

# CITS4009 Project 1: 2021

Code ▼

Amal Seby Chirackal (22829076)

## Introduction

This is a project involving the study of a dataset from the US Department of Labor's "US Accident Injury Dataset", which was downloaded from Data.gov. The total number of records in the collection is 202814, and it spans 15 years (2000 to 2015). This data set includes information about the accidents/illness/injuries occurs in various mines. Nature, cause, reason, time & location of accident, employee experience, days lost as result of accident and various parameters are recorded in the data.

## Exploratory study of data

In order to exploration of data, firstly required libraries and data set is loaded.

Hide

```
#load libraries
library(dplyr)
library(ggplot2)
#load data set in R environment
us_data <- read.csv("us_data.csv")
```

## Using R functions to explore the data

Hide

```
#structure of the data set
str(us_data)
```

```

'data.frame':  202814 obs. of  57 variables:
 $ MINE_ID      : int  100003 100003 100008 100011 100011 100011 100011 100011
100016 100021 ...
 $ CONTROLLER_ID : chr   "41044" "41044" "M31753" "M11763" ...
 $ CONTROLLER_NAME : chr   "Lhoist Group" "Lhoist Group" "Alan B Cheney" "Imerys S
A" ...
 $ OPERATOR_ID   : chr   "L13586" "L13586" "L31753" "L17074" ...
 $ OPERATOR_NAME : chr   "Lhoist North America " "Lhoist North America " "Cheney
Lime & Cement Company" "Imerys Pigments LLC" ...
 $ CONTRACTOR_ID : chr   "" "" "" "" "" ...
 $ DOCUMENT_NO   : num  2.2e+11 2.2e+11 2.2e+11 2.2e+11 2.2e+11 ...
 $ SUBUNIT_CD    : int   3 30 30 30 30 30 30 30 3 3 ...
 $ SUBUNIT       : chr   "STRIP, QUARY, OPEN PIT" "MILL OPERATION/PREPARATION PLA
NT" "MILL OPERATION/PREPARATION PLANT" "MILL OPERATION/PREPARATION PLANT" ...
 $ ACCIDENT_DT   : chr   "14/03/2012" "8/01/2007" "4/07/2009" "26/05/2000" ...
 $ CAL_YR        : int   2012 2007 2009 2000 2005 2006 2008 2012 2000 2006 ...
 $ CAL_QTR       : int   1 1 3 2 1 1 4 2 3 1 ...
 $ FISCAL_YR     : int   2012 2007 2009 2000 2005 2006 2009 2012 2000 2006 ...
 $ FISCAL_QTR    : int   2 2 4 3 2 2 1 3 4 2 ...
 $ ACCIDENT_TIME : int   945 1105 1000 1100 1430 1130 430 930 730 230 ...
 $ DEGREE_INJURY_CD : chr   "5" "6" "3" "5" ...
 $ DEGREE_INJURY : chr   "DAYS RESTRICTED ACTIVITY ONLY" "NO DYS AWY FRM WRK,NO R
STR ACT" "DAYS AWAY FROM WORK ONLY" "DAYS RESTRICTED ACTIVITY ONLY" ...
 $ FIPS_STATE_CD : int   1 1 1 1 1 1 1 1 1 1 ...
 $ UG_LOCATION_CD : chr   "?" "?" "?" "?" ...
 $ UG_LOCATION   : chr   "NO VALUE FOUND" "NO VALUE FOUND" "NO VALUE FOUND" "NO V
ALUE FOUND" ...
 $ UG_MINING_METHOD_CD : chr   "?" "?" "?" "?" ...
 $ UG_MINING_METHOD : chr   "NO VALUE FOUND" "NO VALUE FOUND" "NO VALUE FOUND" "NO V
ALUE FOUND" ...
 $ MINING_EQUIP_CD : chr   "24" "28" "?" "?" ...
 $ MINING_EQUIP   : chr   "Front-end loader, Tractor-shovel, Payloader, Highlift,
Skip loader" "Hand tools (not powered)" "NO VALUE FOUND" "NO VALUE FOUND" ...
 $ EQUIP_MFR_CD   : chr   "119" "121" "?" "?" ...
 $ EQUIP_MFR_NAME : chr   "Not on this list" "Not Reported" "NO VALUE FOUND" "NO V
ALUE FOUND" ...
 $ EQUIP_MODEL_NO : chr   "22321" "" "" "?" ...
 $ SHIFT_BEGIN_TIME : int   600 700 600 700 700 700 2300 700 700 1800 ...
 $ CLASSIFICATION_CD : chr   "12" "10" "18" "9" ...
 $ CLASSIFICATION : chr   "POWERED HAULAGE" "HANDTOOLS (NONPOWERED)" "SLIP OR FALL
OF PERSON" "HANDLING OF MATERIALS" ...
 $ ACCIDENT_TYPE_CD : chr   "21" "8" "30" "27" ...
 $ ACCIDENT_TYPE   : chr   "CGHT I, U, B, MVNG & STTN OBJS" "STRUCK BY, NEC" "OVER-
EXERTION, NEC" "OVER-EXERTION IN LIFTING OBJS" ...
 $ NO_INJURIES    : int   1 1 1 1 1 1 1 1 1 1 ...
 $ TOT_EXPER      : num  4.35 0.02 10 NA 0.87 ...
 $ MINE_EXPER     : num  4.35 0.02 2.15 0.23 0.87 ...
 $ JOB_EXPER      : num  0.67 0.02 2.15 0.23 0.38 ...
 $ OCCUPATION_CD  : chr   "374" "374" "374" "374" ...
 $ OCCUPATION     : chr   "Warehouseman, Bagger, Palletizer/Stacker, Store keeper,
Packager, Fabricator, Cleaning plant operator" "Warehouseman, Bagger, Palletizer/Stac
ker, Store keeper, Packager, Fabricator, Cleaning plant operator" "Warehouseman, Bagg
er, Palletizer/Stacker, Store keeper, Packager, Fabricator, Cleaning plant operator"
"Warehouseman, Bagger, Palletizer/Stacker, Store keeper, Packager, Fabricator, Cleani
ng plant operator" ...

```

```

$ ACTIVITY_CD      : chr  "28" "30" "13" "28" ...
$ ACTIVITY         : chr  "HANDLING SUPPLIES/MATERIALS" "HAND TOOLS (NOT POWERED)"
"CLIMB SCAFFOLDS/LADDERS/PLATFORMS" "HANDLING SUPPLIES/MATERIALS" ...
$ INJURY_SOURCE_CD : chr  "76" "46" "117" "4" ...
$ INJURY_SOURCE    : chr  "SURFACE MINING MACHINES" "AXE,HAMMER,SLEDGE" "GROUND"
"BAGS" ...
$ NATURE_INJURY_CD : chr  "160" "180" "330" "330" ...
$ NATURE_INJURY    : chr  "CONTUSN,BRUISE,INTAC SKIN" "CUT,LACER,PUNCT-OPN WOUND"
"SPRAIN,STRAIN RUPT DISC" "SPRAIN,STRAIN RUPT DISC" ...
$ INJ_BODY_PART_CD : chr  "700" "100" "520" "420" ...
$ INJ_BODY_PART    : chr  "MULTIPLE PARTS (MORE THAN ONE MAJOR)" "HEAD,NEC" "ANKL
E" "BACK (MUSCLES/SPINE/S-CORD/TAILBONE)" ...
$ SCHEDULE_CHARGE  : int   0 0 0 NA 0 0 0 0 NA 0 ...
$ DAYS_RESTRICT    : int   8 0 0 5 5 3 0 21 10 19 ...
$ DAYS_LOST        : int   0 0 9 NA 0 0 0 0 NA 13 ...
$ TRANS_TERM       : chr   "N" "N" "N" "N" ...
$ RETURN_TO_WORK_DT : chr   "03/26/2012" "1/09/2007" "07/14/2009" "6/01/2000" ...
$ IMMED_NOTIFY_CD  : chr   "? " "? " "? " "13" ...
$ IMMED_NOTIFY     : chr   "NO VALUE FOUND" "NO VALUE FOUND" "NO VALUE FOUND" "NOT
MARKED" ...
$ INVEST_BEGIN_DT  : chr   "" "" "" "" ...
$ NARRATIVE        : chr   "Employee was cleaning up at the Primary Crusher with th
e Dingo skid steer. The employee slipped and fell while"| __truncated__ "Handle of s
ledgehammer broke and head of hammer hit employee in the forehead." "EMPLOYEE WAS CLI
MBING DOWN A LADDER AND WHEN HE STEPPED TO THE GROUND HE SLIPPED AND SPRAINED HIS LEF
T ANKLE." "HE PULLED A BACK MUSCLE WHILE STACKING BAGS OF MATERIAL." ...
$ CLOSED_DOC_NO    : num   NA NA 3.2e+11 3.2e+11 NA ...
$ COAL_METAL_IND   : chr   "M" "M" "M" "M" ...

