## 5.  The normal approximation to discrete distributions, and assessing normality

Readings:        Rosner: 5.6-5.9
                 OpenIntro Statistics: 3.1-3.2
                 Chihara and Hesterberg: 4.3


R:               stem, boxplot, qqnorm, qqline, shapiro.test, lillie.test, cvm.test, ad.test


SAS:             PROC UNIVARIATE


Homework:     Homework 2 due by noon on September 17
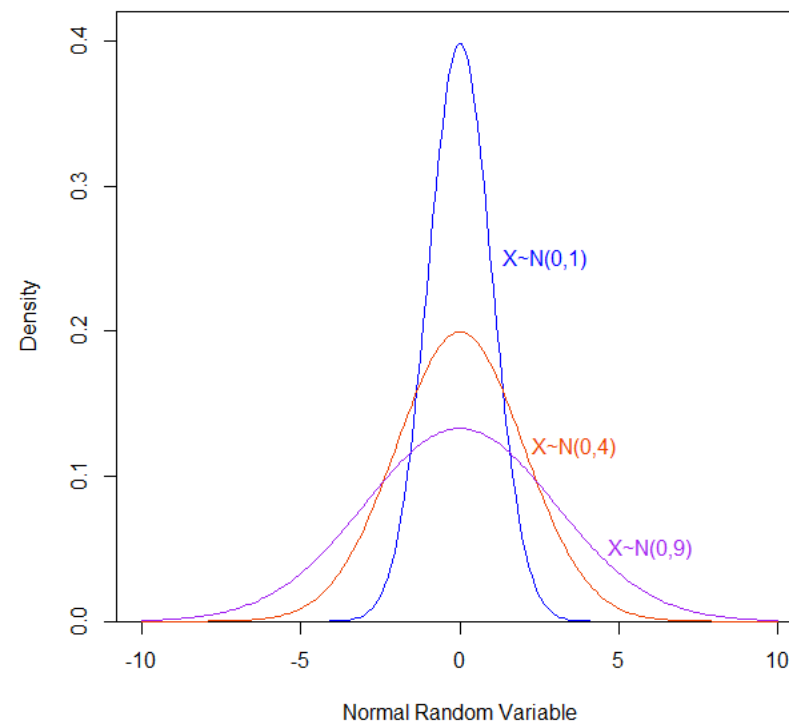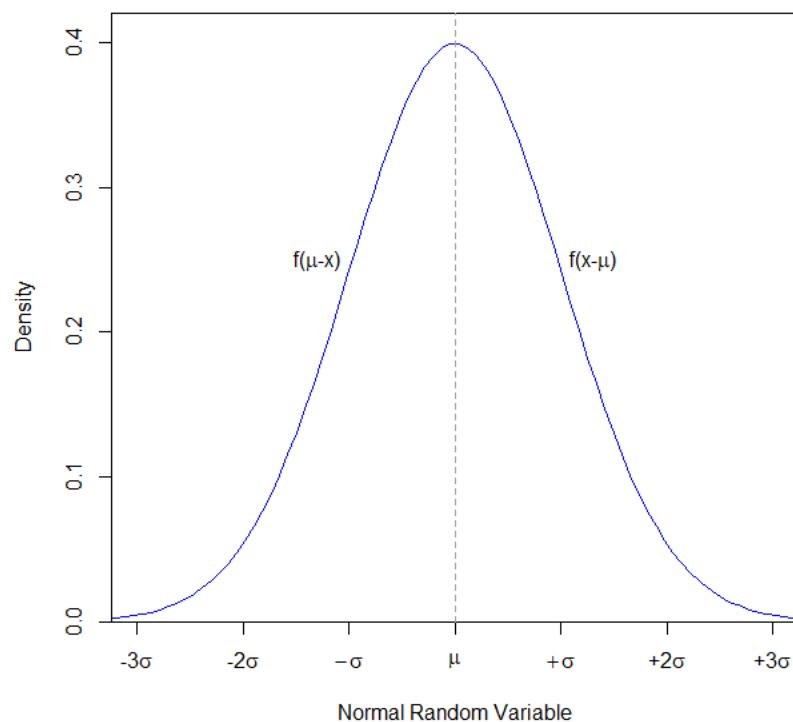                 Homework 3 due by noon on September 24

