

ACM 157 Set 1

1. x_1, \dots, x_n \bar{x} \tilde{x} s_x IQR_x $y_i = \alpha + \beta x_i$

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^n (\alpha + \beta x_i) = \frac{1}{n} \cdot n\alpha + \frac{1}{n} \cdot \beta \sum_{i=1}^n x_i = \alpha + \frac{\beta}{n} \sum_{i=1}^n x_i$$

$$\tilde{y} = \begin{cases} y_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ \frac{1}{2}(y_{(\frac{n}{2})} + y_{(\frac{n}{2}+1)}) & \text{if } n \text{ even} \end{cases}$$

$$= \begin{cases} \alpha + \beta x_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ \frac{1}{2}(\alpha + \beta x_{(\frac{n}{2})} + \alpha + \beta x_{(\frac{n}{2}+1)}) & \text{if } n \text{ even} \end{cases}$$

$$= \begin{cases} \alpha + \beta x_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ \alpha + \beta \frac{1}{2}(x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}) & \text{if } n \text{ even} \end{cases}$$

Since $\tilde{x} = \begin{cases} x_{(\frac{n+1}{2})} & \text{if } n \text{ odd} \\ \frac{1}{2}(x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}) & \text{if } n \text{ even} \end{cases}$

$$= \alpha + \beta \tilde{x}$$

$$s_y = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (\alpha + \beta x_i - \alpha - \beta \tilde{x})^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n \beta^2 (x_i - \tilde{x})^2}$$

$$= \sqrt{\frac{\beta^2}{n} \sum_{i=1}^n (x_i - \tilde{x})^2} = |\beta| s_x$$

We could also see that $\tilde{y} = \alpha + \beta \tilde{x}$, because we are just linearly scaling x_i 's to get y_i 's, so the median stays the same data point. Similarly we can see this will be the case for Q_{1y} and Q_{3y} , except if β is negative as the order of the data points switch. Therefore we get:

$$IQR_y = Q_{3y} - Q_{1y}$$

$$= \begin{cases} \alpha + \beta Q_{3x} - (\alpha + \beta Q_{1x}) & \text{if } \beta \geq 0 \\ \alpha + \beta Q_{1x} - (\alpha + \beta Q_{3x}) & \text{if } \beta < 0 \end{cases}$$

$$= \begin{cases} \beta(Q_{3x} - Q_{1x}) & \text{if } \beta \geq 0 \\ \beta(Q_{1x} - Q_{3x}) & \text{if } \beta < 0 \end{cases}$$

$$= \begin{cases} \beta IQR_x & \text{if } \beta \geq 0 \\ -\beta IQR_x & \text{if } \beta < 0 \end{cases}$$

$$= |\beta| IQR_x$$

2. $\bar{x} = \arg\min_{\alpha} \sum_{i=1}^n (x_i - \alpha)^2$ $\tilde{x} = \arg\min_{\alpha} \sum_{i=1}^n |x_i - \alpha|$

$\sum_{i=1}^n (x_i - \alpha)^2$ minimized when $\frac{d}{d\alpha} \sum_{i=1}^n (x_i - \alpha)^2 = 0$ & $\frac{d^2}{d\alpha^2} \sum_{i=1}^n (x_i - \alpha)^2 > 0$

$$\frac{d}{d\alpha} \sum_{i=1}^n (x_i - \alpha)^2 = 0$$

$$\sum_{i=1}^n 2(x_i - \alpha) \cdot (-1) = 0$$

$$\frac{d^2}{d\alpha^2} \sum_{i=1}^n (x_i - \alpha)^2 > 0$$

$$\frac{d}{d\alpha} (2 \sum_{i=1}^n \alpha - 2 \sum_{i=1}^n x_i) > 0$$

$$2 \sum_{i=1}^n \alpha - 2 \sum_{i=1}^n x_i = 0$$

$$2 \sum_{i=1}^n 1 > 0$$

$$2n > 0$$

$$n\alpha = \sum_{i=1}^n x_i$$

$$\alpha = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

From these, we see that α is a minimum and is equal to \bar{x}

To show $\tilde{x} = \arg\min_{\alpha} \sum_{i=1}^n |x_i - \alpha|$, we will show that the median minimizes $\sum_{i=1}^n |x_i - \alpha|$ when n is odd and the median is a minimizer of $\sum_{i=1}^n |x_i - \alpha|$ when n is even.

When n is odd, we can see what happens to $\sum_{i=1}^n |x_i - \alpha|$ when we move it away from \tilde{x} . We know there are an equal amount of points ($\frac{n-1}{2}$) above and below \tilde{x} since \tilde{x} is on the middle point when n is odd. If we move α in either direction from \tilde{x} then we can see that $|x_i - \alpha|$ changes by the same amount for each i , increasing or decreasing depending if α moves away from or towards x_i , respectively. However, if we move α away from \tilde{x} in either direction then we can see that it is moving towards fewer points than it is moving away from (initially $\frac{n-1}{2}$ towards and $\frac{n-1}{2} + 1$ away and gets worse as we move past more points), so we are increasing $|x_i - \alpha|$ for more points than we are decreasing it. Therefore we can see that \tilde{x} minimizes $\sum_{i=1}^n |x_i - \alpha|$ for odd n .