```

Data set contain 57 attributes with 202814 instances. 13 integer, 5 numeric and 39 character variables.

Hide

```

#summary of the data set
summary(us_data)

```

MINE_ID	CONTROLLER_ID	CONTROLLER_NAME	OPERATOR_ID	OPERATOR_NAME	
Min. : 100003	Length:202814	Length:202814	Length:202814	Length:202814	
1st Qu.:1300095	Class :character	Class :character	Class :character	Class :character	
Median :2602512	Mode :character	Mode :character	Mode :character	Mode :character	
Mean :2684336					
3rd Qu.:4400170					
Max. :5500008					
CONTRACTOR_ID	DOCUMENT_NO	SUBUNIT_CD	SUBUNIT	ACCIDENT_ID	
DT	CAL_YR				
Length:202814	Min. :2.200e+11	Min. : 1.000	Length:202814	Length:202814	
Min. :2000					
Class :character	1st Qu.:2.200e+11	1st Qu.: 1.000	Class :character	Class :character	
1st Qu.:2002					
Mode :character	Median :2.201e+11	Median : 3.000	Mode :character	Mode :character	
Median :2006					
	Mean :2.201e+11	Mean : 9.362			
Mean :2006					
	3rd Qu.:2.201e+11	3rd Qu.:17.000			
3rd Qu.:2010					
	Max. :2.202e+11	Max. :99.000			
Max. :2015					
CAL_QTR	FISCAL_YR	FISCAL_QTR	ACCIDENT_TIME	DEGREE_INJURY_CD	DEGREE_INJURY
Min. :1.000	Min. :2000	Min. :1.000	Min. : 1	Length:202814	Length:202814
1st Qu.:1.000	1st Qu.:2003	1st Qu.:2.000	1st Qu.: 845	Class :character	Class :character
Median :2.000	Median :2006	Median :3.000	Median :1230	Mode :character	Mode :character
Mean :2.453	Mean :2006	Mean :2.574	Mean :1886		
3rd Qu.:3.000	3rd Qu.:2010	3rd Qu.:4.000	3rd Qu.:1730		
Max. :4.000	Max. :2015	Max. :4.000	Max. :9999		
			NA's :1		
FIPS_STATE_CD	UG_LOCATION_CD	UG_LOCATION	UG_MINING_METHOD_CD	UG_MINING_METHOD	
Min. : 1.00	Length:202814	Length:202814	Length:202814	Length:202814	
1st Qu.:18.00	Class :character	Class :character	Class :character	Class :character	
Median :32.00	Mode :character	Mode :character	Mode :character	Mode :character	
Mean :32.31					
3rd Qu.:49.00					
Max. :78.00					
MINING_EQUIP_CD	MINING_EQUIP	EQUIP_MFR_CD	EQUIP_MFR_NAME	EQUIP_MODEL_NO	
Length:202814	Length:202814	Length:202814	Length:202814	Length:202814	

Class :character haracter	Class :character	Class :character	Class :character	Class :c
Mode :character haracter	Mode :character	Mode :character	Mode :character	Mode :c
SHIFT_BEGIN_TIME	CLASSIFICATION_CD	CLASSIFICATION	ACCIDENT_TYPE_CD	ACCIDENT_T
Min. : 1	Length:202814	Length:202814	Length:202814	Length:202814
1st Qu.: 630	Class :character	Class :character	Class :character	Class :cha
Median : 700	Mode :character	Mode :character	Mode :character	Mode :cha
Mean :1384				
3rd Qu.:1545				
Max. :9999				
NA's :991				
NO_INJURIES	TOT_EXPER	MINE_EXPER	JOB_EXPER	OCCUPATION_CD
OCCUPATION				
Min. : 0.0000	Min. : 0.01	Min. : 0.01	Min. : 0.01	Length:202814
Length:202814				
1st Qu.: 1.0000	1st Qu.: 2.00	1st Qu.: 0.69	1st Qu.: 1.00	Class :character
Class :character				
Median : 1.0000	Median : 7.00	Median : 2.77	Median : 3.23	Mode :character
Mode :character				
Mean : 0.9039	Mean :11.32	Mean : 6.63	Mean : 6.98	
3rd Qu.: 1.0000	3rd Qu.:20.00	3rd Qu.: 8.77	3rd Qu.:10.00	
Max. :36.0000	Max. :65.00	Max. :65.00	Max. :65.00	
	NA's :37400	NA's :34325	NA's :33746	
ACTIVITY_CD	ACTIVITY	INJURY_SOURCE_CD	INJURY_SOURCE	NATURE_I
NJURY_CD				
Length:202814	Length:202814	Length:202814	Length:202814	Length:202814
Class :character	Class :character	Class :character	Class :character	Class :c
haracter				
Mode :character	Mode :character	Mode :character	Mode :character	Mode :c
haracter				
NATURE_INJURY	INJ_BODY_PART_CD	INJ_BODY_PART	SCHEDULE_CHARGE	DAYS_REST
RICT				
Length:202814	Length:202814	Length:202814	Min. : 0.00	Min. :
0.00				
Class :character	Class :character	Class :character	1st Qu.: 0.00	1st Qu.:
0.00				
Mode :character	Mode :character	Mode :character	Median : 0.00	Median :
0.00				
			Mean : 71.93	Mean :
7.98				
			3rd Qu.: 0.00	3rd Qu.:
4.00				

First 5 rows of the data set to have glance of the data are bellow:

```
head(us_data, 5)
```

# Data cleaning

## Dealing with NA values

## Checking for missing values in each column of data set

6/18

```
sapply(us_data,function(x)sum(is.na(x)))
```

MINE_ID	CONTROLLER_ID	CONTROLLER_NAME	OPERATOR_ID	
OPERATOR_NAME				
0	0	0	0	
0				
CONTRACTOR_ID	DOCUMENT_NO	SUBUNIT_CD	SUBUNIT	
ACCIDENT_DT				
0	0	0	0	
0				
CAL_YR	CAL_QTR	FISCAL_YR	FISCAL_QTR	
ACCIDENT_TIME				
0	0	0	0	
1				
DEGREE_INJURY_CD	DEGREE_INJURY	FIPS_STATE_CD	UG_LOCATION_CD	
UG_LOCATION				
0	0	0	0	
0				
UG_MINING_METHOD_CD	UG_MINING_METHOD	MINING_EQUIP_CD	MINING_EQUIP	
EQUIP_MFR_CD				
0	0	0	0	
0				
EQUIP_MFR_NAME	EQUIP_MODEL_NO	SHIFT_BEGIN_TIME	CLASSIFICATION_CD	
CLASSIFICATION				
0	48	991	0	
0				
ACCIDENT_TYPE_CD	ACCIDENT_TYPE	NO_INJURIES	TOT_EXPER	
MINE_EXPER				
0	0	0	37400	
34325				
JOB_EXPER	OCCUPATION_CD	OCCUPATION	ACTIVITY_CD	
ACTIVITY				
33746	0	0	0	
0				
INJURY_SOURCE_CD	INJURY_SOURCE	NATURE_INJURY_CD	NATURE_INJURY	IN
J_BODY_PART_CD				
0	0	0	0	
0				
INJ_BODY_PART	SCHEDULE_CHARGE	DAYS_RESTRICT	DAYS_LOST	
TRANS_TERM				
0	65006	54856	39677	
0				
RETURN_TO_WORK_DT	IMMED_NOTIFY_CD	IMMED_NOTIFY	INVEST_BEGIN_DT	
NARRATIVE				
0	0	0	0	
0				
CLOSED_DOC_NO	COAL_METAL_IND			
116621	0			

Variables **EQUIP\_MODEL\_NO**, **SHIFT\_BEGIN\_TIME**, **TOT\_EXPER**, **MINE\_EXPER**, **JOB\_EXPER**, **SCHEDULE\_CHARGE**, **DAYS\_RESTRICT**, **DAYS\_LOST**, **CLOSED\_DOC\_NO** have missing values.

# Dealing with Anamolies

checking for outliers in the data.

