

# Midterm Report

Haode Qi

March 21, 2017

lets load the library for data cleaning and read relevant data. Since we trying to figure where is the most dangerous places to work in Masschusetts, osha table and accident table will be the most relevant data we need for the purpose.

```
library(foreign)
library(dplyr)
library(magrittr)
library(tidyr)
```

```
osha <- read.dbf("osha.DBF")
```

```
accid<- read.dbf("accid.dbf")
```

Next step, we need to examine the data along with the documentation provided.

```
length(accid$ACTIVITYNO)
```

```
## [1] 2147
```

```
# we can see that there are 2000+ observations
```

```
length(unique(accid$ACTIVITYNO))
```

```
## [1] 1570
```

```
# but only 1570 unique observations, we need to look at the duplicated data and see what we can do about
```

```
#accid$ACTIVITYNO[duplicated(accid$ACTIVITYNO)]
```

```
noneunique<- subset(accid, duplicated(accid$ACTIVITYNO) ==TRUE)
```

```
# for the first four record, all the entry is exactly the same
```

```
head(noneunique,5)
```

```
##      ACTIVITYNO SITESTATE NAME  RELINSP  SEX DEGREE NATURE BODYPART SOURCE
## 2      10096592      MA <NA> 10096592 <NA>    3    21      04      16
## 3      10096592      MA <NA> 10096592 <NA>    3    21      04      16
## 4      10096592      MA <NA> 10096592 <NA>    3    21      04      16
## 5      10096592      MA <NA> 10096592 <NA>    3    21      04      16
## 15     10109288      MA <NA> 10109288 <NA>    2    12     21      29
##      EVENT ENVIRON HUMAN TASK HAZSUB OCC_CODE AGE
## 2      08      07    06    2    <NA>    000    0
## 3      08      07    06    2    <NA>    000    0
## 4      08      07    06    2    <NA>    000    0
## 5      08      07    06    2    <NA>    000    0
## 15     01      06    09    2    <NA>    000    0
```

```

# but if we examine the other record, not all the entries with the same activity number are the same
# the fact that this may occur because there may be multiple individuals involved in any accident and t
# we cannot simply removed the duplicated records

# a very important indicator of the seriousness of an accident is the degree. 1 for fatality, 2 for hosp
# it will probably be worthwhile to see if the same activityno also entails the same degree of seriousn
#if they are duplicated in the exact same way, the above statement is true

any((duplicated(accid$ACTIVITYNO) == duplicated(accid$DEGREE)) == FALSE)

```

```
## [1] TRUE
```

This table tells us that it is not

One way we can eliminate the duplication record is to accumulate the degree of seriousness and number of individuals involved in each incident.

```

accid$DEGREE%<>% as.character() %>% as.numeric()

uniqueaccid <- accid %>% group_by(ACTIVITYNO)

test3 <- accid %>% group_by(ACTIVITYNO) %>% filter(DEGREE == 3)
test2 <- accid %>% group_by(ACTIVITYNO) %>% filter(DEGREE == 2)
test1 <- accid %>% group_by(ACTIVITYNO) %>% filter(DEGREE == 1)

#fatality
test1 %<>% summarise(count = n())
colnames(test1) <- c("ACTIVITYNO", "Fatality")

#hospitality
test2 %<>% summarise(count = n())
colnames(test2) <- c("ACTIVITYNO", "Hospitality")

#non-hospitality
test3 %<>% summarise(count = n())
colnames(test3) <- c("ACTIVITYNO", "Non-hospitality")

uniqueaccid %<>% summarise(cumulativeDegree = sum(DEGREE), count = n())

#join them one by one
uniqueaccid <- uniqueaccid %>% left_join(test1, by = "ACTIVITYNO")
uniqueaccid <- uniqueaccid %>% left_join(test2, by = "ACTIVITYNO")
uniqueaccid <- uniqueaccid %>% left_join(test3, by = "ACTIVITYNO")

head(uniqueaccid,5)