## Overview


A) Normal distribution and properties
B) Normal approximation to the binomial and Poisson
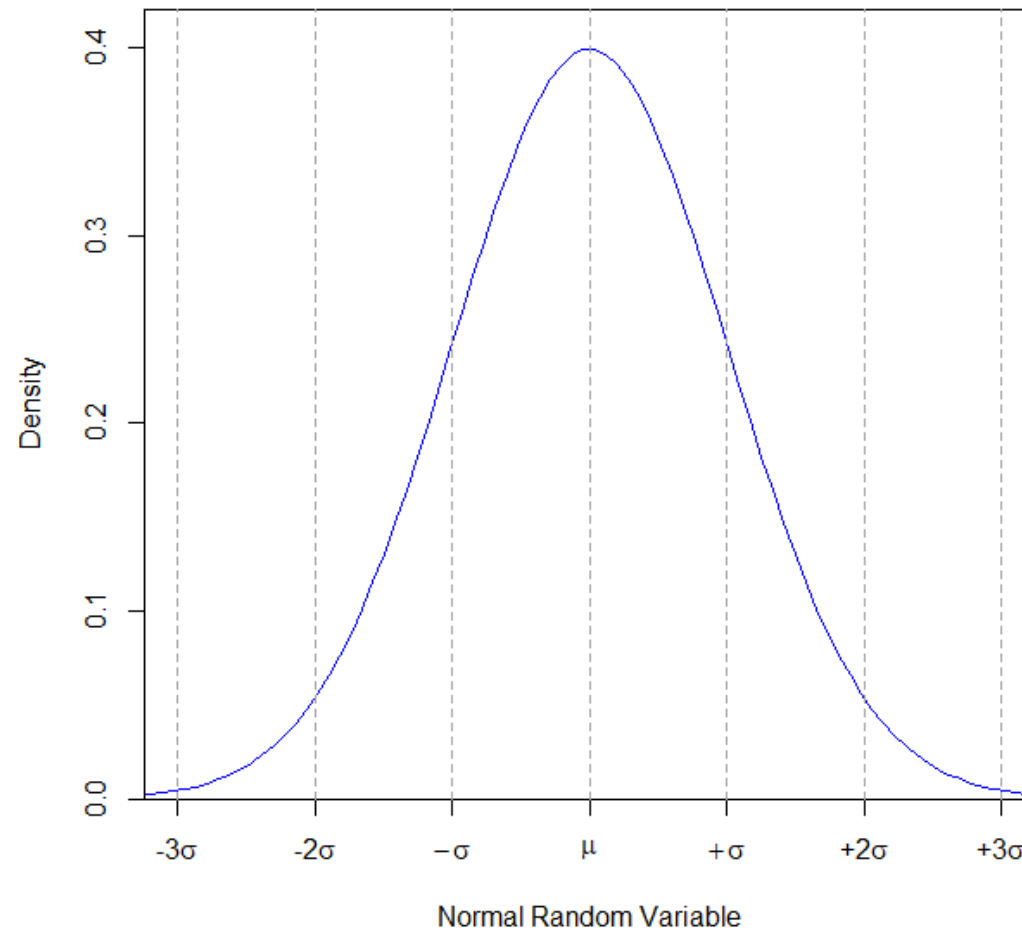C) Assessing normality

## A) Normal distribution properties

The normal distribution is defined by its *probability density function* (pdf), which is given as:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right), -\infty < x < \infty$$

for parameters $\mu$, $\sigma^2 > 0$. And the distribution is symmetric about $\mu$: $f(\mu+x) = f(\mu-x)$.

$P(\mu - \sigma < X < \mu + \sigma) = 0.6827 =$ (about 68% of area lies between $\pm$ 1 s.d.)

$P(\mu - 2\sigma < X < \mu + 2\sigma) = 0.9545$ (about 95% of area lies between $\pm$ 2 s.d.)

$P(\mu - 3\sigma < X < \mu + 3\sigma) = 0.9973$ (about 99.7% of area lies between $\pm$ 3 s.d.)

## Population Moments

The population moments can describe the location, variability, skewness, and kurtosis of a population, just as the corresponding sample moments describe a sample.

