# Statistics 135 – Lab Project

April 29, 2015

## 1 Background

```r
library("ggplot2")
library("grid")
library("gridExtra")

data.df <- na.omit(read.csv("diabetes.csv"))
```
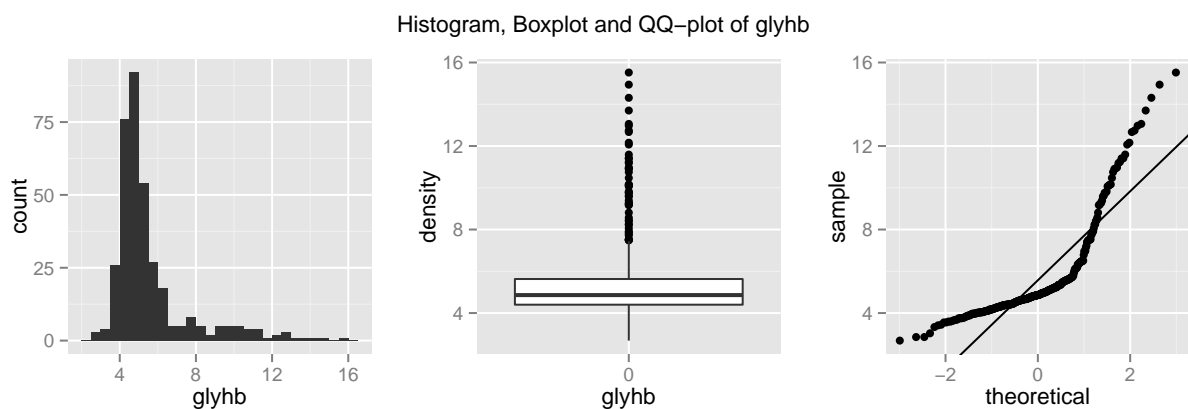
## 2 Accessing Data, Visualization and Summarization

1.
```r
glyhb.histogram <- ggplot(data.df) +
  geom_histogram(aes(x=glyhb), binwidth=0.5)

glyhb.boxplot <- ggplot(data.df) +
  geom_boxplot(aes(x=factor(0), y=glyhb)) +
  labs(x="glyhb", y="density")

glyhb.qqplot <- ggplot(data.df) +
  stat_qq(aes(sample=glyhb)) +
  geom_abline(aes(intercept=mean(glyhb), slope=sd(glyhb)))

grid.arrange(glyhb.histogram, glyhb.boxplot, glyhb.qqplot, ncol=3,
             main="Histogram, Boxplot and QQ-plot of glyhb")
```
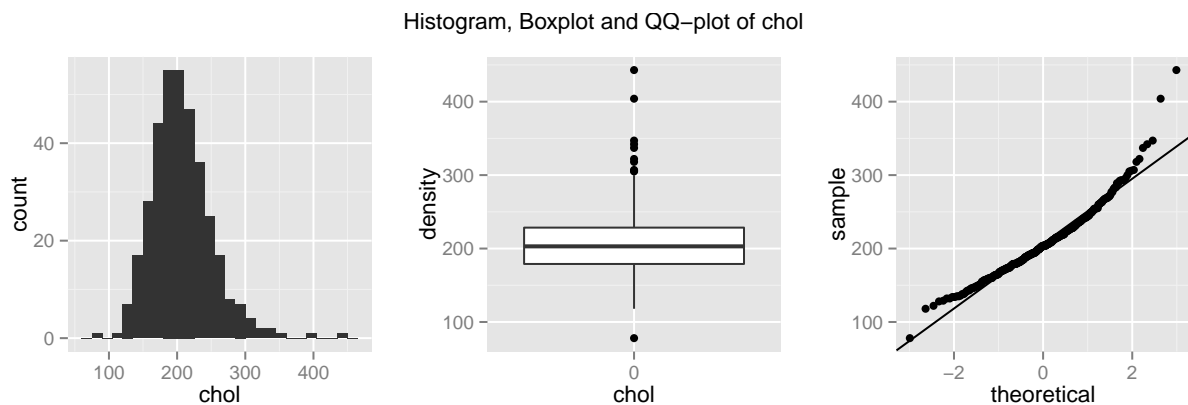