```

```

## # A tibble: 5 × 6
##   ACTIVITYNO cumulativeDegree count Fatality Hospitality `Non-hospitality`
##       <int>           <dbl> <int>   <int>         <int>          <int>
## 1    141879             1     1     1           NA             NA
## 2    142349             1     1     1           NA             NA
## 3    142455             1     1     1           NA             NA

```

```
## 4      142737      1      1      1      NA      NA
## 5      159020      1      1      1      NA      NA
```

```
uniqueaccid %<>% mutate(averageDegree = round(cumulativeDegree/count,3))
```

*# the NAs in each type of injury simply means that is zero number of people involved in it. We should r*

```
uniqueaccid[is.na(uniqueaccid)] <-0
```

*#just to test to see if we do it correctly*

```
uniqueaccid[uniqueaccid$ACTIVITYNO == "10096592",]
```

```
## # A tibble: 1 × 7
```

```
##   ACTIVITYNO cumulativeDegree count Fatality Hospitality `Non-hospitality`
```

```
##   <int>          <dbl> <int>    <dbl>         <dbl>          <dbl>
```

```
## 1   10096592         15     5      0          0              5
```

```
## # ... with 1 more variables: averageDegree <dbl>
```

now we have the list of unique accidents and the cumulative seriousness and the associated number of individuals involved. the address information is stored in osha and we can retrieve the information by matching the activity number

In the osha documentation, the following are the address info. States are all MA, but i prefer to keep state and it become relevant later on in mapping. SITE STREET Street Address SITE STATE State Code (alpha postal abbreviation) SITE ZIP United States Postal Zip Code SITE CITY CODE Department of Commerce City Code SITE CNTY CODE

Since the primary purpose of this step is to retrieve the relevant information, NAs and duplications are no less of a concern here. we will worry about it later. Though some of the colnames seems to be overlylong like SITESTATE OR SITEADD, it is necessary to keep them the way it is incase we need to join any information with other table with the same colnames.

```
#glimpse(osha)
```

*#we may want to keep the date information. I actually did examine all the different dates available but*

```
#str(osha$OSHA1MOD)
```

```
#str(osha$OPENDATE)
```

```
#str(osha$CLOSEDT)
```

```
#str(osha$CLOSEDT2)
```

```
#str(osha$CLOSEDATE)
```

```
#str(osha$CLOSEDATE2)
```

```
#str(osha$OPENDT)
```

Next, we move on to the address information and match it by activity number.

```
oshaaddress <- osha[,c("ACTIVITYNO" , "SITEADD", "SITESTATE", "SITEZIP", "SITECITY", "SITECNTY" ) ]
```

```
uniqueaccid %<>% left_join(oshaaddress,by = "ACTIVITYNO")
```

*#just to illustrate some of the things we can do with the cleaned data, the following a frequency table*

```
summarise(group_by(uniqueaccid,SITEZIP),count= n())
```

```
## # A tibble: 383 × 2
##   SITEZIP count
##   <fctr> <int>
## 1    01001     5
## 2    01002     1
## 3    01005     1
## 4    01007     2
## 5    01008     1
## 6    01010     2
## 7    01011     1
## 8    01013     4
## 9    01020     6
## 10   01022     2
## # ... with 373 more rows
```