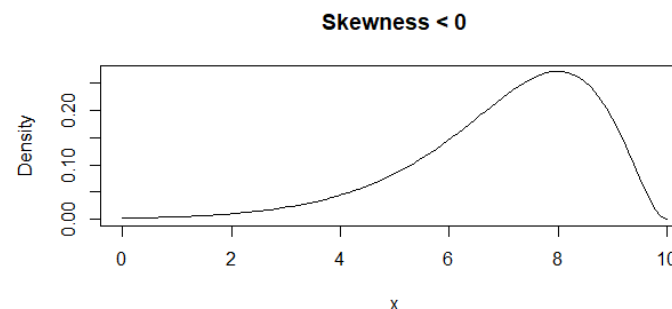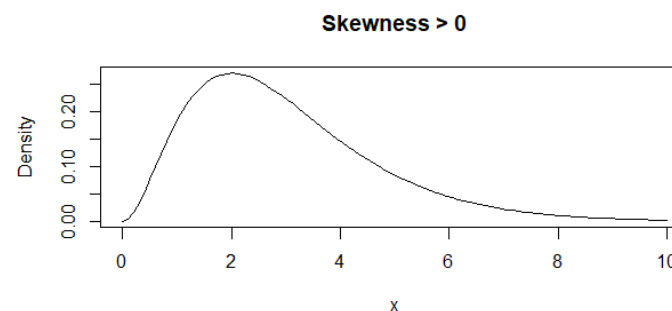$E(X) = \mu$, $E(X^2)$ are population moments about zero

$E[(X-\mu)] = 0$, $E[(X-\mu)^2] = \sigma^2$, are population moments about $\mu$ (these are called *central moments*)

**Skewness** describes symmetry of distribution. The 3rd central moment of the data, as with the 1st central moment, will balance out from left to right if the data are symmetric.  With normally distributed data, we expect skewness to be 0.

$(E[(X-\mu)^3])$ is standardized by $s^3 = \frac{1}{n}\sum_{i=1}^{n}(X_i - \bar{X})^3$

If skewness is > 0:   positive skew; skewed to the right;
     more common
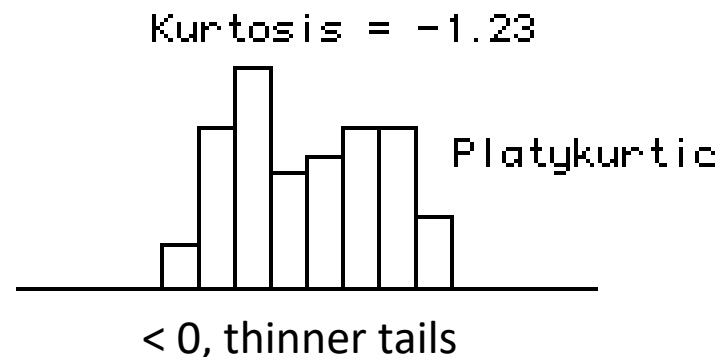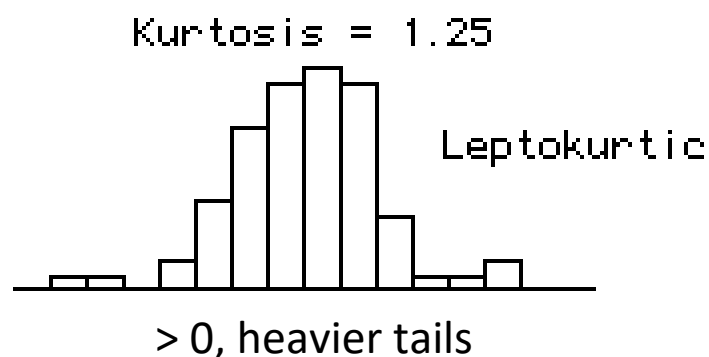If skewness if < 0:   negative skew; skewed to the left

Skewness > 0

Skewness < 0

**Kurtosis**: $4^{th}$ central moment standardized by $s^4 - 3$, where $s^4 = \frac{1}{n}\sum_{i=1}^{n}(X_i - \bar{X})^4$ and 3 is
the kurtosis for a true normal distribution

kurtosis = 0, tails just like a normal distribution
kurtosis > 0 heavier tails than a normal distribution (*leptokurtosis*)
kurtosis < 0 lighter tails than a normal distribution (*platykurtosis*)