Hide

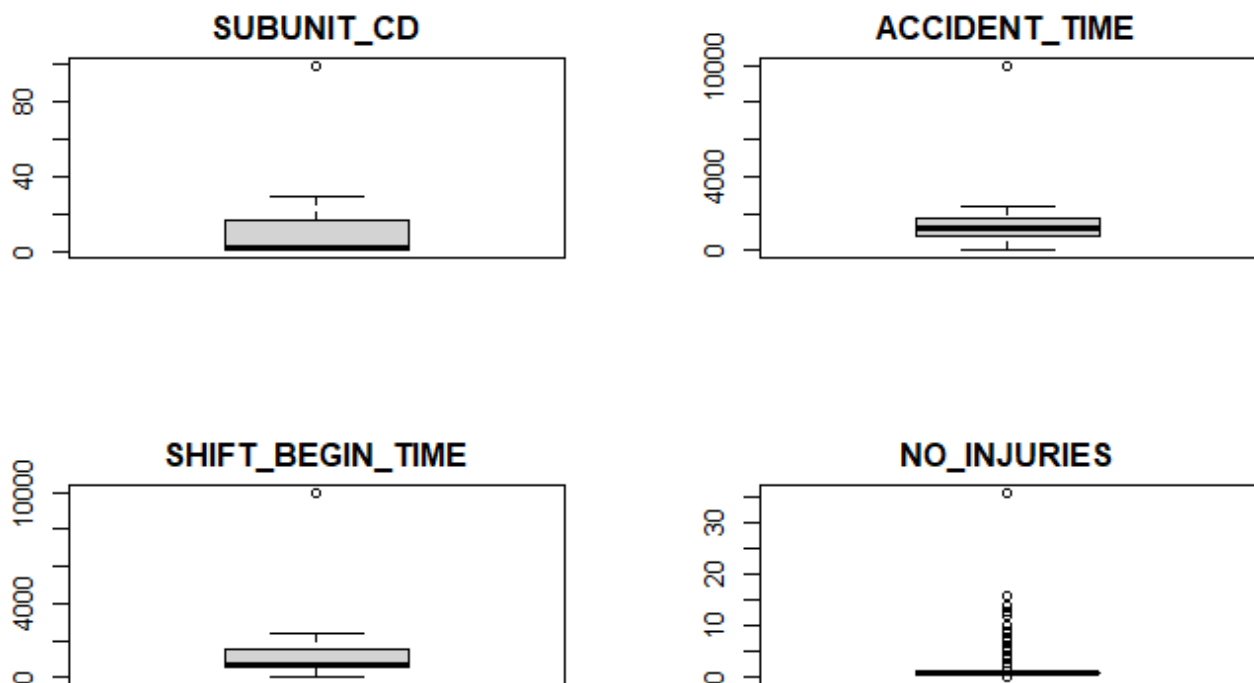
```
par(mfrow = c(2, 2))
boxplot(us_data$SUBUNIT_CD, main= "SUBUNIT_CD")
boxplot(us_data$ACCIDENT_TIME, main= "ACCIDENT_TIME")
```

Hide

```
boxplot(us_data$SHIFT_BEGIN_TIME, main= "SHIFT_BEGIN_TIME")
boxplot(us_data$NO_INJURIES, main= "NO_INJURIES")
```

Hide

```
par(mfrow = c(2, 2))
```



Hide

```
barplot(us_data$TOT_EXPER, main= "TOT_EXPER")
barplot(us_data$MINE_EXPER, main= "MINE_EXPER")
```

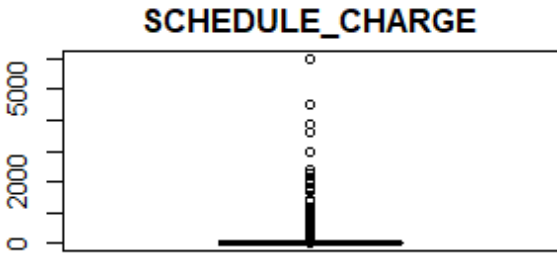
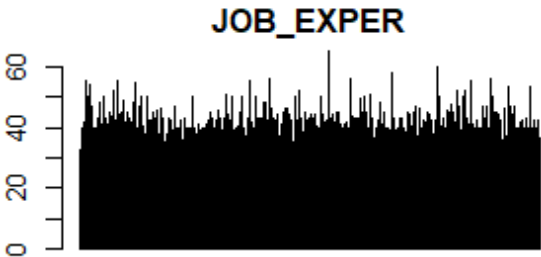
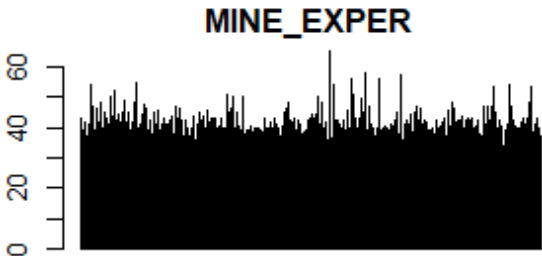
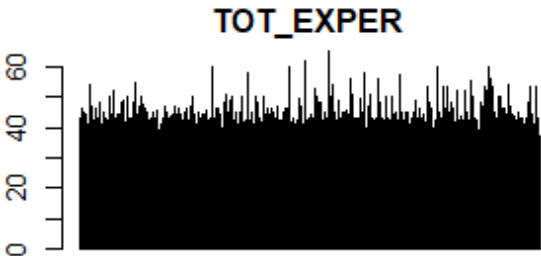
Hide

```
barplot(us_data$JOB_EXPER, main= "JOB_EXPER")
boxplot(us_data$SCHEDULE_CHARGE, main= "SCHEDULE_CHARGE")
```

Hide

```
par(mfrow = c(2, 2))
```



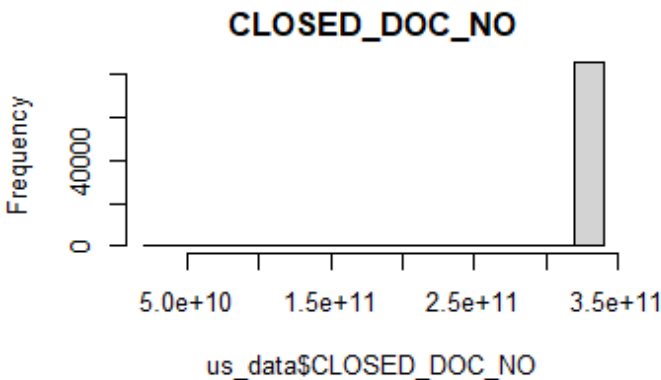
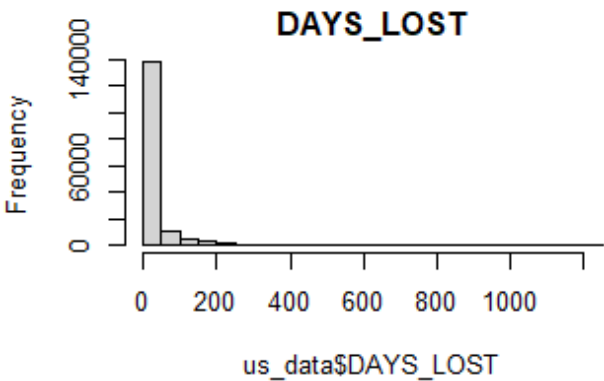
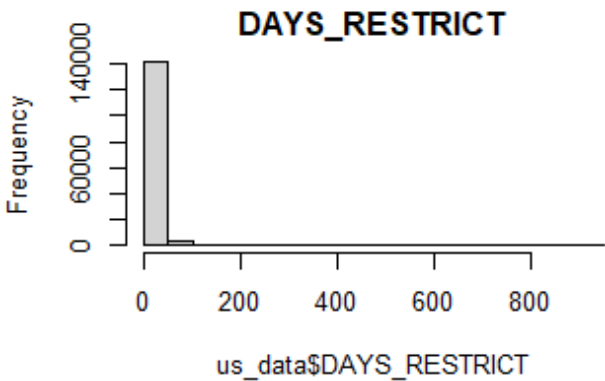


Hide

```
hist(us_data$DAYS_RESTRICT, main= "DAYS_RESTRICT")
hist(us_data$DAYS_LOST, main= "DAYS_LOST")
```

Hide

```
hist(us_data$CLOSED_DOC_NO, main= "CLOSED_DOC_NO")
```



SUBUNIT\_CD, ACCIDENT\_TIME, SHIFT\_BEGIN\_TIME, NO\_INJURIES, TOT\_EXPER, MINE\_EXPER, JOB\_EXPER, SCHEDULE\_CHARGE, DAYS\_RESTRICT, DAYS\_LOST and CLOSED\_DOC\_NO have outliers.