If we look at the lookup database, we can probably retrieve the city name information instead of encoding.

```
scc <- read.dbf("lookups/scc.dbf")

#by looking at scc and do some googling, the names are indeed city names, since some of the encoding ma
#glimpse(scc)
#head(scc)

colnames(scc) <- c("TYPE", "SITESTATE", "SITECNTY", "SITECITY", "NAME")

testname <- scc[,c("SITESTATE","SITECITY","NAME")]

uniqueaccid %<>% left_join(testname,by = c("SITESTATE","SITECITY"))

head(uniqueaccid)
```

```
## # A tibble: 6 × 13
##   ACTIVITYNO cumulativeDegree count Fatality Hospitality `Non-hospitality`
##   <int>          <dbl> <int>    <dbl>         <dbl>          <dbl>
## 1    141879          1     1        1           0            0
## 2    142349          1     1        1           0            0
## 3    142455          1     1        1           0            0
## 4    142737          1     1        1           0            0
## 5    159020          1     1        1           0            0
## 6    159319          1     1        1           0            0
## # ... with 7 more variables: averageDegree <dbl>, SITEADD <fctr>,
## #   SITESTATE <chr>, SITEZIP <fctr>, SITECITY <chr>, SITECNTY <fctr>,
## #   NAME <fctr>
```

```
dangerouscity<-summarise(group_by(uniqueaccid,NAME),incidence = n(), total= sum(count))

arrange(dangerouscity,desc(total))
```

```
## # A tibble: 257 × 3
##   NAME incidence total
##   <fctr>      <int> <int>
```

```
## 1      BOSTON      162    198
## 2      DANVERS      12     86
## 3      CAMBRIDGE     61     66
## 4      HAVERHILL     42     52
## 5      LAWRENCE      22     52
## 6      WORCESTER     44     52
## 7      SPRINGFIELD   29     48
## 8      BEVERLY      20     41
## 9      WINTHROP      13     39
## 10     WALTHAM       31     38
## # ... with 247 more rows
```

Not surprisingly, Boston has the highest number of accidents. Next, we may want to graph the findings in our data. And more intuitively, spatial visualization may be our best option. The following code will plot all the incidents and only the fatalities on two separate graphs.

```
library(ggmap)
library(ggplot2)

mapdata <- uniqueaccid[,c("ACTIVITYNO", "count", "Fatality", "SITESTATE", "NAME")]

head(mapdata)
```

```
## # A tibble: 6 × 5
##   ACTIVITYNO count Fatality SITESTATE      NAME
##   <int> <int>   <dbl>   <chr>   <fctr>
## 1    141879     1     1      MA    BRIMFIELD
## 2    142349     1     1      MA    WORCESTER
## 3    142455     1     1      MA    TURNERS FALLS
## 4    142737     1     1      MA    TEMPLETON
## 5    159020     1     1      MA    HOPKINTON
## 6    159319     1     1      MA    BROCKTON
```

```
mapdata$EXACT <- paste(mapdata$SITESTATE, mapdata$NAME)
```

```
test <- summarise(group_by(mapdata, EXACT), total = sum(count))
```

```
### this piece of code allows us to retrieve the longitude and latitude of the cities.
```

```
#It takes around 20 minutes to run and if you can run it yourself or simply run the file I wrote and save it
```

```
#ma.location<- geocode(test$EXACT)
```

```
#write.table(ma.location, file = "MA_Location")
```

```
###
```

```
# in case we lost the data
```

```
x <- read.table("MA_Location")
```

```
head(x)
```

```
##      lon      lat
## 1 -70.94532 42.10482
## 2 -71.43284 42.48509
```

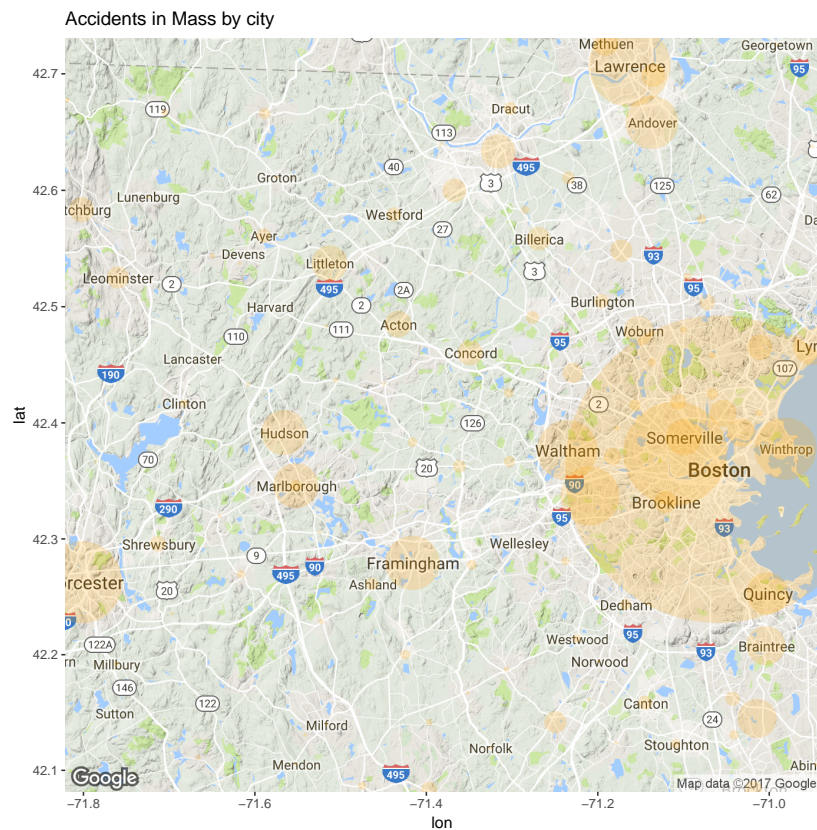
```
## 3 -70.89523 41.72237
## 4 -73.11743 42.62423
## 5 -72.61481 42.06954
## 6 -70.93004 42.85839
```

```
test$lon <- x$lon
test$lat <- x$lat
```

```
mass<- get_map("Massachusetts")
```

