill will only work with “NA” cells.

```
## Loading required package: bitops
```

## Reshaping and Cleaning the Data

After the data has been imported, I begin by naming the airline and flight\_status columns as they were blank in the original dataset. Next I remove any lines that don't have data in them. Using the melt function I convert the dataframe from a wide format to a long format which makes it much easier to analyze. Next, I use the fill function to pull the airlines down to the blank cells below them. Finally, I rename the auto-populated variable and value fields and clean up the City field names. Now I am ready to analyze this data set.

```
##
## Attaching package: 'tidyr'

## The following object is masked from 'package:Rcurl':
##
##     complete

##
## 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
```

```
##
## Attaching package: 'reshape'

## The following object is masked from 'package:dplyr':
##
##      rename

## The following objects are masked from 'package:tidyr':
##
##      expand, smiths
```

## Reshaping Analyzing the Data

After I have transformed the dataset I will analyze it to determine which airlines are the most frequently on time and if they have any problems being on time for certain destination cities. First I create an aggregate table of flight information. I then perform a chi-square test via the `prop.test` function and find that the airlines are significantly different in the proportion of times they are on time with AM WEST being on time 89% of the time vs 86% for ALASKA. I also see that San Francisco is the destination city that most often has delays of flights. The worst Airline and destination combination is AM WEST and San Francisco at 71% on time. This is somewhat surprising because AM WEST is more often on time than ALASKA. This leads me to believe that the difficulty of having San Francisco as destination has caused this on time percentage to be surprisingly low and in turn has dragged AM WEST overall on time percentage down.

```
## Loading required package: gsubfn

## Loading required package: proto

## Loading required package: RSQLite

##
## 2-sample test for equality of proportions with continuity correction
##
## data:  table(air$total_airline_on_time, air$total_airline_flights)
## X-squared = 16.2, df = 1, p-value = 5.699e-05
## alternative hypothesis: two.sided
## 95 percent confidence interval:
##  0.9 1.0
## sample estimates:
## prop 1 prop 2
##      1      0
```

## Conclusion

After reshaping and analyzing this data set I have determined that AM WEST is more often on time than ALASKA airlines and that San Francisco is the destination city that most often leads to delays. To expand on this dataset it would be interesting to get the detail of each flight rather than a summary of on time and delayed flights. It would then be interesting to compare this data to external data such as time of year and weather to determine if those are having more of an impact on one airline over another. Additionally, it would be interesting to get a more robust dataset with more cities and airlines to determine if these are outliers or if in reality they are close to each other in performance when compared with all airlines.