When n is even, we use the same logic as above to show that $\alpha = \tilde{x}$ is in the range of values that minimize $\sum_{i=1}^n |x_i - \alpha|$. From the previous case, we know that if we are moving α towards fewer points than we are moving it away from, then $\sum_{i=1}^n |x_i - \alpha|$ increases. Unlike the previous case where $\alpha = \tilde{x}$ is the only spot with an equal number of points on either side of α , since n is even, the entire range of values between the two center points of the data set will have an equal number of points on either side. Since \tilde{x} is the average of the two center points, we know it will be in the range of minimum values, $\tilde{x} \in [x_{(\frac{n}{2})}, x_{(\frac{n}{2}+1)}]$, so is a minimizer of $\sum_{i=1}^n |x_i - \alpha|$ for even n .

$$3. x_1, \dots, x_n \quad \{(z_{n+1}^u, x_{(k)})\}$$

If the collection of points falls on the line $y=ax+b$, then we know it still is a normal distribution since departures from normality are indicated by departures from a straight line. We can also see this is the case since $y=ax+b$ is just scaling and shifting $y=x$. This tells us that the distribution is $N(b, a^2)$ since the mean will be shifted by b and the standard deviation scaled by a .

$$6. P = \{1, \dots, N\} \quad S = \{s_1, \dots, s_n\}$$

$$a) P(s_1=N), \dots, P(s_n=N)$$

$$P(s_1=N) = \frac{1}{N} \quad \text{since all values of } P \text{ are unique}$$

$$P(s_2=N) = \frac{1}{N-1} \cdot \frac{N-1}{N} = \frac{1}{N} \quad \text{probability it gets selected given it}$$

$$P(s_3=N) = \frac{1}{N-2} \cdot \frac{N-2}{N-1} \cdot \frac{N-1}{N} = \frac{1}{N} \quad \text{wasn't already picked}$$

$$\text{For all } i, P(s_i=N) = \frac{1}{N}$$

$$b) P(N^{\text{th}} \text{ unit is in sample})$$

We can see it is the sum of probabilities from (a)

$$P(s_1=N) + P(s_2=N) + \dots + P(s_n=N)$$

$$= \frac{1}{N} + \frac{1}{N} + \dots + \frac{1}{N}$$

$$= \frac{n}{N}$$

$$c) E[s_i]$$

$$E[s_i] = \mu = \frac{1}{N} \sum_{i=1}^N x_i = \frac{1}{N} (1+2+\dots+N) = \frac{1}{N} \cdot \frac{1}{2}(N+1)N$$

$$= \frac{1}{2}(N+1)$$

$$d) P(s_1=N, s_2=1)$$

$$= P(s_1=N) \cdot P(s_2=1 | s_1=N)$$

$$= \frac{1}{N} \cdot \frac{1}{N-1} = \frac{1}{N(N-1)}$$

$$e) P(s_i=i, \text{ for all } i=1, \dots, n)$$

$$= P(s_1=1) \cdot P(s_2=2 | s_1=1) \cdot \dots \cdot P(s_n=n | s_1=1, s_2=2, \dots, s_{n-1}=n-1)$$

$$= \frac{1}{N} \cdot \frac{1}{N-1} \cdot \dots \cdot \frac{1}{N-n+1}$$

$$= \frac{(N-n)!}{N!}$$

$$7. \bar{X}_n^w = \sum_{i=1}^n w_i x_i \quad \mu = \frac{1}{N} \sum_{i=1}^N x_i$$

$$a) \bar{X}_n^w \text{ unbiased estimate of } \mu$$

$$E[\bar{X}_n^w] = \mu$$

$$E[\bar{X}_n^w] = E\left[\sum_{i=1}^n w_i x_i\right] = \sum_{i=1}^n w_i E[x_i] = \sum_{i=1}^n w_i \mu = \mu \sum_{i=1}^n w_i$$

$$\sum_{i=1}^n w_i = 1$$

b) Minimize $V[\bar{X}_n^w]$ w/ $\sum_{i=1}^n w_i = 1$

$$V[\bar{X}_n^w] = V[\sum_{i=1}^n w_i x_i]$$

$$\text{Since } V[\sum_{i=1}^n x_i] = \sum_{i=1}^n V[x_i] + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n \text{Cov}(x_i, x_j)$$

$$\begin{aligned} V[\sum_{i=1}^n w_i x_i] &= \sum_{i=1}^n V[w_i x_i] + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n \text{Cov}(w_i x_i, w_j x_j) \\ &= \sum_{i=1}^n w_i^2 V[x_i] + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n w_i w_j \text{Cov}(x_i, x_j) \\ &= \sum_{i=1}^n w_i^2 \sigma^2 + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n w_i w_j \left(-\frac{\sigma^2}{N-1}\right) \\ &= \sum_{i=1}^n w_i^2 \sigma^2 - \frac{\sigma^2}{N-1} 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n w_i w_j \end{aligned}$$

$$\left(\sum_{i=1}^n w_i\right)^2 = \sum_{i=1}^n w_i^2 + 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n w_i w_j \quad \& \quad \sum_{i=1}^n w_i = 1$$

$$\Rightarrow 2 \sum_{i=1}^{n-1} \sum_{j=i+1}^n w_i w_j = 1 - \sum_{i=1}^n w_i^2$$

$$= \sum_{i=1}^n w_i^2 \sigma^2 - \frac{\sigma^2}{N-1} \left(1 - \sum_{i=1}^n w_i^2\right)$$

$$= \left(\sigma^2 + \frac{\sigma^2}{N-1}\right) \sum_{i=1}^n w_i^2 - \frac{\sigma^2}{N-1}$$

Since $(\sigma^2 + \frac{\sigma^2}{N-1}) > 0$ and $(\sigma^2 + \frac{\sigma^2}{N-1})$ and $(-\frac{\sigma^2}{N-1})$ are both constants, we can minimize:

$$\sum_{i=1}^n w_i^2 \quad \text{w/} \quad \sum_{i=1}^n w_i = 1$$

Using Lagrange multipliers:

$$\frac{\partial}{\partial w_j} \left(\sum_{i=1}^n w_i^2 + \lambda \left(\sum_{i=1}^n w_i - 1 \right) \right) = 0$$

$$2w_j + \lambda = 0 \quad w_j = -\frac{\lambda}{2} \quad \forall j$$

Since all the weights are equal to the same constant and $\sum_{i=1}^n w_i = 1$, then $w_i = \frac{1}{n} \quad \forall i$

We can see $\frac{1}{n}$ is the minimum when $\sum_{i=1}^n w_i = 1$ by looking at the Cauchy-Schwartz inequality:

$$\left(\sum_{i=1}^n u_i v_i\right)^2 \leq \left(\sum_{i=1}^n u_i^2\right) \left(\sum_{i=1}^n v_i^2\right)$$

$$\left(\sum_{i=1}^n w_i\right)^2 \leq \left(\sum_{i=1}^n w_i^2\right) \left(\sum_{i=1}^n 1^2\right)$$

$$1 \leq \left(\sum_{i=1}^n w_i^2\right) n$$

$$\sum_{i=1}^n w_i^2 \geq \frac{1}{n} \quad \text{we see that } w_i = \frac{1}{n} \text{ is the minimum} \quad \searrow$$