*# rmarkdown has an inherent issue with graphs, and the size of the dot is much larger than it should be  
#the graph more visually appealing*

```
ggmap(mass) +geom_point(aes(x=lon, y=lat), data=test, col="orange", alpha=0.2, size=(test$total)/2) + g
```



*#lets only concern about fatality.*

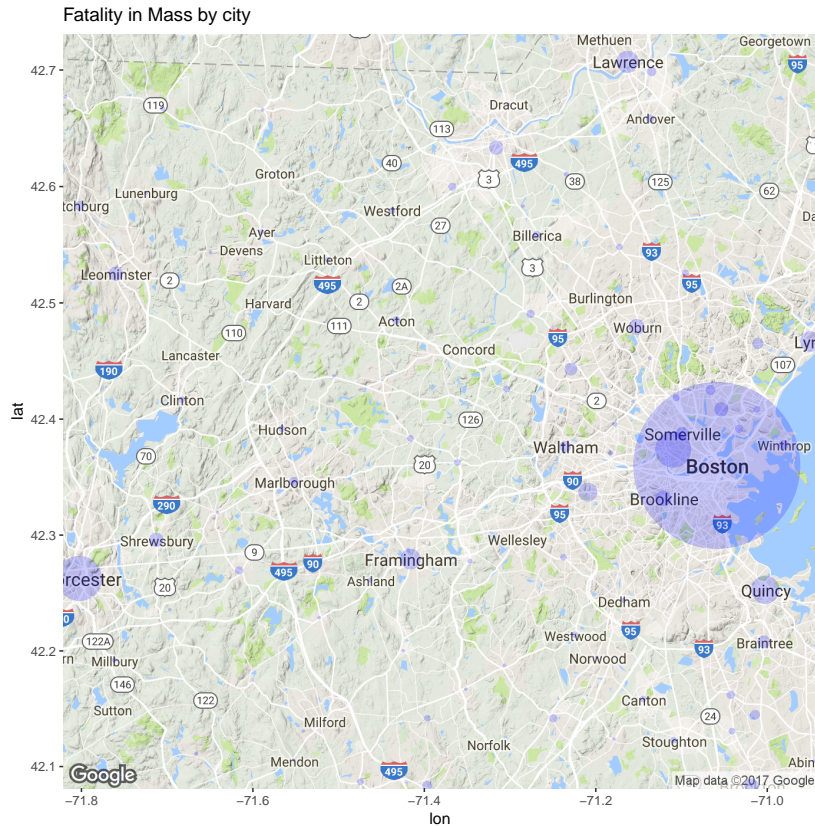
```
map_fatality <- summarise(group_by(mapdata,EXACT), total = sum(Fatality))
```

```
map_fatality$lon <- x$lon
map_fatality$lat <- x$lat
```

*# fatality rate is much lower and ,therefore, the size can be keep as the same.*

```
ggmap(mass) +geom_point(aes(x=lon, y=lat), data=map_fatality, col="blue", alpha=0.2, size=map_fatality$
```