2.
```r
chol.histogram <- ggplot(data.df) +
  geom_histogram(aes(x=chol), binwidth=15)

chol.boxplot <- ggplot(data.df) +
```

```
  geom_boxplot(aes(x=factor(0), y=chol)) +
  labs(x="chol", y="density")

chol.qqplot <- ggplot(data.df) +
  stat_qq(aes(sample=chol)) +
  geom_abline(aes(intercept=mean(chol), slope=sd(chol)))

grid.arrange(chol.histogram, chol.boxplot, chol.qqplot, ncol=3,
             main="Histogram, Boxplot and QQ-plot of chol")
```
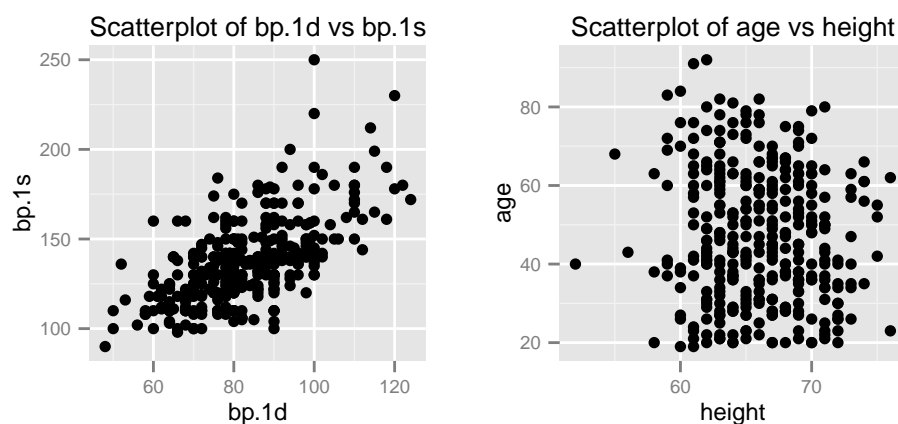
Histogram, Boxplot and QQ–plot of chol