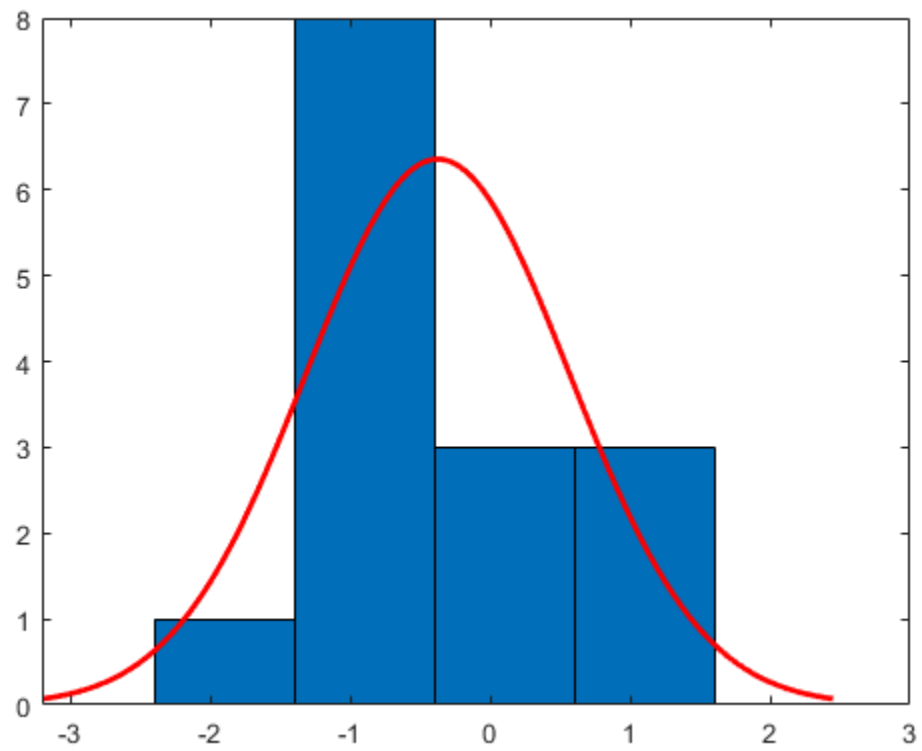
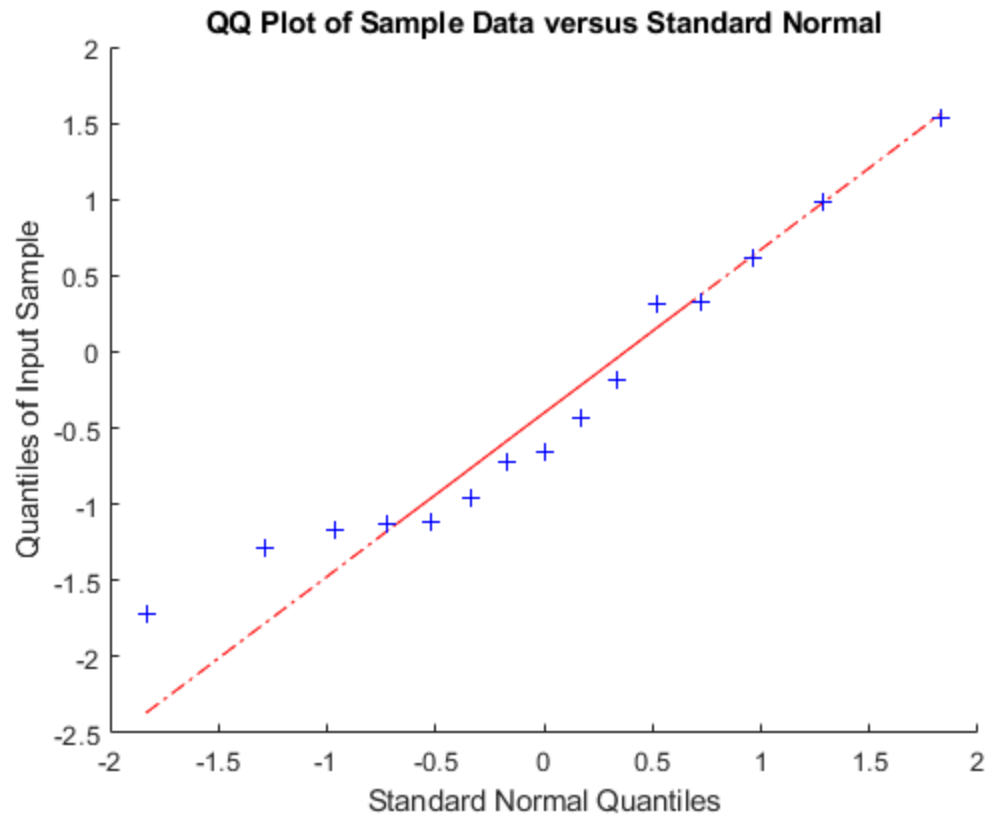
$$\sum_{i=1}^n \left(\frac{1}{n}\right)^2 = \frac{1}{n^2} \sum_{i=1}^n 1 = \frac{1}{n^2} n = \frac{1}{n} \geq \frac{1}{n}$$