Now we know how exactly does the mapping work. Let's have some fun and graph the accidents in Boston at a street level. Lets grab the data we needed and put it on a map.

```
#testing street level
```

```
testinfor <- uniqueaccid[,c( "ACTIVITYNO", "SITEADD", "SITESTATE", "NAME" )]
```

```
head(testinfor)
```

```
## # A tibble: 6 × 4
##   ACTIVITYNO SITEADD SITESTATE NAME
##   <int>      <fctr> <chr>   <fctr>
## 1    141879 OFF MONSON RD      MA    BRIMFIELD
## 2    142349 25 SAGAMORE RD      MA    WORCESTER
## 3    142455 CANAL STREET      MA    TURNERS FALLS
## 4    142737 MAIN ST      MA    TEMPLETON
## 5    159020 MESERVE ST      MA    HOPKINTON
## 6    159319 OAK HILL WAY      MA    BROCKTON
```

```
testinfor <- subset(testinfor, testinfor$NAME == "BOSTON")
```

```
head(testinfor)
```

```
## # A tibble: 6 × 4
##   ACTIVITYNO SITEADD SITESTATE NAME
##   <int>      <fctr> <chr>   <fctr>
```

```
## 1      922690      CAUSEWAY & LOMASNEY WAY      MA BOSTON
## 2      954883      CAUSEWAY & LOMASNEY WAY      MA BOSTON
## 3      954891 GALLIVAN, MORRISSEY, & NEPONSET      MA BOSTON
## 4      1803592 ST MARYS SCHOOL WARREN & PARK      MA BOSTON
## 5      1804889              31 BRIMMER ST      MA BOSTON
## 6      1808666              260 FRANKLIN STREET      MA BOSTON
```

```
testinfor$EXACT <- paste(testinfor$SITEADD,testinfor$NAME)

head(testinfor)
```

```
## # A tibble: 6 × 5
##   ACTIVITYNO      SITEADD SITESTATE  NAME
##   <int>          <fctr>    <chr> <fctr>
## 1      922690      CAUSEWAY & LOMASNEY WAY      MA BOSTON
## 2      954883      CAUSEWAY & LOMASNEY WAY      MA BOSTON
## 3      954891 GALLIVAN, MORRISSEY, & NEPONSET      MA BOSTON
## 4      1803592 ST MARYS SCHOOL WARREN & PARK      MA BOSTON
## 5      1804889              31 BRIMMER ST      MA BOSTON
## 6      1808666              260 FRANKLIN STREET      MA BOSTON
## # ... with 1 more variables: EXACT <chr>
```

```
testinfo <- summarise(group_by(testinfor,EXACT), count = n())
```

```
### again this piece of code is for loading data from online and I save a table just to save the time.
```

```
#boston_location <- geocode(testinfo$EXACT)
```

```
#write.table(boston_location, file = "Boston_location")
```

```
###
```

```
bos <- read.table("Boston_location")
```