3.
```
bp.1d.bp.1s.scatterplot <- ggplot(data.df) +
  geom_point(aes(x=bp.1d, y=bp.1s)) +
  labs(title="Scatterplot of bp.1d vs bp.1s") +
  theme(text=element_text(size=8.5))

height.age.scatterplot <- ggplot(data.df) +
  geom_point(aes(x=height, y=age)) +
  labs(title="Scatterplot of age vs height") +
  theme(text=element_text(size=8.5))

grid.arrange(bp.1d.bp.1s.scatterplot, height.age.scatterplot, ncol=2)
```
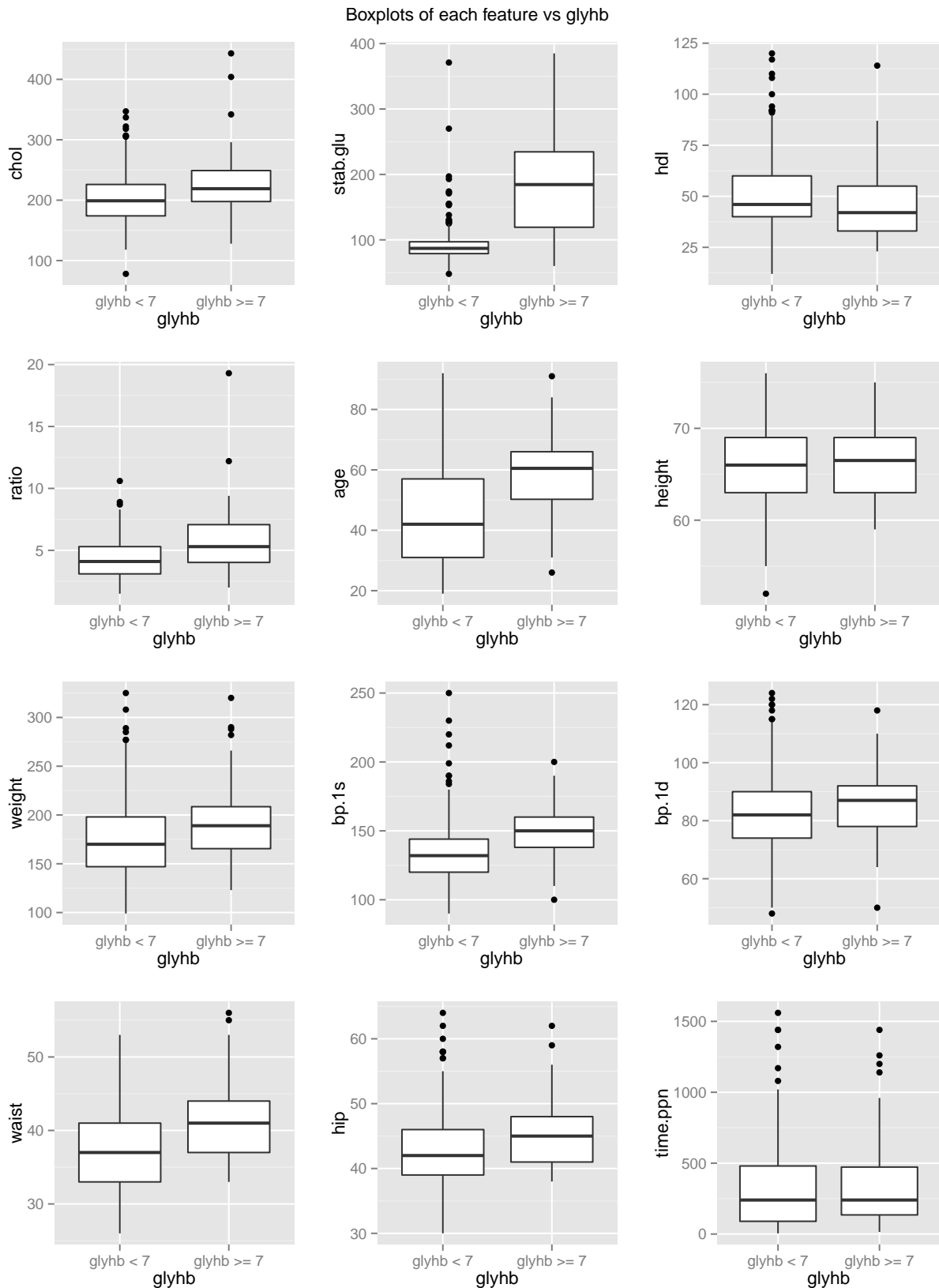