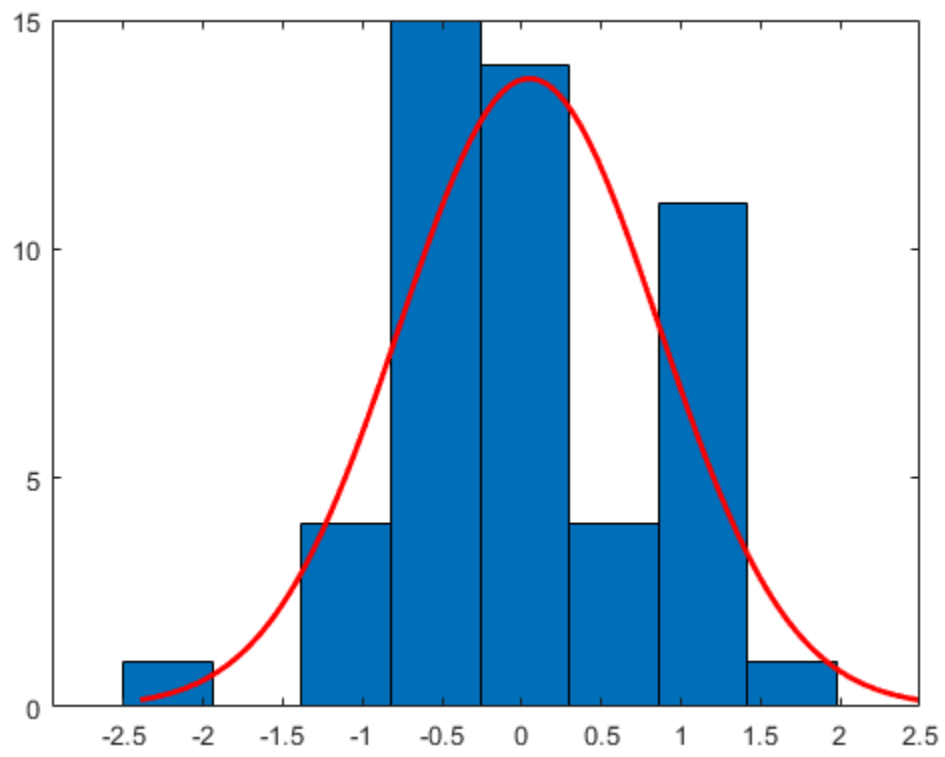
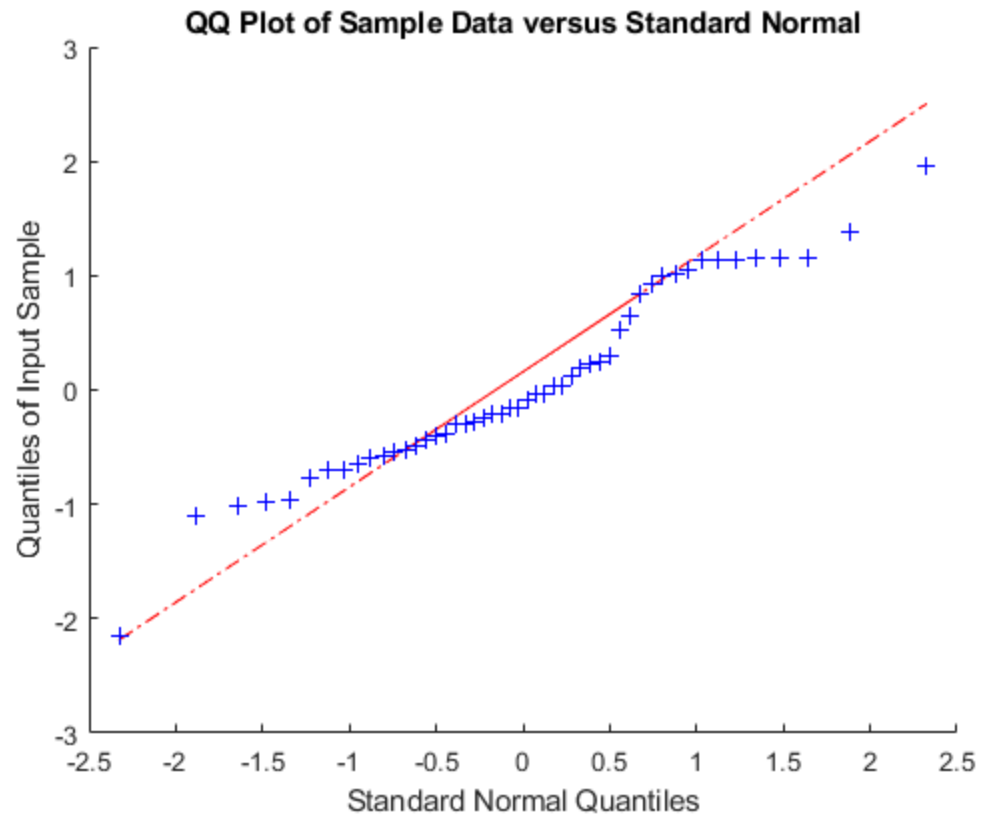
Most efficient estimate is when $w_i = \frac{1}{n}$
or \bar{X}_n