Hide

```
#replacing outliers with missing values
for (i in c(8,15,28,33,34,35,36,47,48,49,56)){
  qntile <- quantile(us_data[,i], probs=c(.25, .75),na.rm = TRUE)
  H <- 1.5 * IQR(us_data[,i], na.rm = T)
  us_data[,i][us_data[,i] < (qntile[1] - H)] <- NA
  us_data[,i][us_data[,i] > (qntile[2] + H)] <- NA
}
```

Hide

```
#replacing "?"
us_data$MINING_EQUIP_CD[us_data$MINING_EQUIP_CD == "?"] <- NA
us_data$EQUIP_MFR_CD[us_data$EQUIP_MFR_CD == "?"] <- NA
```

Hide

```
#removing NA values
us_data <- na.omit(us_data)
```

## Data transformation

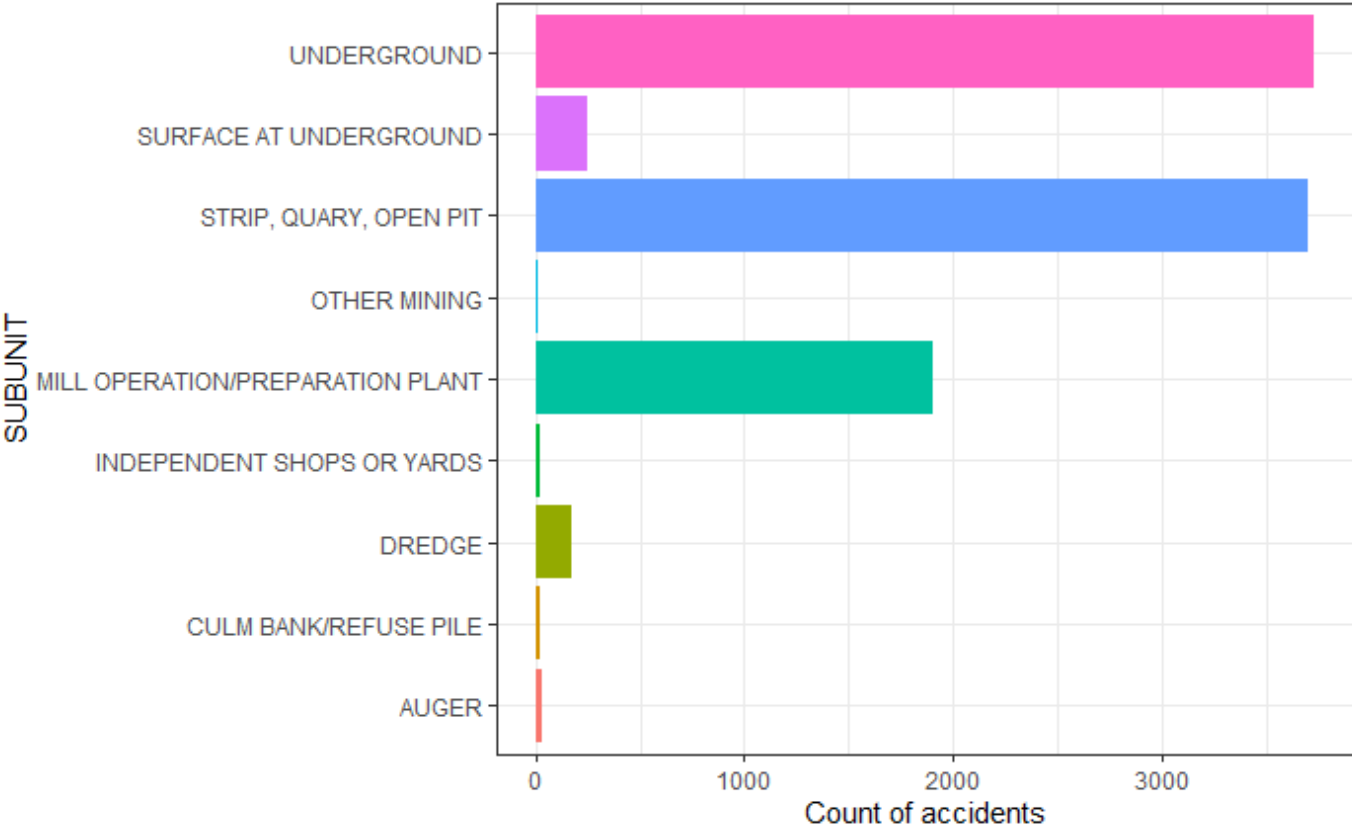
Hide

```
#transform character type variables to factor type
us_data[sapply(us_data, is.character)] <- lapply(us_data[sapply(us_data, is.character)], as.factor)
```

## Visualization

Hide

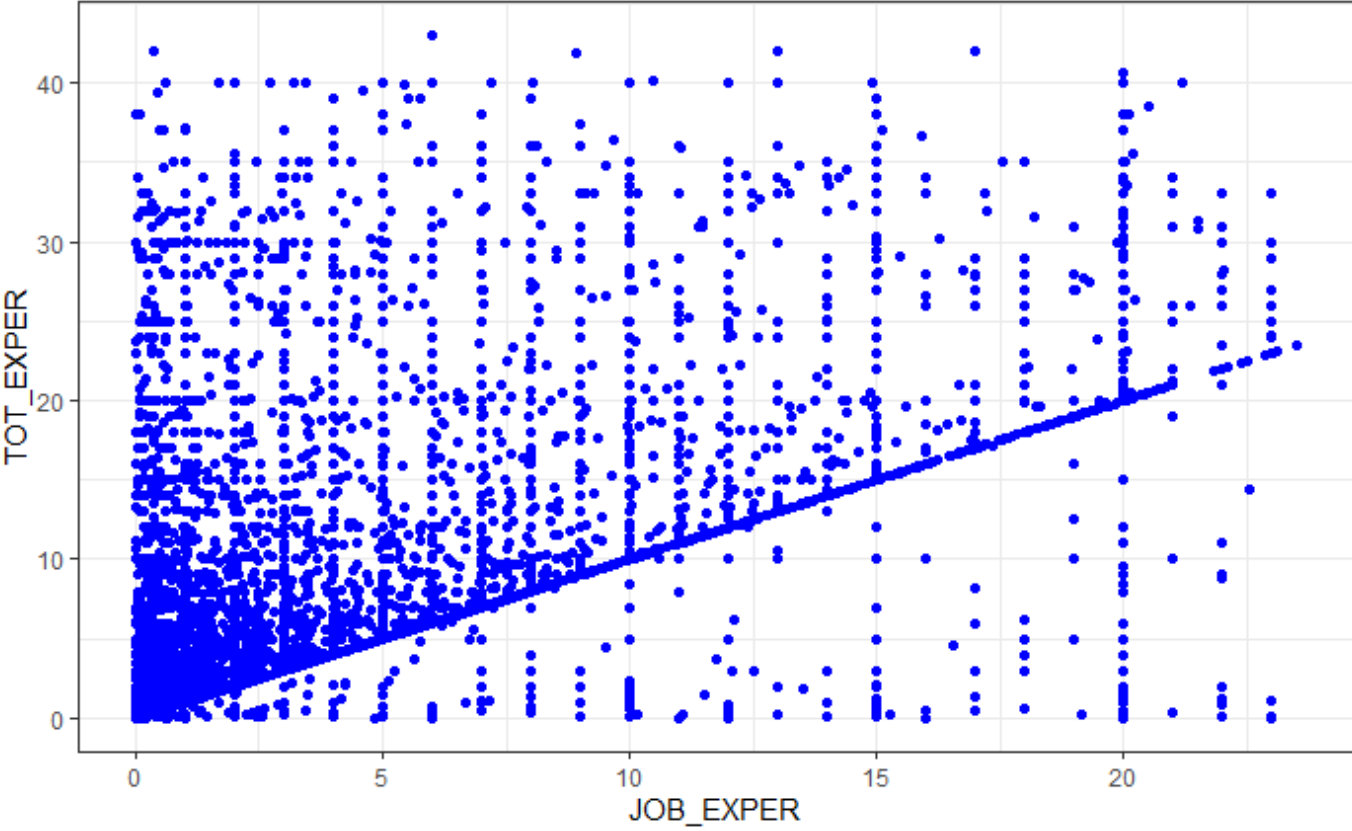
```
# location within a mine where the accident/injury/illness occurred
ggplot(data = us_data,aes(SUBUNIT,fill=SUBUNIT )) +geom_bar() + coord_flip() +theme_bw() +
  theme(legend.position = "none") +ylab("Count of accidents")
```



underground,strip and open pit locations have more number of accidents as compare to other locations.