4.
```
data.df$glyhb.cond <- NA
data.df[data.df$glyhb>=7, ]$glyhb.cond <- "glyhb >= 7"
data.df[data.df$glyhb<7, ]$glyhb.cond <- "glyhb < 7"
```

```r
features <- c("chol", "stab.glu", "hdl", "ratio", "age", "height", "weight",
              "bp.1s", "bp.1d", "waist", "hip", "time.ppn")
feature.boxplots = list()
for (i in 1:length(features)) {
  feature <- features[i]
  feature.boxplot <- ggplot(data.frame(glyhb.cond=data.df$glyhb.cond,
                                        feature=data.df[, feature])) +
    geom_boxplot(aes(x=glyhb.cond, y=feature)) +
    labs(x="glyhb", y=feature)
  feature.boxplots[[i]] <- feature.boxplot
}

do.call(grid.arrange, c(feature.boxplots, ncol=3,
                        main="Boxplots of each feature vs glyhb"))
```
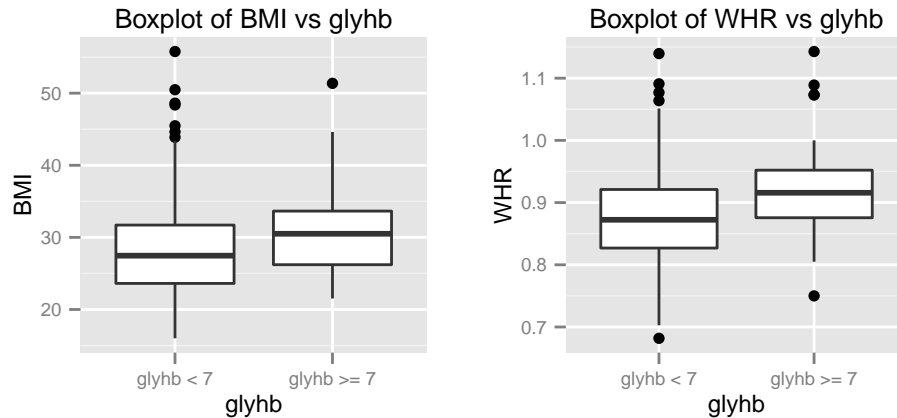
Boxplots of each feature vs glyhb

```
5. data.df$BMI <- 703*data.df$weight/data.df$height^2
   data.df$WHR <- data.df$waist/data.df$hip

   features <- c("BMI", "WHR")
   feature.boxplots = list()
   for (i in 1:length(features)) {
     feature <- features[i]
```

```
    feature.boxplot <- ggplot(data.frame(glyhb.cond=data.df$glyhb.cond,
                                    feature=data.df[, feature])) +
      geom_boxplot(aes(x=glyhb.cond, y=feature)) +
      labs(title=paste(c("Boxplot of", feature, "vs glyhb"), collapse=" "),
            x="glyhb", y=feature) +
      theme(text=element_text(size=8.5))
    feature.boxplots[[i]] <- feature.boxplot
}

do.call(grid.arrange, c(feature.boxplots, ncol=2))
```



6.

# 3  Parametric Inferece

1.

$$X \sim Gamma(\alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x}$$

$$E(X) = \frac{\alpha}{\beta}$$

$$E(X^2) = Var(X) + [E(X)]^2$$
$$= \frac{\alpha}{\beta^2} + \left(\frac{\alpha}{\beta}\right)^2$$
$$= \frac{\alpha(\alpha+1)}{\beta^2}$$

$$\begin{cases} E(X) = \frac{\alpha}{\beta} \\ E(X^2) = \frac{\alpha(\alpha+1)}{\beta^2} \end{cases} \implies \begin{cases} \alpha = \frac{[E(X)]^2}{Var(x)} \\ \beta = \frac{E(X)}{Var(x)} \end{cases} \implies \begin{cases} \hat{\alpha}_{MOM} = \frac{\overline{X}_n}{\frac{1}{n}\sum_{i=1}^n (X_i - \overline{X}_n)^2} \\ \hat{\beta}_{MOM} = \frac{\overline{X}_n^2}{\frac{1}{n}\sum_{i=1}^n (X_i - \overline{X}_n)^2} \end{cases}$$

```
alpha.mom <- mean(data.df$BMI)^2/var(data.df$BMI)
beta.mom <- mean(data.df$BMI)/var(data.df$BMI)

BMI.bootstrap <- sapply(1:1000, function(x) {
  BMI.sample <- sample(data.df$BMI, nrow(data.df), replace=TRUE)
  alpha.sample <- mean(BMI.sample)^2/var(BMI.sample)
  beta.sample <- mean(BMI.sample)/var(BMI.sample)
```

```
  return(c(alpha.sample, beta.sample))
})

CI.BMI <- sapply(1:2, function(x) {
  CI <- quantile(BMI.bootstrap[x, ], probs=c(0.025, 0.975))
  return(CI)
})
colnames(CI.BMI) <- c("alpha", "beta")

################################################################################
CI.BMI

##          alpha      beta
## 2.5%   15.91447 0.5477096
## 97.5% 21.79548 0.7626395


################################################################################
```
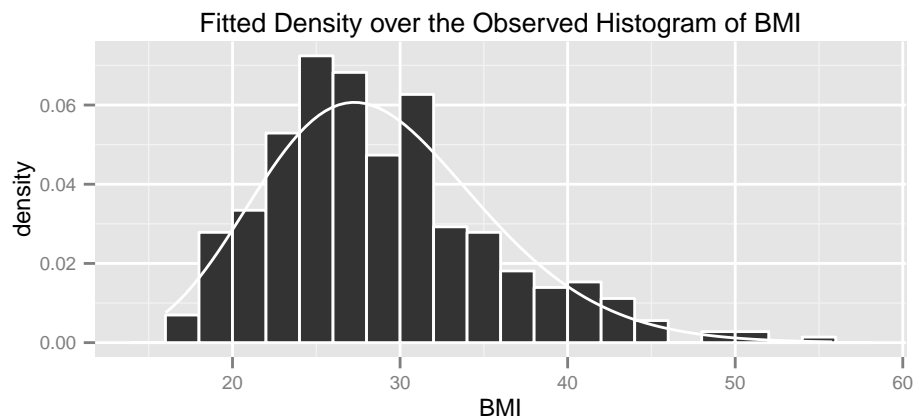
```
ggplot(data.df) +
  geom_histogram(aes(x=BMI, y=..density..), binwidth=2, col="white") +
  stat_function(fun=function(x)
    dgamma(x, shape=alpha.mom, rate=beta.mom), col="white") +
  labs(title="Fitted Density over the Observed Histogram of BMI") +
  theme(text=element_text(size=8.5))
```