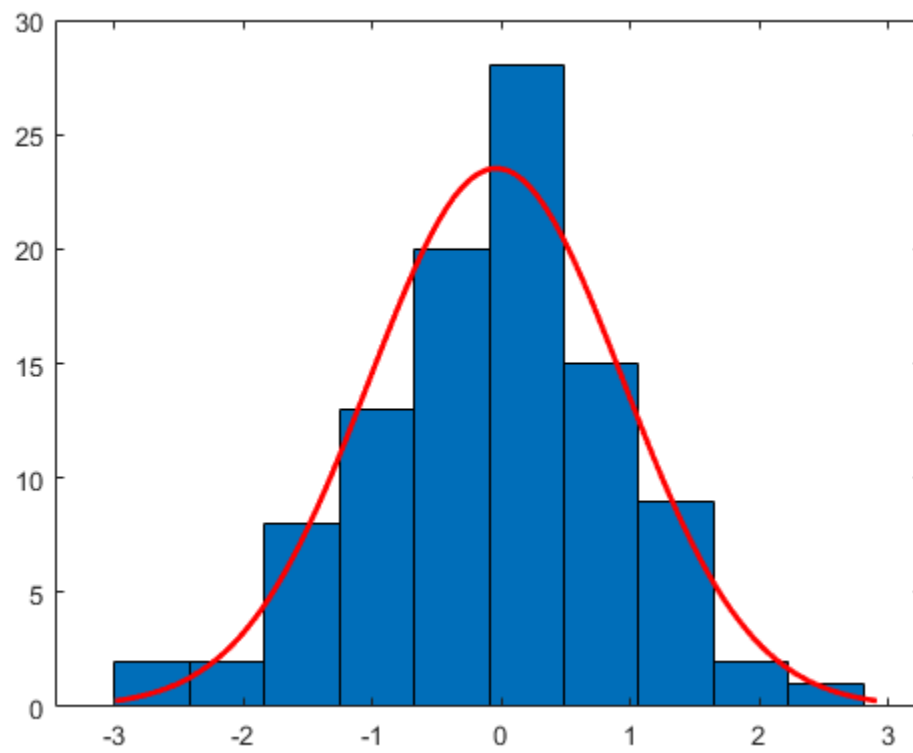
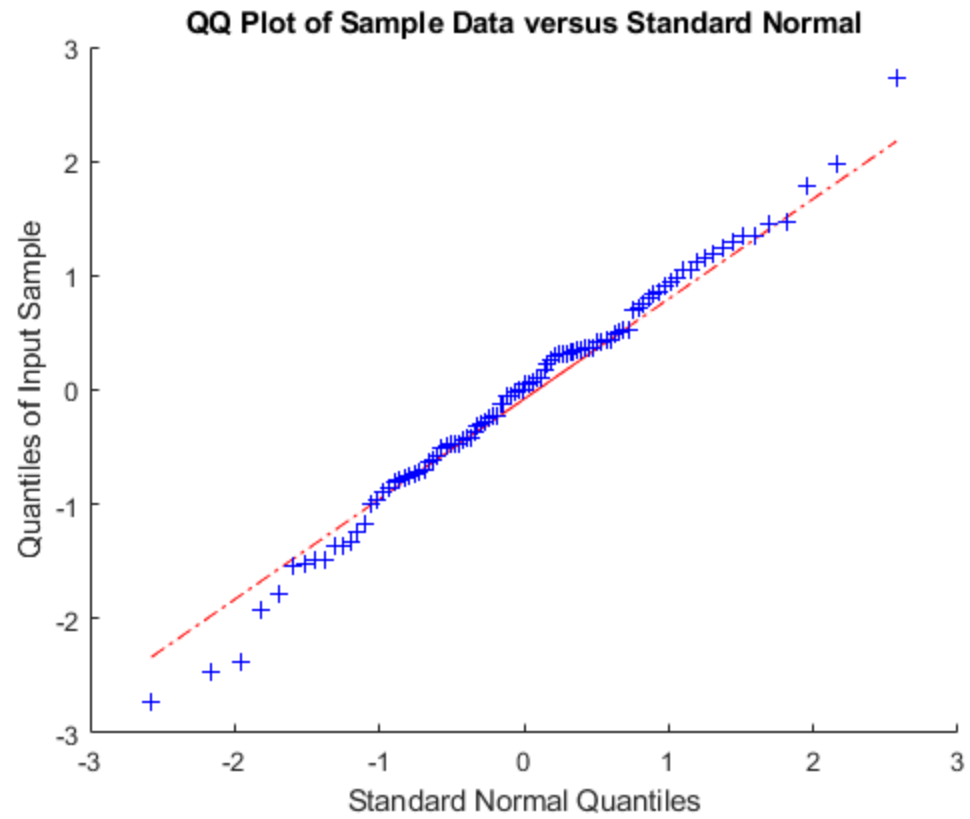
```
r = normrnd(0,1,[1 15]);
figure(1);
qqplot(r)
figure(2);
histfit(r)
% (a) The points on the QQ plot appear to fall around a straight line
% with
% some deviation. The histogram is unimodal but is not symmetric and
% is
% only very roughly bell-shaped. In general, the QQ plot is fairly
% close to
% the line but deviates further, the less symmetric and bell-shaped
% the
% histogram is.
r = normrnd(0,1,[1 50]);
figure(3);
qqplot(r)
figure(4);
histfit(r)

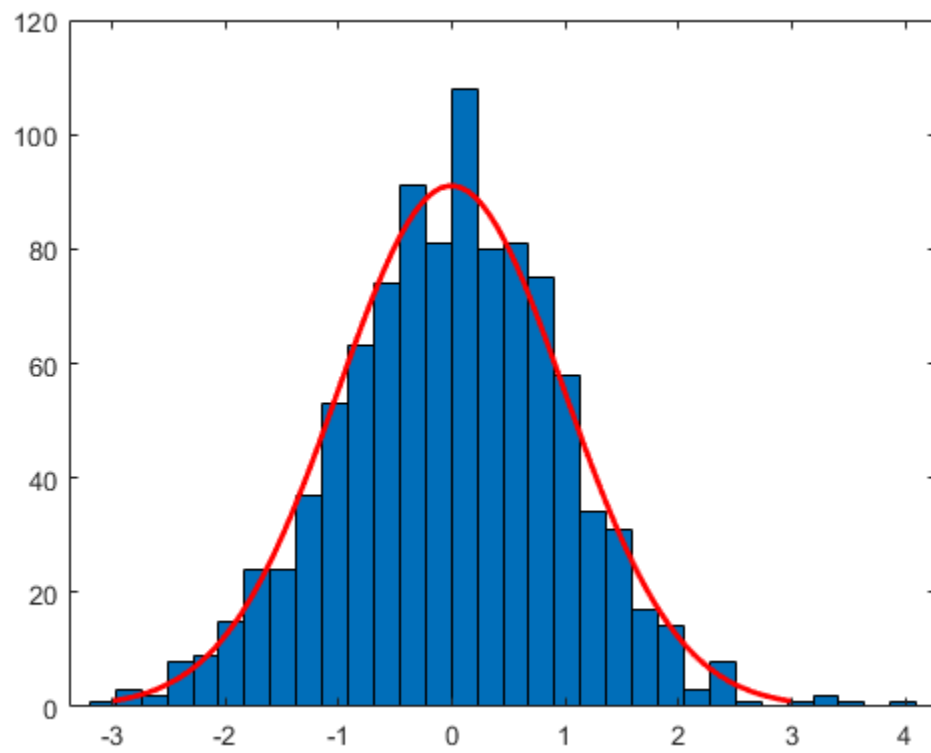