Kurtosis = 1.25                                    Kurtosis = -1.23

Leptokurtic                                       Platykurtic

> 0, heavier tails                                < 0, thinner tails

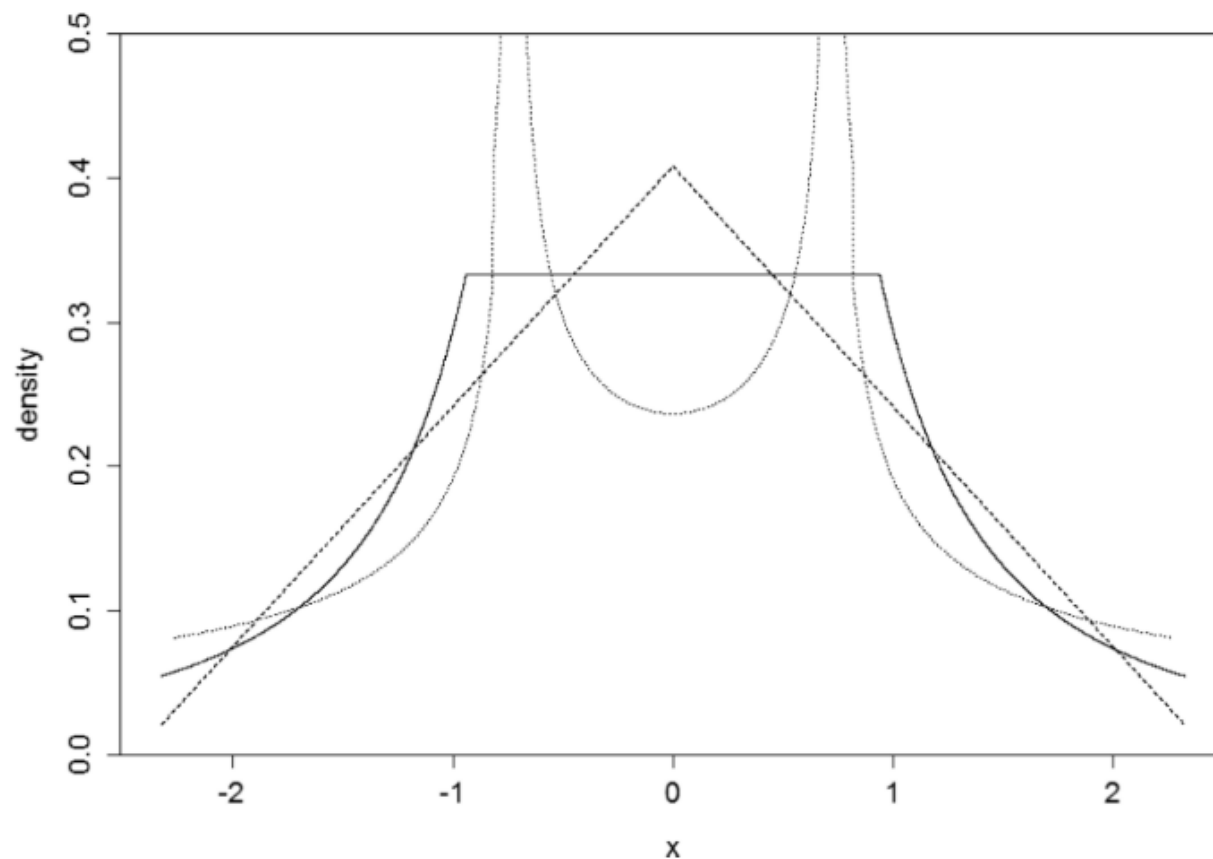## Contrary to popular belief, kurtosis has pretty much nothing to do with the peakedness of a distribution!



**Figure 2.**
Distributions with identical kurtosis = 2.4: solid = devil's tower, dashed = triangular, dotted = slip-dress.

Ref: Kurtosis as Peakedness, 1905-2014, *R.I.P.* Peter H. Westfall, *The American Statistician*, August, 68:191-195.