Hide

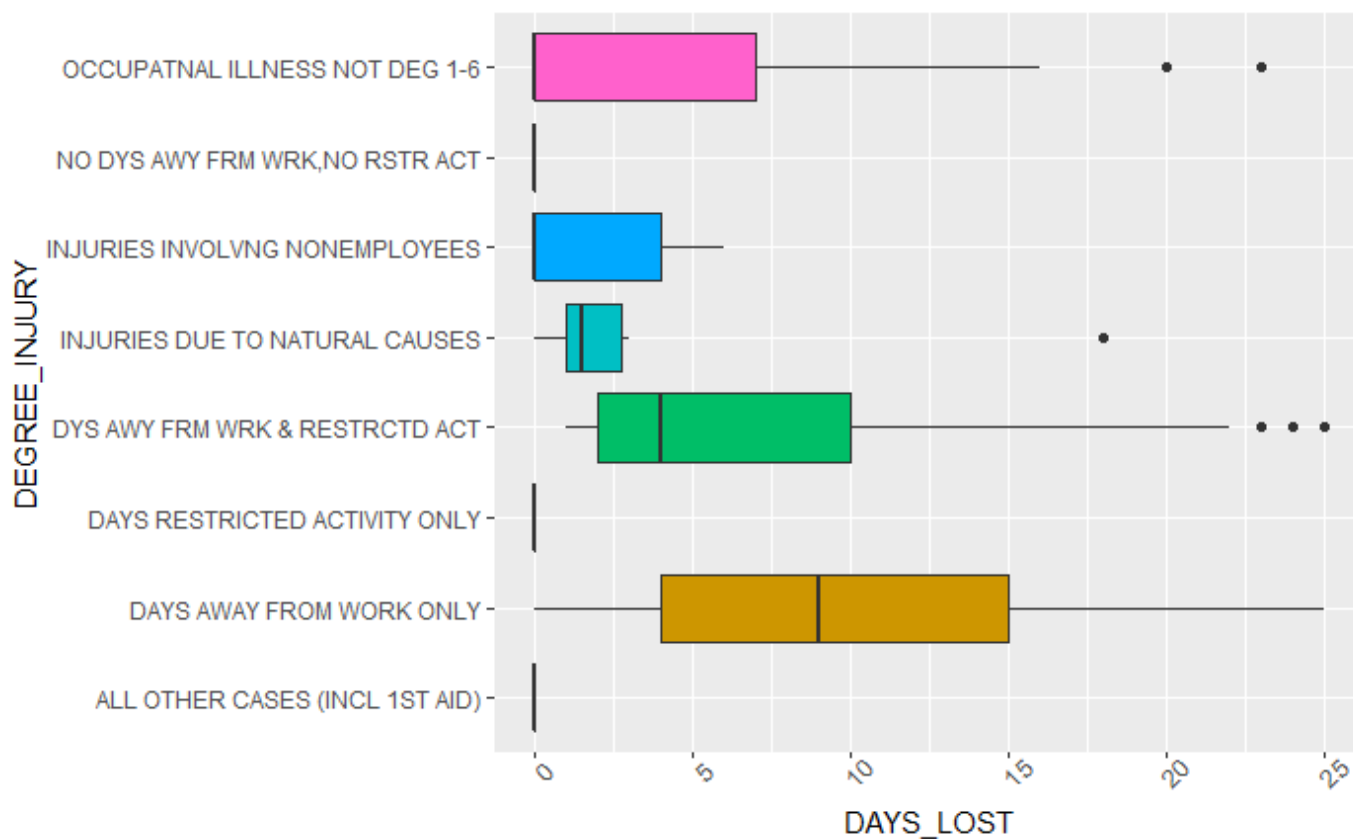
```
#Relation of total mining experience of person with Experience in the job title of th
e person
ggplot(data = us_data,aes(JOB_EXPER,TOT_EXPER)) + geom_point(col= "blue") +
  theme_bw()
```



There is a quite vague relation between JOB\_EXPER & TOT\_EXPE, but most of the employees have more Experience in the job title with more total mining experience.

[Hide](#)

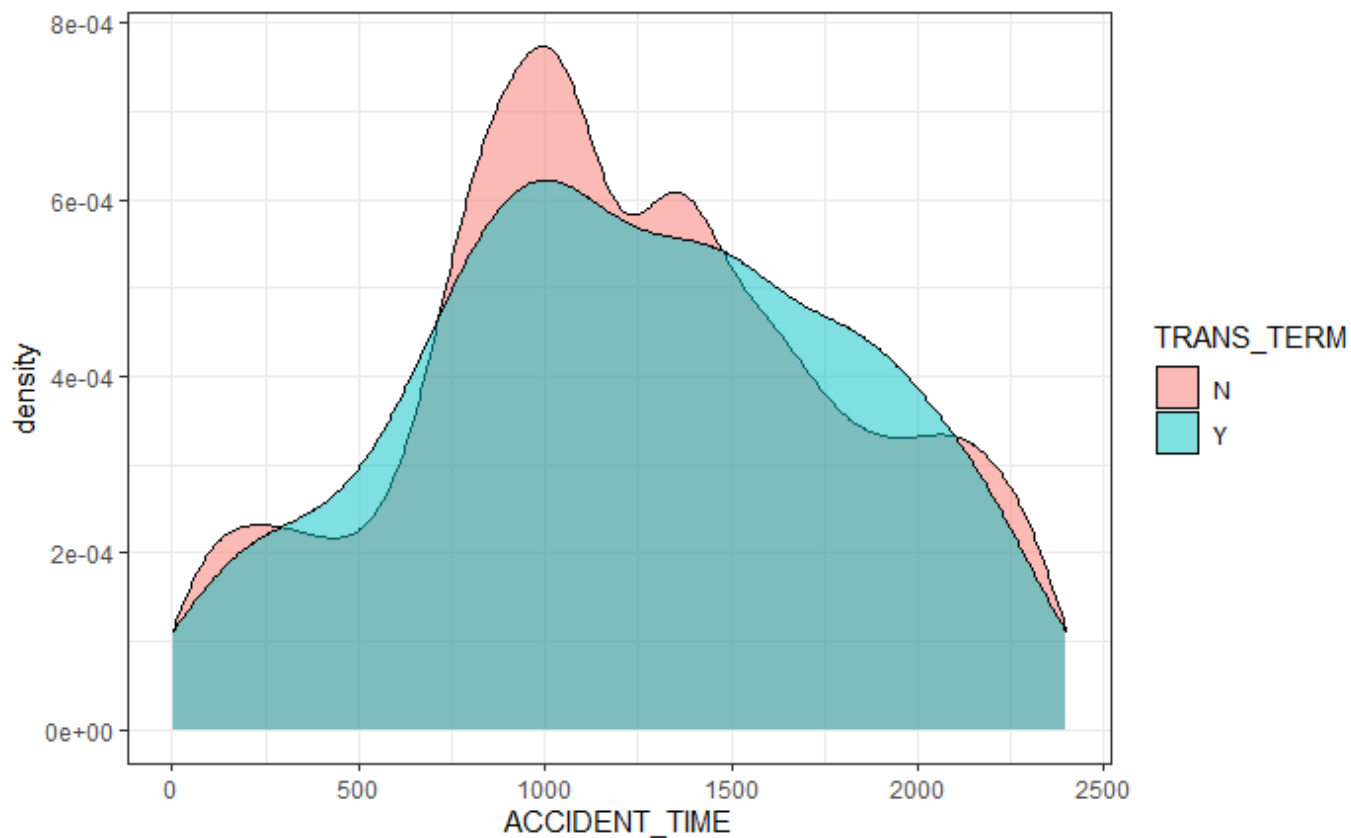
```
#Actual days lost from work due to the injury/illness
ggplot(data = us_data,aes(DEGREE_INJURY,DAYS_LOST,fill= DEGREE_INJURY)) + geom_boxplot() +
  theme(axis.text.x = element_text(angle=45), legend.position = "none") +
  ylim(0,25) +coord_flip()
```