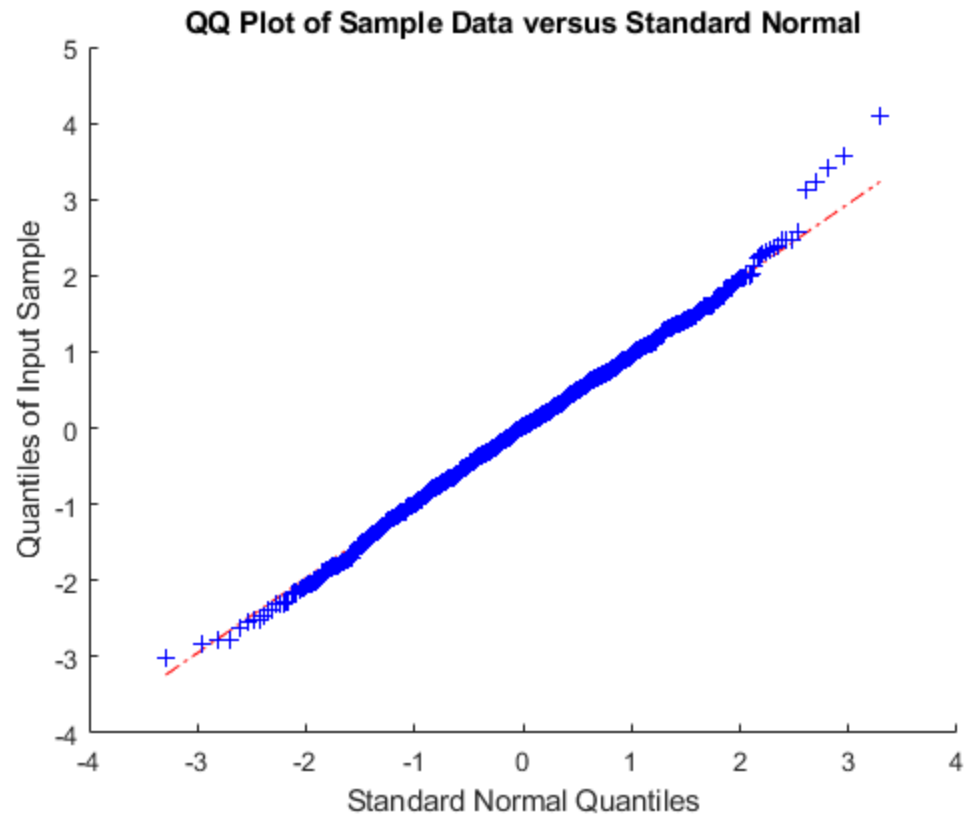
r = normrnd(0,1,[1 100]);
figure(5);
qqplot(r)
figure(6);
histfit(r)

r = normrnd(0,1,[1 1000]);
figure(7);
qqplot(r)
figure(8);
histfit(r)
% (b) Similar to (a), the points on the QQ plot appear to fall on the
% line
% with less deviation as we increase the number of points. Similarly,
% the
% histogram is more symmetric and bell-shaped with the more points
% used.
% (c) I would estimate the critical sample size to be around 1000 as
% at
% this number of points the plot doesn't deviate substantially from
% linearity.
```