## B) Normal Approximation to Binomial and Poisson Distributions

Let's now see how the normal distribution is a limiting distribution for (or approximation to) the discrete distributions that we've talked about, the binomial and Poisson. This helps us to understand the how these distributions are similar and how they are different.

We know that the Normal distribution is a symmetric distribution with certain probabilities in specific parts of the distribution. Because of this, the normal will only approximate binomial and Poisson distributions that are, as you would suspect, fairly *symmetric*.
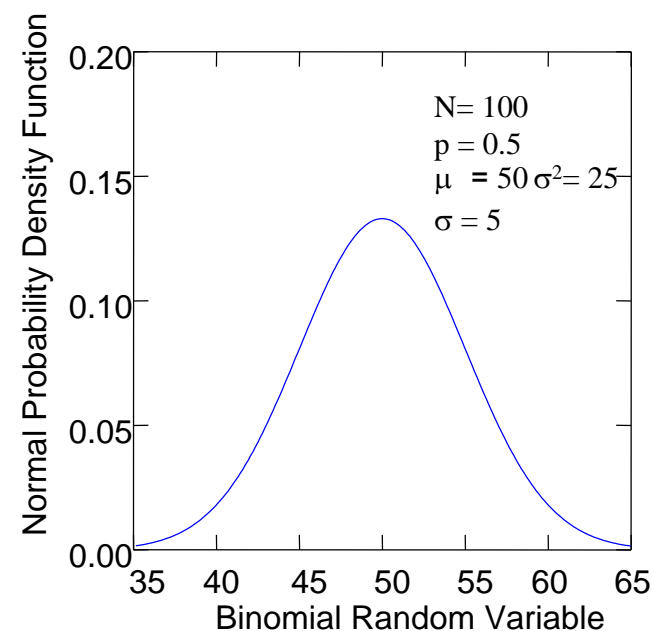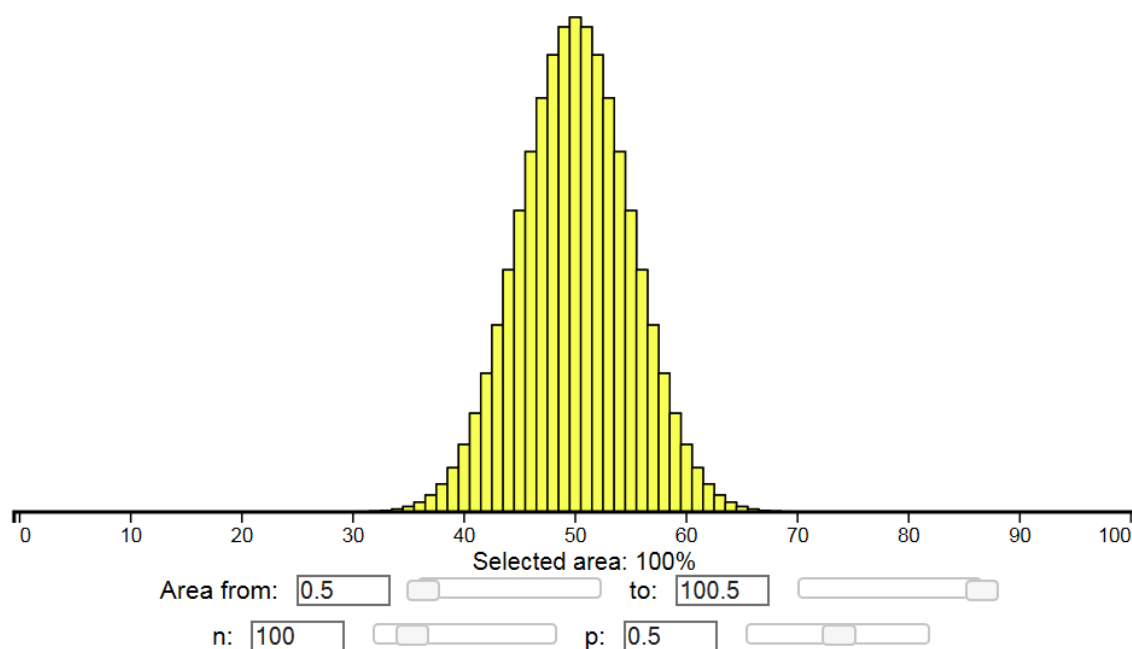
Recall that the binomial distribution has parameters *n* and *p* that completely define it, and that $\mu = np$, and $\sigma^2 = np(1-p)$. We saw that the variance of the binomial is maximized at p = 0.5. Binomial distributions with p close to 0.5 are the closest to symmetry. In general, if $np(1-p) \geq 5$, the binomial distribution looks fairly symmetric and is a possible candidate for a normal approximation.