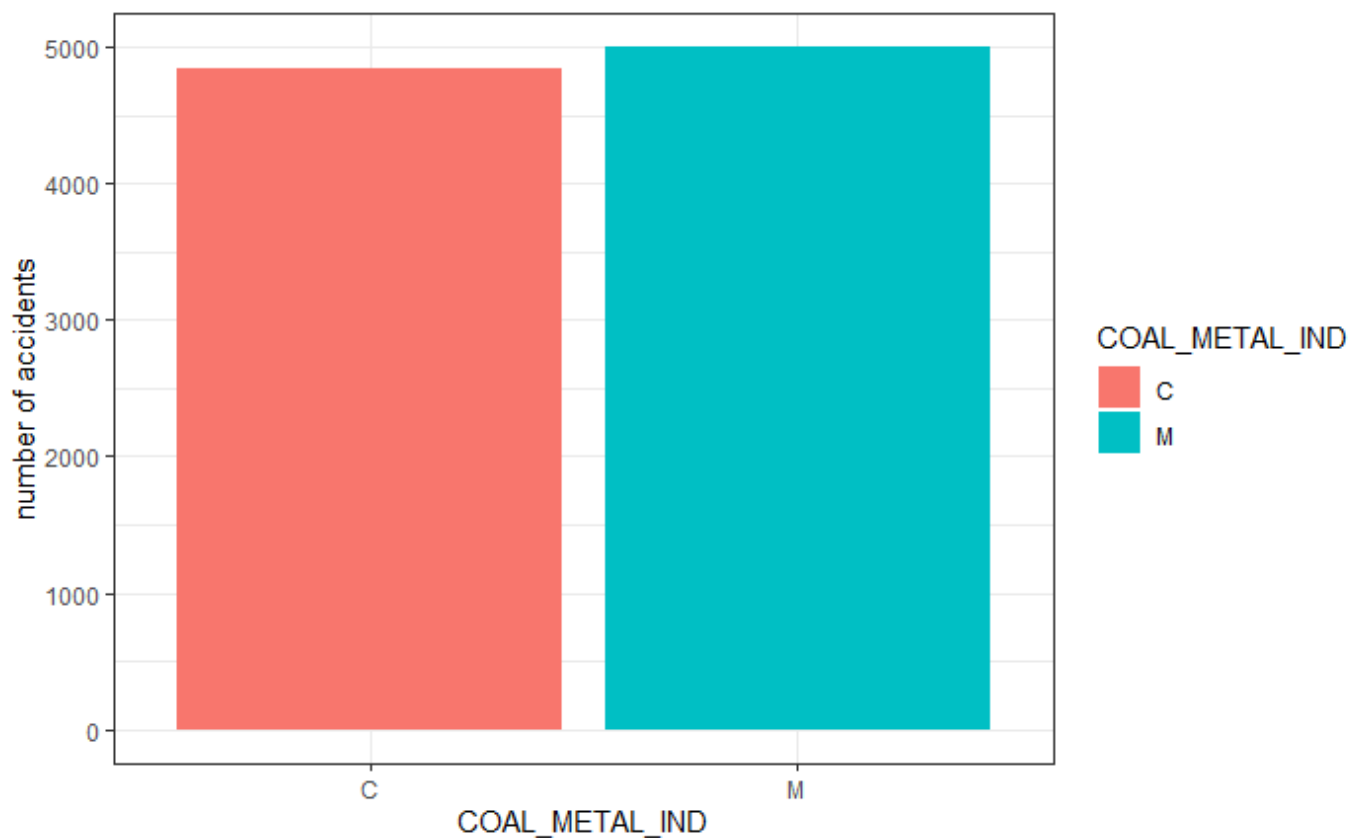
Degree of injury named Days away from work only has more Days lost.

[Hide](#)

```
#Distribution of accident time for employees who are terminated or not
ggplot(data = us_data,aes(ACCIDENT_TIME, fill= TRANS_TERM)) +
  geom_density(alpha= 0.5) +theme_bw()
```

[Hide](#)

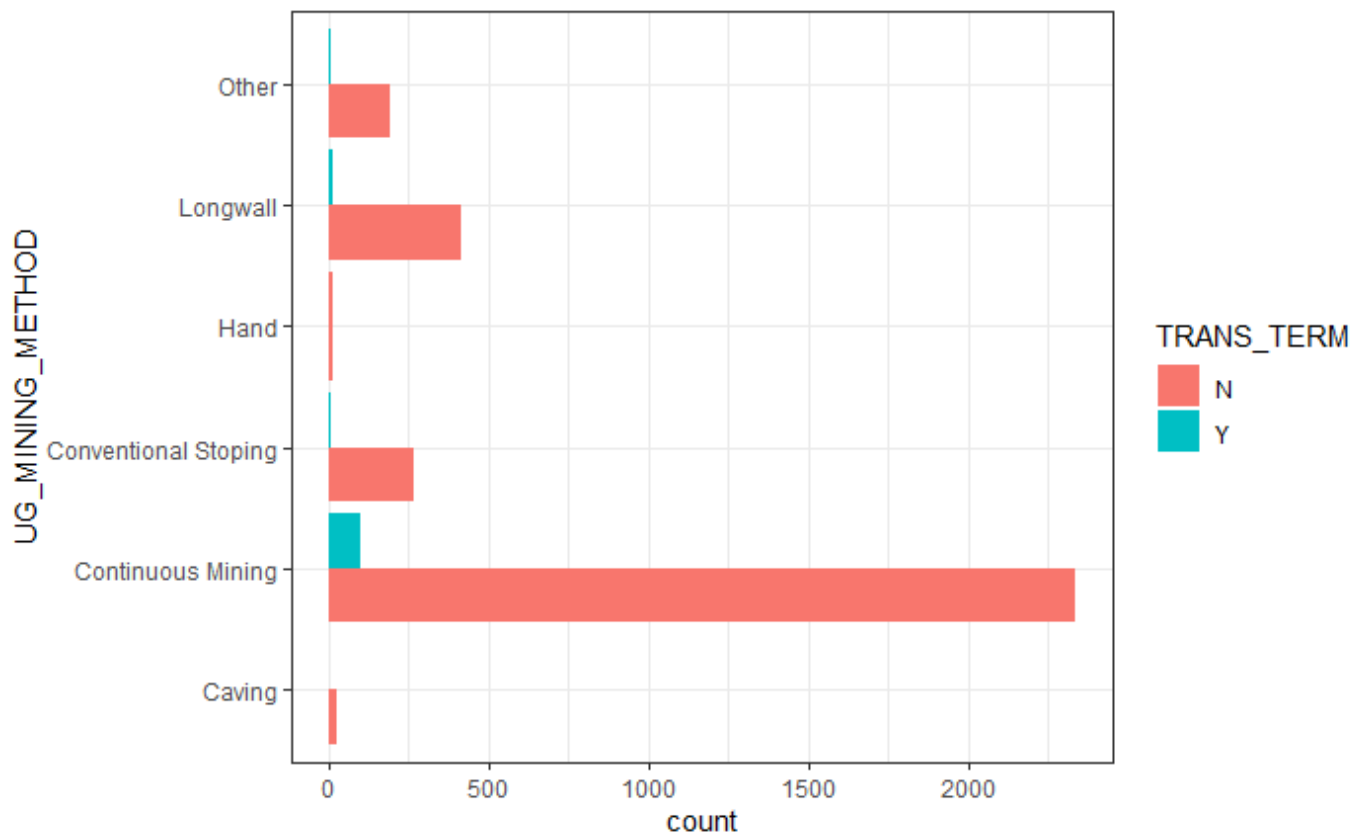
```
# Identifies if the accident occurred at a Coal or Metal/Non-Metal mine.  
ggplot(data = us_data, aes(COAL_METAL_IND, fill= COAL_METAL_IND)) + geom_histogram(stat  
="count")+  
  theme_bw() + ylab("number of accidents")
```



Metal/Non-Metal mines have slightly more accidents as compare to coal mines.

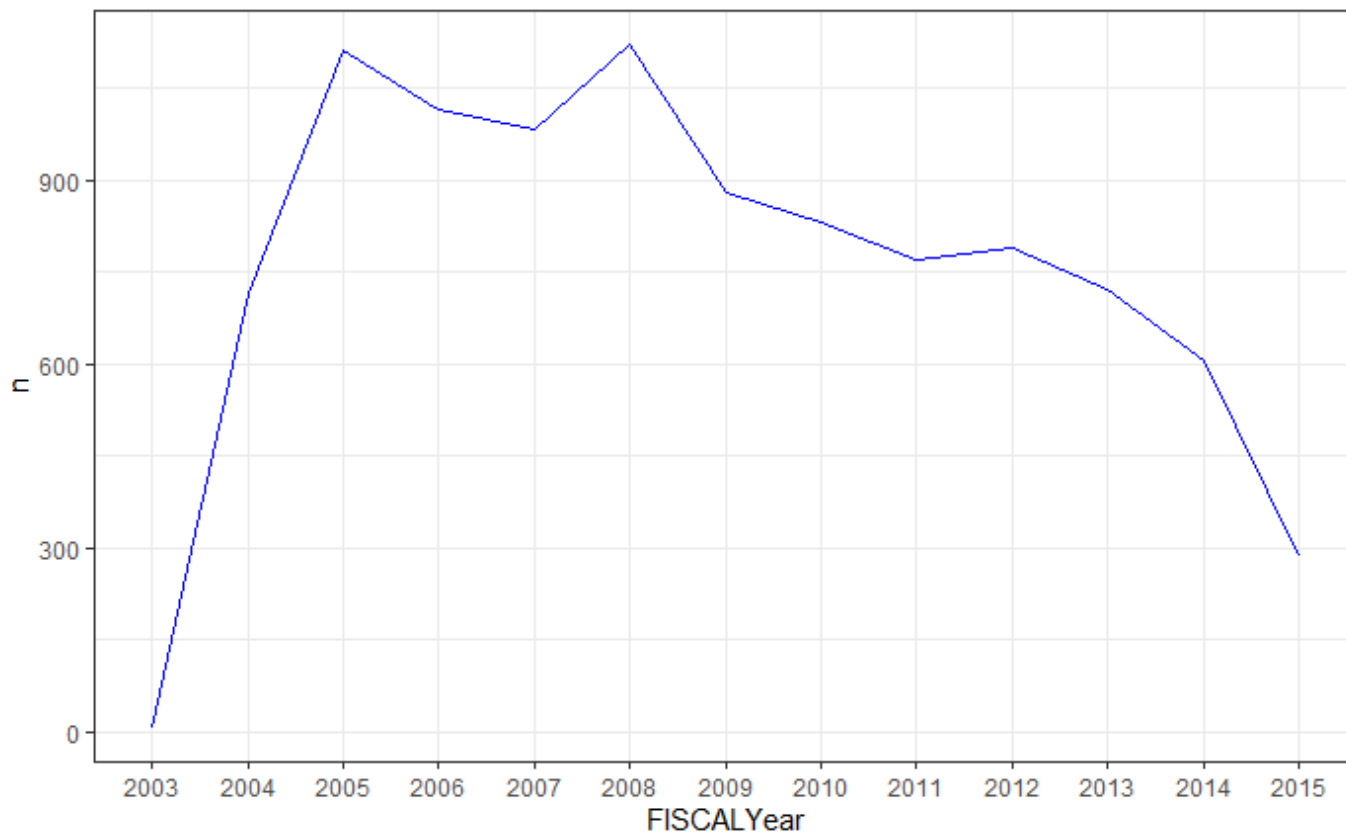
Hide

```
#Mining method using which employees face accident and are permanently transferred or terminated
us_data %>% filter(UG_MINING_METHOD != "NO VALUE FOUND") %>%
  ggplot(aes(UG_MINING_METHOD,fill= TRANS_TERM)) +
  geom_histogram(stat="count",position = "dodge") + coord_flip()+
  theme_bw()
```