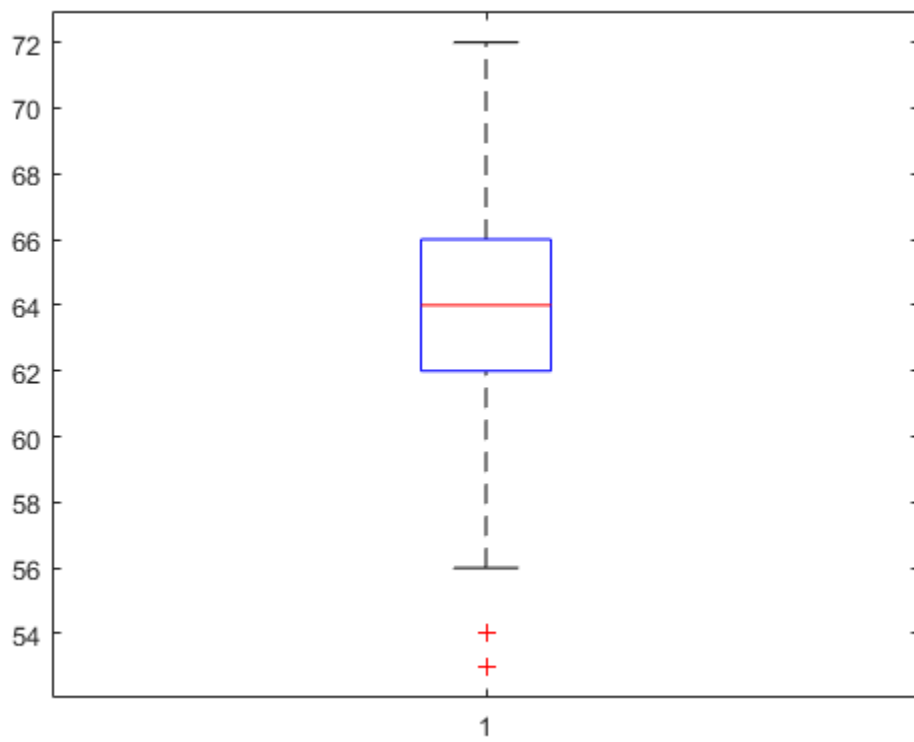
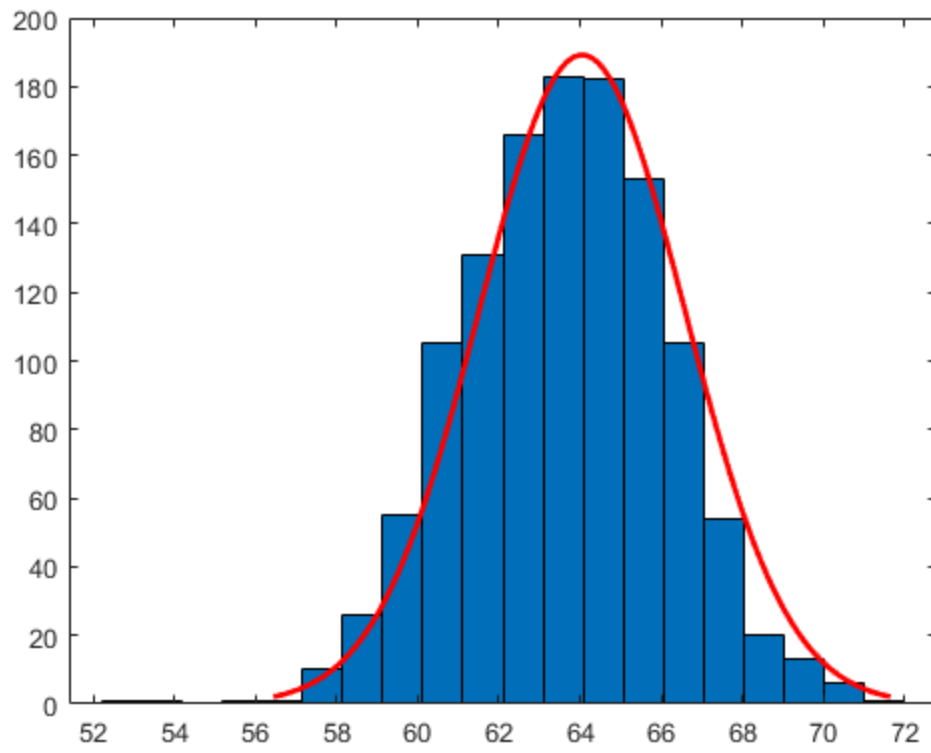
Published with MATLAB® R2021a

