# Lab 1 VHE

#### Victoria Eastman

September 19, 2018

#### Initial EDA

```
Problem statement:
```

```
# Import libraries
suppressPackageStartupMessages(library(dplyr))
suppressPackageStartupMessages(library(Hmisc))
setwd("/home/victoriaeastman/berkeley/w271/w271 lab1")
data <- read.csv("challenger.csv")</pre>
glimpse(data)
## Observations: 23
## Variables: 5
## $ Flight <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16...
## $ Temp
           <int> 66, 70, 69, 68, 67, 72, 73, 70, 57, 63, 70, 78, 67, 5...
## $ Pressure <int> 50, 50, 50, 50, 50, 50, 100, 100, 200, 200, 200, 200,...
## $ 0.ring <int> 0, 1, 0, 0, 0, 0, 0, 1, 1, 1, 0, 0, 2, 0, 0, 0, ...
## $ Number
           describe(data)
## data
##
  5 Variables
                  23 Observations
## Flight
##
       n missing distinct Info Mean Gmd
                                                   .05
                                                            .10
##
       23
          0 23
                           1
                                    12
                                            8
                                                    2.1
                                                            3.2
##
      .25
             .50
                     .75
                            .90
                                    .95
                  17.5
##
      6.5
            12.0
                            20.8
                                    21.9
##
## lowest : 1 2 3 4 5, highest: 19 20 21 22 23
## Temp
##
       n missing distinct
                           Info
                                  Mean
                                           Gmd
                                                   .05
                                                           .10
##
       23 0 16 0.992
                                   69.57
                                        7.968
                                                   57.1
                                                           59.0
##
      .25
             .50
                    .75
                            .90
                                    .95
     67.0 70.0 75.0
                           77.6
##
                                    78.9
##
## Value
                 57
              53
                        58
                             63
                                  66
                                       67
                                            68
                                                      70
                                                            72
                                                 69
## Frequency
             1
                             1
                                  1
                                       3
                  1
                        1
                                             1
                                                  1
## Proportion 0.043 0.043 0.043 0.043 0.043 0.130 0.043 0.043 0.174 0.043
##
## Value
              73
                   75
                        76
                             78
                                  79
## Frequency
           1 2 2 1
                                  1
                                       1
## Proportion 0.043 0.087 0.087 0.043 0.043 0.043
```

```
## Pressure
##
             missing distinct
                                                         Gmd
                                    Info
                                              Mean
          n
                                                       67.59
##
         23
                    0
                                   0.706
                                             152.2
##
## Value
                  50
                       100
                              200
                         2
## Frequency
                   6
                               15
## Proportion 0.261 0.087 0.652
## O.ring
##
             missing distinct
                                    Info
                                              Mean
                                                         Gmd
##
         23
                    0
                              3
                                   0.654
                                            0.3913
                                                     0.6087
##
                   0
                                2
## Value
                          1
## Frequency
                  16
                          5
## Proportion 0.696 0.217 0.087
##
## Number
             missing distinct
##
                                    Info
                                              Mean
                                                         Gmd
##
         23
                                                           0
                    0
                                       0
                                                 6
                              1
##
## Value
                6
## Frequency
## Proportion 1
# I'm curious about the value counts for o-ring failures
table(data$0.ring)
##
```

Initial findings:

5

## 0 1 2

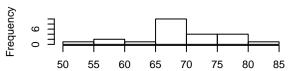
## 16

- 23 data points with no missing values for any variables
  - Dependent variable, O.ring, is categorical and takes three values: 0, 1, and 2 representing the number of o-ring failures on space launches. The mean value is 0.3913 which means the data is skewed towards 0 o-ring failures. Futher investigation shows there were 2 flights with 2 o-ring failures, 5 with 1 failure, and 16 with no failures.
  - The explanatory variables are as follows:
    - Temp: temperature at launch (degrees F)
    - Pressure: Combustion pression (psi)

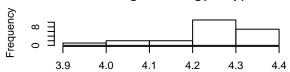
```
par(mfrow = c(3,2))

# histogram of explanatory variables
for (i in 2:4){
   hist(as.numeric(data[,i]), main=paste0("Histogram of ", colnames(data)[i]), xlab=NA)
   hist(as.numeric(log(data[,i])), main=paste0("Histogram of log(", colnames(data)[i], ")"), xlab=NA)
}
```

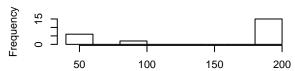
#### **Histogram of Temp**



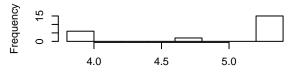
#### Histogram of log(Temp)



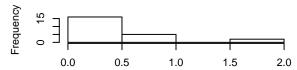
#### **Histogram of Pressure**



## Histogram of log(Pressure)



#### **Histogram of O.ring**



### Histogram of log(O.ring)