Hide

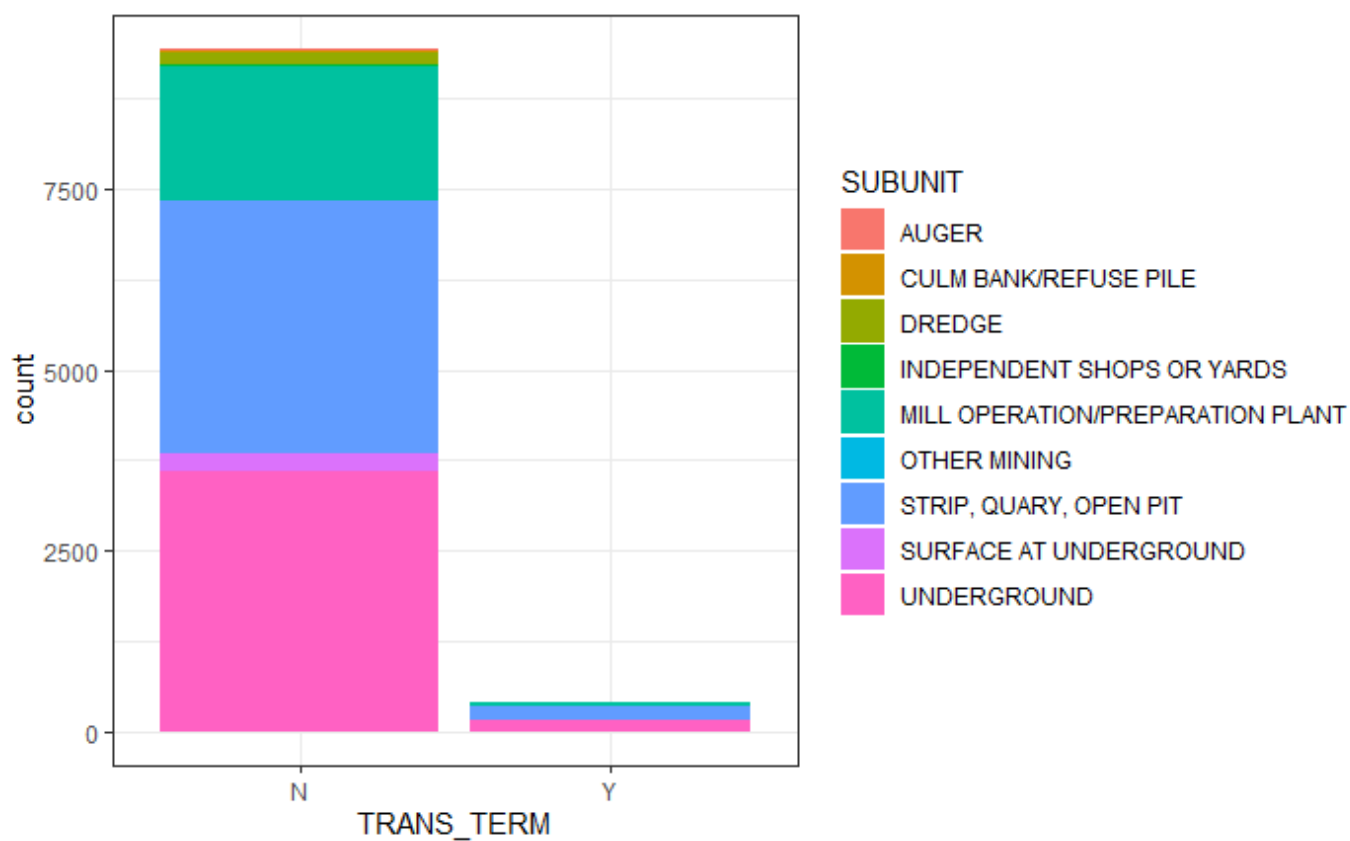
```
#number of accidents/illness in each year
us_data %>% group_by(FISCAL_YR) %>% count() %>%
  ggplot(aes(as.factor(FISCAL_YR),n)) + geom_line(col= "blue",group=1) +
  theme_bw() + xlab("FISCALYear")
```



2005 and 2008 have more cases of accidents and illness as compare to other years.

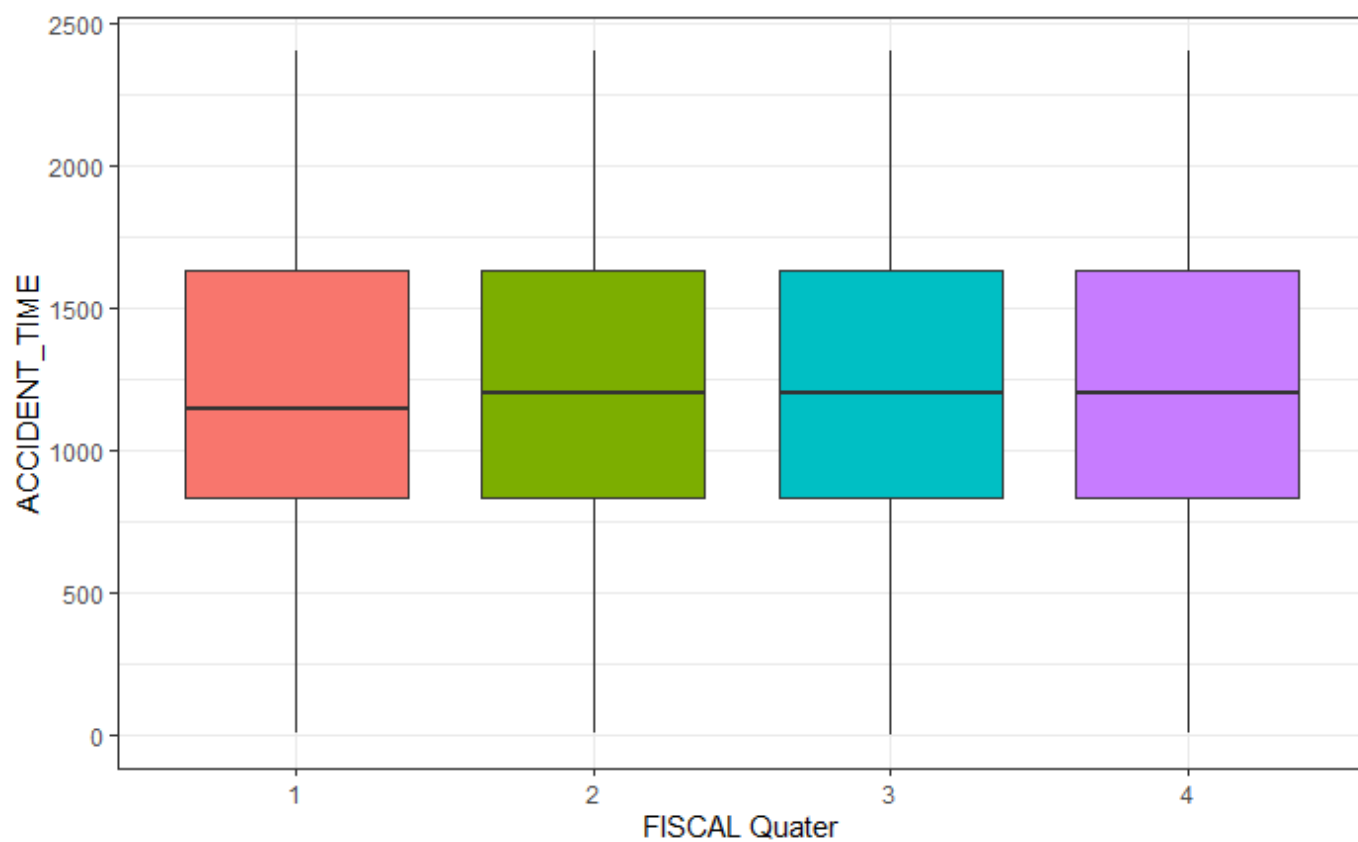
Hide

```
#number of employees transferred/terminated, working in sub units who faced accident
s
ggplot(data = us_data,aes(TRANS_TERM,fill=SUBUNIT)) +geom_bar()+
  theme_bw()
```