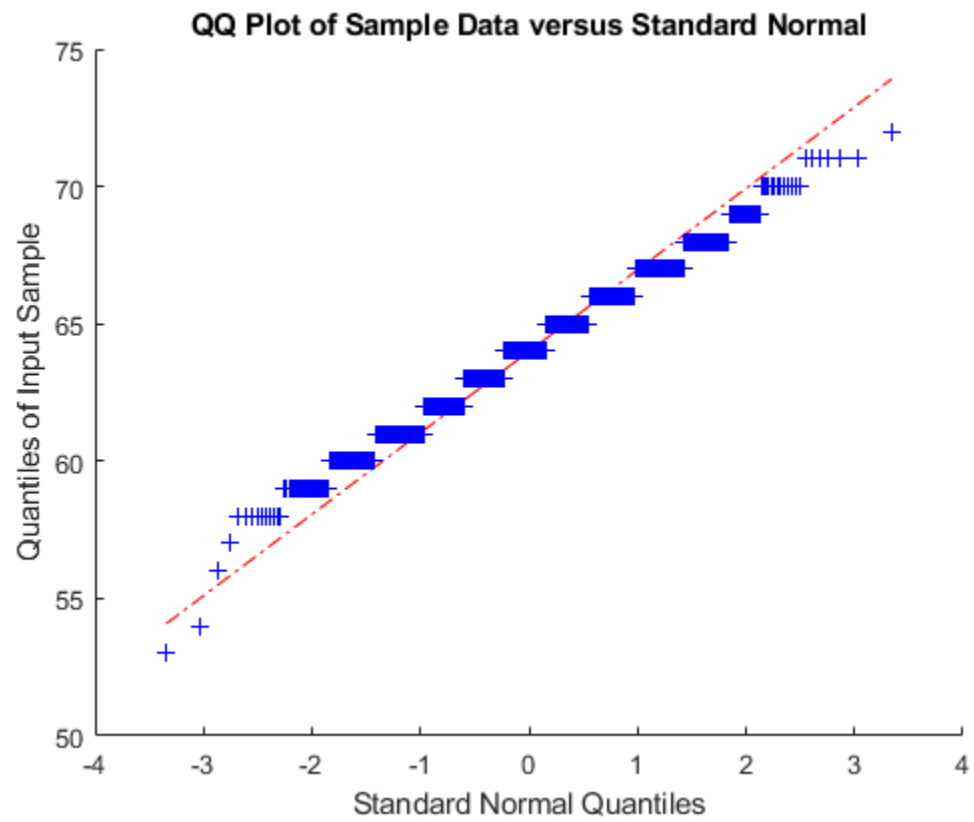
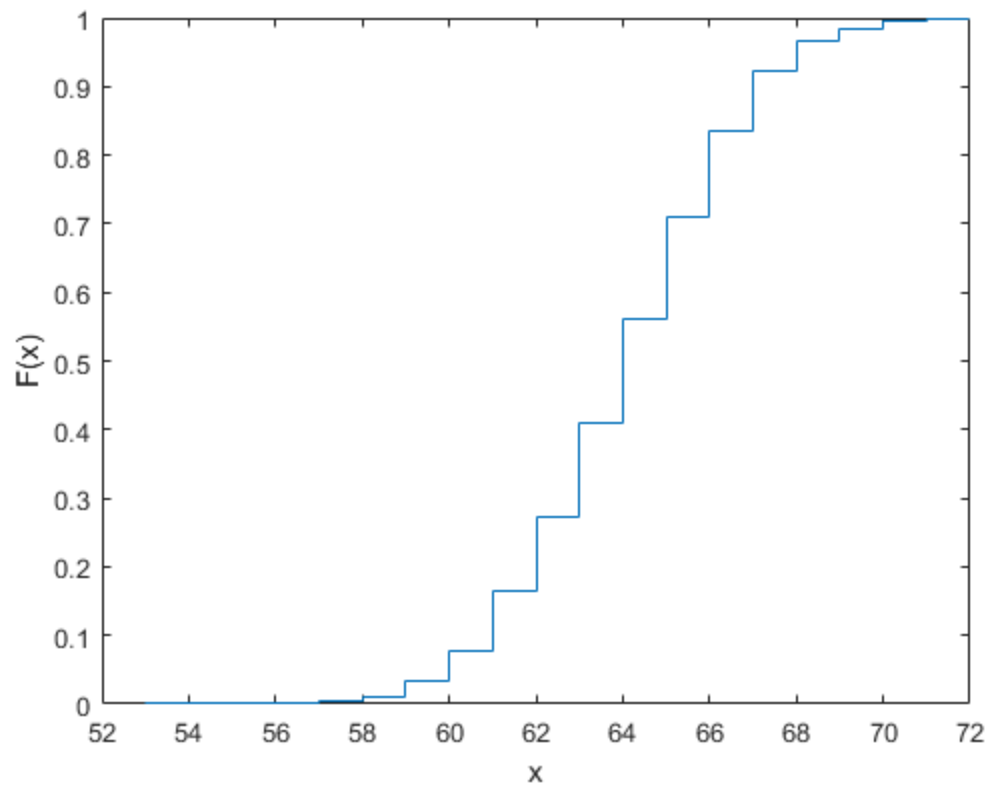
```

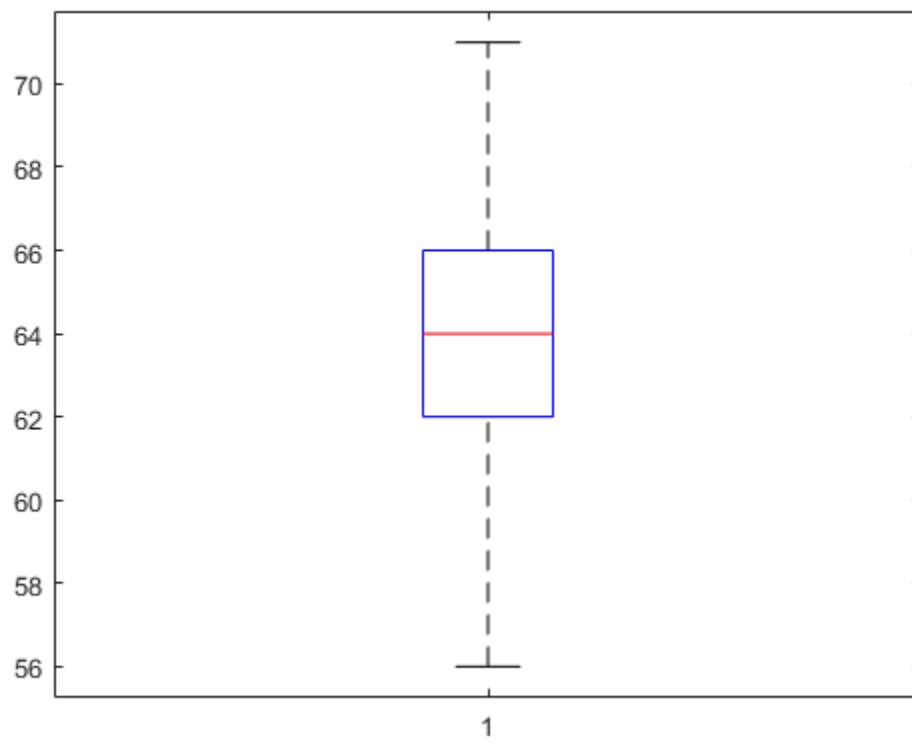
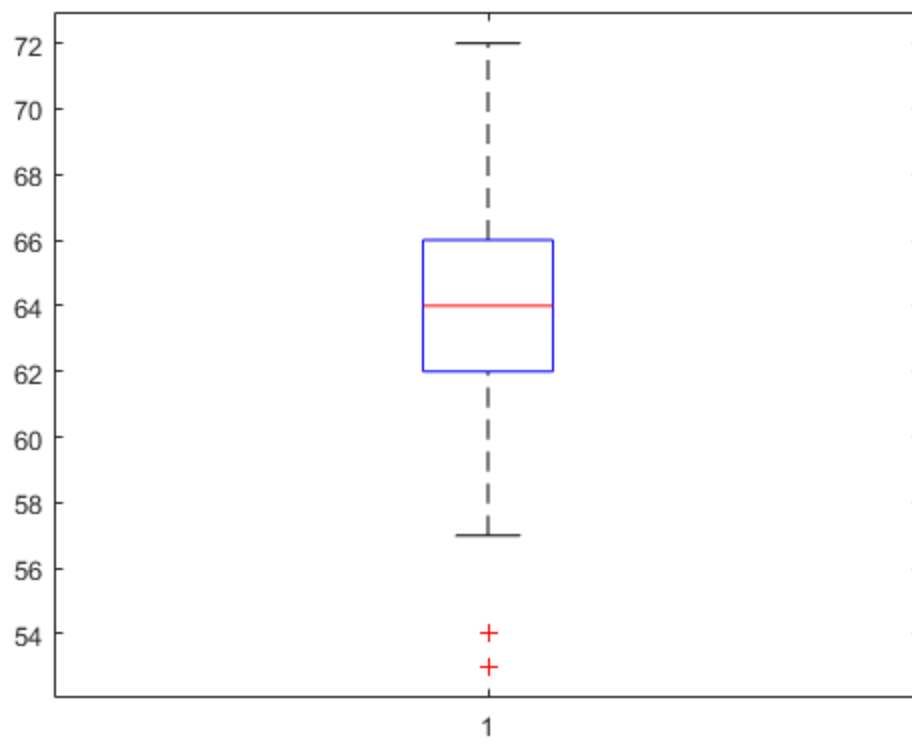
A = importdata("birth.txt");
A(A(:,5) == 99, :) = [];
H = A(:,5);
figure(1);
histfit(H,20)
% (a) 20 bins seems to be optimal for representing the shape of the
% distribution.
disp([mean(H) median(H) std(H) iqr(H)])
% (b) Since the mean and median are very close to each other, the
% center of
% the sample is well-defined.
figure(2);
boxplot(H)
figure(3);
ecdf(H)
figure(4);
qqplot(H)
% (c) I would consider the sample as being approximately normal with  $\mu$ 
% =
% 64.05 and  $\sigma^2 = 6.42$ 
S=A;
N=A;
S(S(:,7) == 9, :) = [];
S(S(:,7) == 0, :) = [];
figure(5);
boxplot(S(:,5))
N(N(:,7) == 9, :) = [];
N(N(:,7) == 1, :) = [];
figure(6);
boxplot(N(:,5))
% (d) On average, the two groups seem to be around the same height as
% the
% median and quartiles are at about the same value.

64.0478    64.0000    2.5334    4.0000

```







Published with MATLAB® R2021a