2.

$$\hat{\mu}_{MLE} = \overline{X}_n$$

$$\hat{\sigma}^2_{MLE} = \frac{1}{n} \sum_{i=1}^{n} (X_i - \overline{X}_n)^2$$

```
mu.mle <- mean(data.df$WHR)
sigma.mle <- sqrt(var(data.df$WHR))

WHR.bootstrap <- sapply(1:1000, function(x) {
  WHR.sample <- sample(data.df$WHR, nrow(data.df), replace=TRUE)
  mu.sample <- mean(WHR.sample)
  sigma.sample <- sqrt(var(WHR.sample))
```

```
    return(c(mu.sample, sigma.sample))
})

CI.WHR <- sapply(1:2, function(x) {
  CI <- quantile(WHR.bootstrap[x, ], probs=c(0.025, 0.975))
  return(CI)
})
colnames(CI.WHR) <- c("mu", "sigma")

################################################################################
CI.WHR

##               mu        sigma
## 2.5%   0.8746646 0.06791852
## 97.5% 0.8896728 0.07939450


################################################################################
```
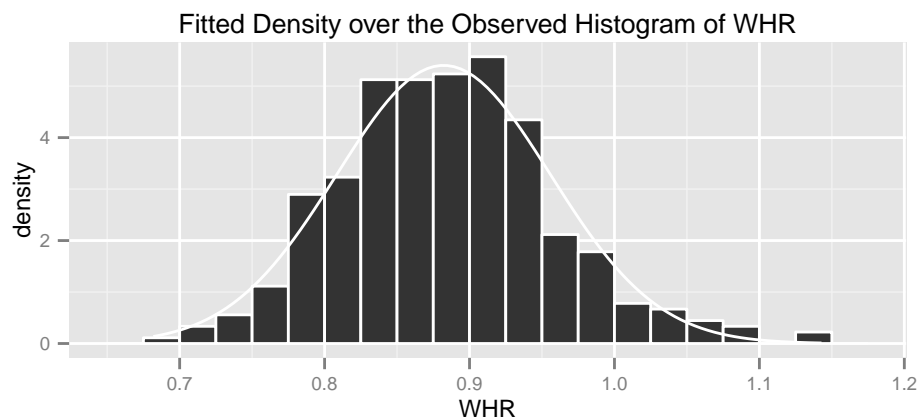
```
ggplot(data.df) +
  geom_histogram(aes(x=WHR, y=..density..), binwidth=0.025, col="white") +
  stat_function(fun=function(x)
    dnorm(x, mean=mu.mle, sd=sigma.mle), col="white") +
  labs(title="Fitted Density over the Observed Histogram of WHR") +
  theme(text=element_text(size=8.5))
```



3.

# 4 Testing

1.

2.

3.

4.

5.

6.

# 5 Regression

1.

2.

3.

4.

5.

6.

7.