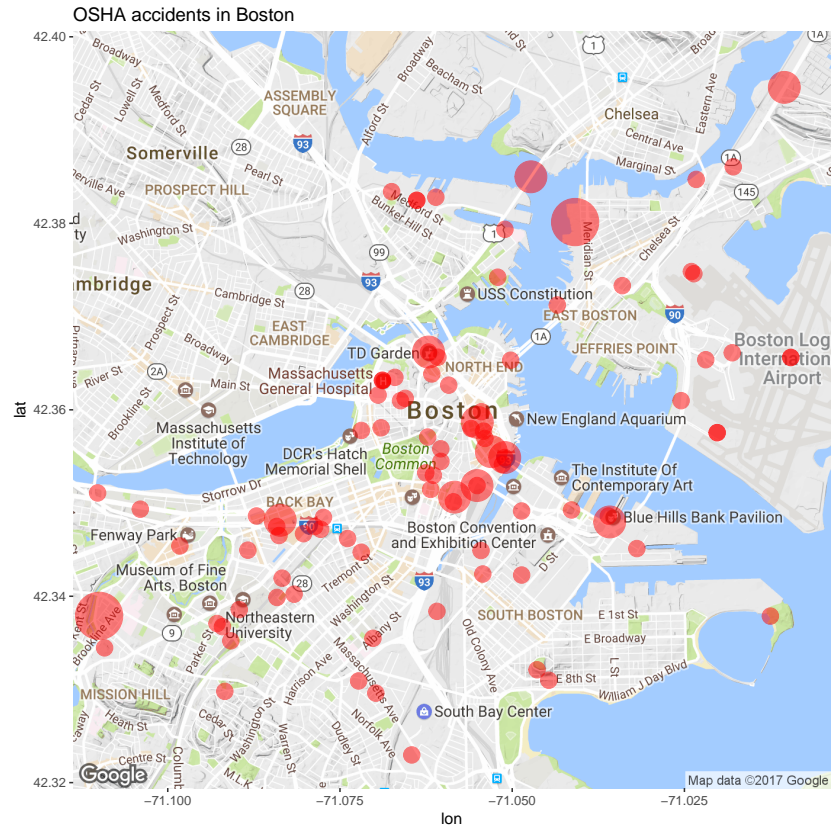
```
testinfo$lon <- bos$lon
```

```
testinfo$lat<- bos$lat
```

```
boston <- get_map("Boston", zoom =13)
```

```
ggmap(boston) + geom_point(aes(x = lon, y= lat), data = testinfo, col = "red" , alpha = .5 , size = tes
```





In the end, with our current data structure. We can easily join data by Activity number and do some interesting analysis. Below is an example of the most commonly types of injury.

```
# bar graph for jury types
```

```
acc<- read.dbf("lookups/acc.dbf")
```

```
ok <- subset(acc, acc$CATEGORY == "SOURC-INJ")
```

```
unique(ok$VALUE)
```

## [1] AIRCRAFT	AIR PRESSURE	ANIMAL/INS/REPT/ETC.
## [4] BOAT	BODILY MOTION	BOILER/PRESS VESSEL
## [7] BOXES/BARRELS, ETC.	BUILDINGS/STRUCTURES	CHEM LIQUIDS/VAPORS
## [10] CLEANING COMPOUND	COLD (ENVIR/MECH)	DIRT/SAND/STONE
## [13] DRUGS/ALCOHOL	DUST/PARTICLES/CHIPS	ELEC APPARAT/WIRING
## [16] FIRE/SMOKE	FOOD	FURNITURE/FURNISHING
## [19] GASES	GLASS	HAND TOOL (POWERED)
## [22] HAND TOOL (MANUAL)	HEAT (ENVIR/MECH)	HOISTING APPARATUS
## [25] LADDER	MACHINE	MATERIALS HANDLG EQ.
## [28] METAL PRODUCTS	MOTOR VEHICLE (HWY)	MOTOR VEHICLE(INDUS)
## [31] MOTORCYCLE	WIND/LIGHTNING, ETC.	FIREARM
## [34] PERSON	PETROLEUM PRODUCTS	PUMP/PRIME MOVER
## [37] RADIATION	TRAIN/RAILROAD EQUIP	VEGETATION
## [40] WASTE PRODUCTS	WATER	WORKING SURFACE
## [43] OTHER	FUME	MISTS
## [46] VIBRATION	NOISE	BIOLOGICAL AGENT

```
## 149 Levels: ABDOMEN ABSORPTION AIR PRESSURE AIRCRAFT ... WRIST(S)
```

```
colnames(ok) <- c("CATEGORY", "SOURCE", "VALUE")

ok$CATEGORY <- NULL

accid_injury<- left_join(accid, ok, by = "SOURCE")

accid_injury<- accid_injury[,c("ACTIVITYNO", "VALUE")]

accid_injury<- summarise(group_by(accid_injury, VALUE), count = n())

accid_injury$VALUE <- droplevels(accid_injury$VALUE)

ggplot(data = accid_injury, aes (count, fill = VALUE), ylab = "injury") + geom_bar(width =5) + coord_fl.
```