Hide

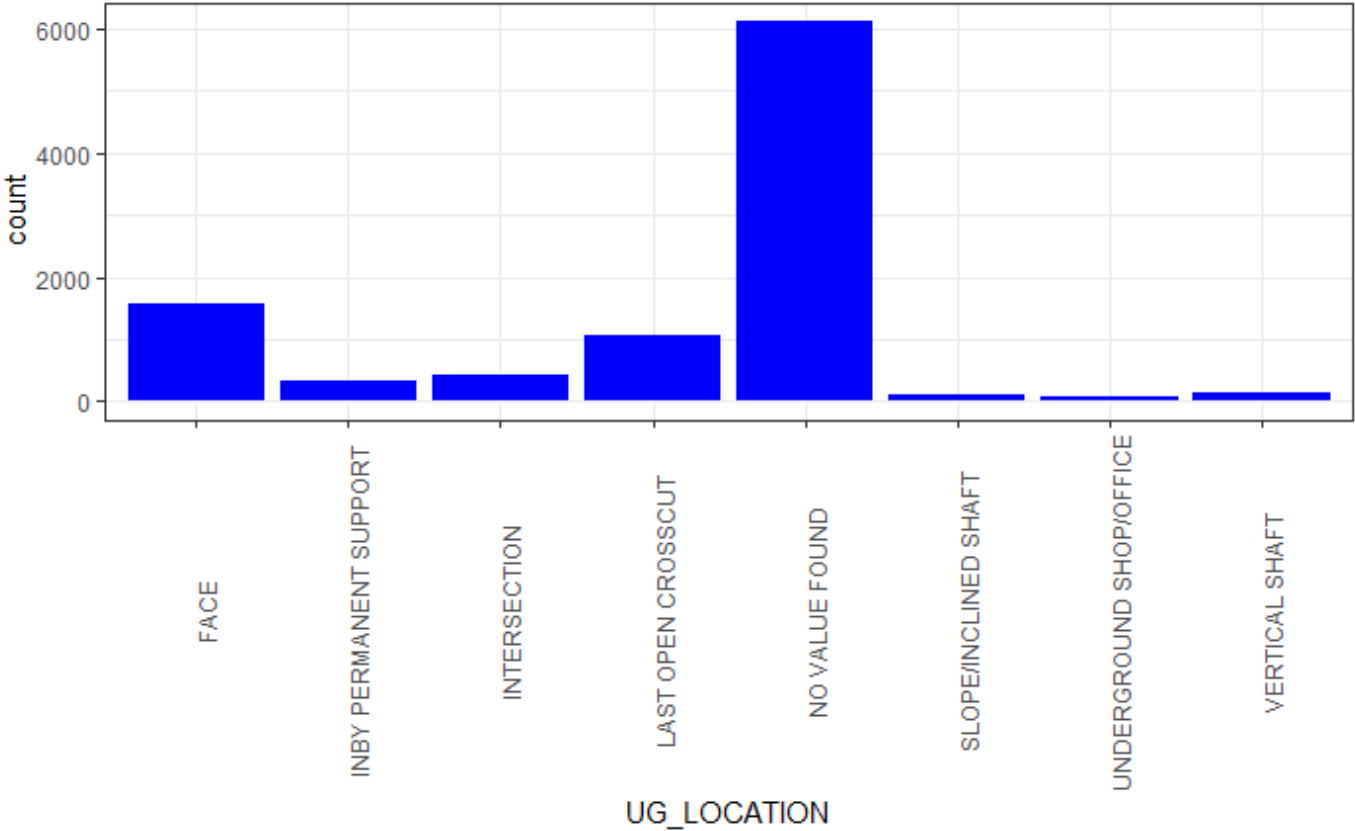
```
ggplot(data = us_data ,aes(as.factor(FISCAL_QTR),ACCIDENT_TIME,fill=as.factor(FISCAL_QTR))) +  
  geom_boxplot() + theme_bw() + theme( legend.position = "none") + xlab("FISCAL Quater  
r")
```



Hide

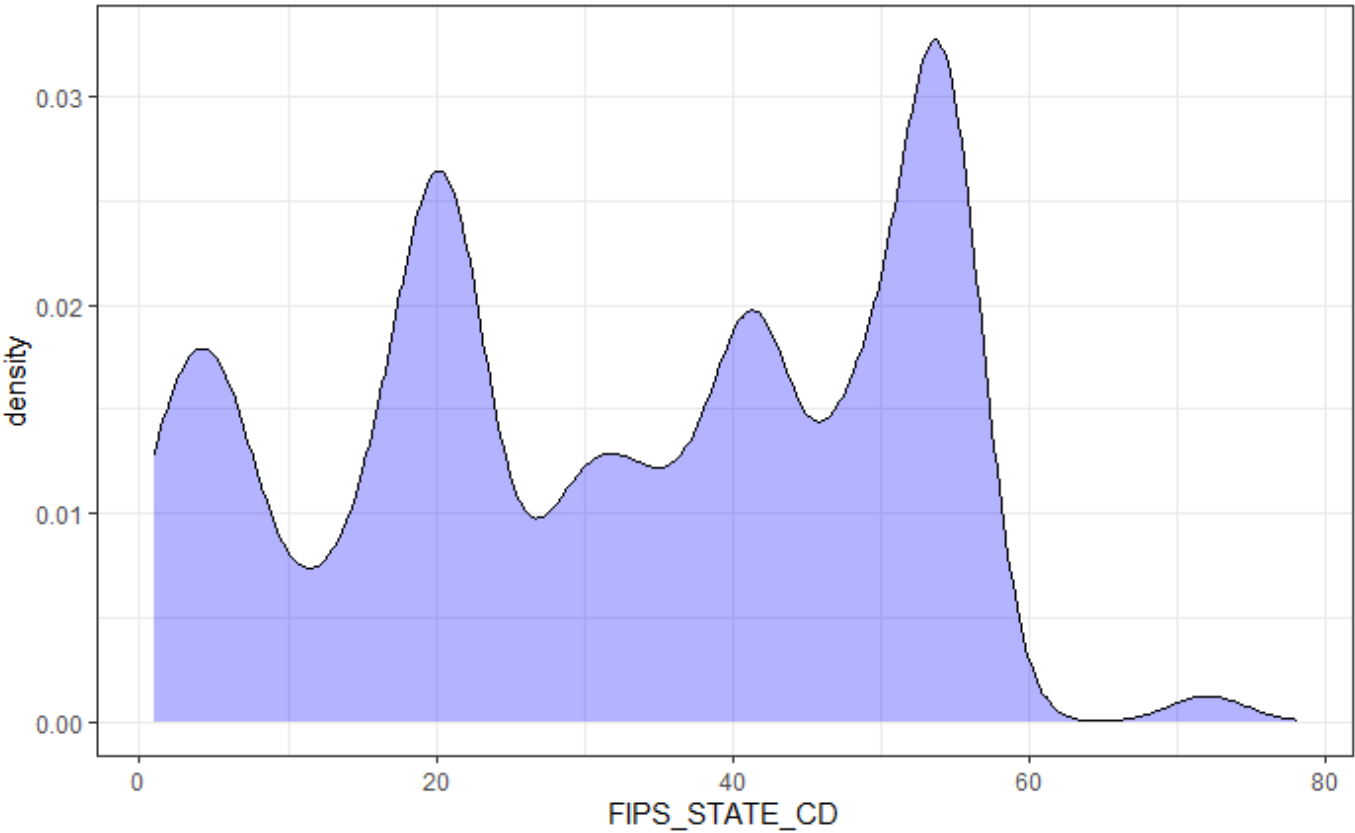
```
ggplot(data = us_data,aes(UG_LOCATION)) + geom_histogram(fill="blue",stat="count") +  
  theme_bw() +  
  theme(axis.text.x = element_text(angle=90))
```





Hide

```
ggplot(us_data,aes(FIPS_STATE_CD)) + geom_density(fill= "blue",alpha=0.3) +theme_bw()
```



Hide

```
ggplot(us_data,aes(TRANS_TERM,fill=TRANS_TERM)) + geom_bar()+  
  scale_fill_manual(values=c("#56B4E9", "#E69F00")) +  
  theme_bw() + theme( legend.position = "none")
```