Since the binomial is discrete, but the normal is continuous, how do we spread out the probability mass in the binomial over an infinity of points? We apply what's called a *correction for continuity*.

Here's how it works:

In the binomial, we can find P(X = x), but in the normal we know that this must somehow become an interval of values P(a $\leq$ X $\leq$ b). The interval that's used in practice is: P(x-0.5 $\leq$ X $\leq$ x+0.5), we take a half step up from *x* and a half step down from *x* to capture the probability area in the normal distribution that belongs to the value *x* in the binomial distribution. In this way we can capture every bit of probability under the curve.

https://www.stat.berkeley.edu/~stark/Java/Html/BinHist.htm

Now, if we want to find probabilities associated with the binomial, and the normal approximation conditions are met, we can use the Z-transformation, $\frac{X-\mu}{\sigma}$, where Z ~ Normal(0, 1).

e.g. Consider a sample of 100 adult where the probability of having 1+ colds in a year is 0.5. What is the probability that 50 adults have 1+ cold in a year?

We know that $np$ = $\mu$ = 100*0.50 = 50 and $\sigma$ = $\sqrt{(np(1-p)}$ = $\sqrt{(100*0.5*0.5)}$ = 5.

Applying the continuity correction, our probability becomes

$$
\begin{aligned}
P(50 - 0.5 \leq X \leq 50 + 0.5) \quad &= P(49.5 \leq X \leq 50.5) \\
&= P\left(\frac{49.5 - 50}{5} \leq Z \leq \frac{50.5 - 50}{5}\right) \\
&= P(-0.1 \leq Z \leq 0.1) \\
&= 0.07965567
\end{aligned}
$$

```
pnorm(0.1)-pnorm(-0.1)
[1] 0.07965567
```

How does this compare to its exact Binomial probability?

```
dbinom(50, 100, 0.5)
[1] 0.07958924
```

(note that accuracy is within +/- 0.001)

What if we <u>didn't</u> apply the correction for continuity? P(X=50) = ____

What about P(X $\geq$ 55)?

This becomes P (X $\geq$ ____ ) = P( Z $\geq$ _____ ) = P(Z $\geq$ 0.9) = 0.184061

How does this compare to its exact Binomial probabilities?

```
1 - pnorm(0.9)
[1] 0.1840601

sum(dbinom(55:100, 100, 0.5))
[1] 0.1841008
```

What if we <u>didn't</u> apply the correction for continuity?
P(X > 55) = P( Z > (55-50)/5 ) = P(Z > 1) = 0.1586553
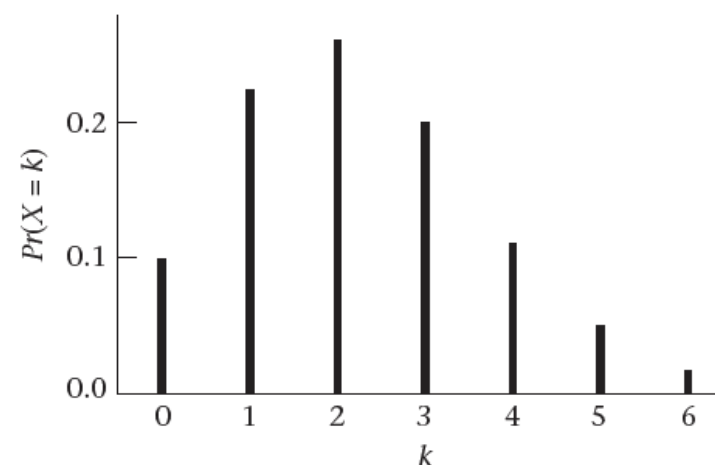
```
1 - pnorm(1.0)
[1] 0.1586553
```

What about P(X $\leq$ 40)?  (obtain on your own …)
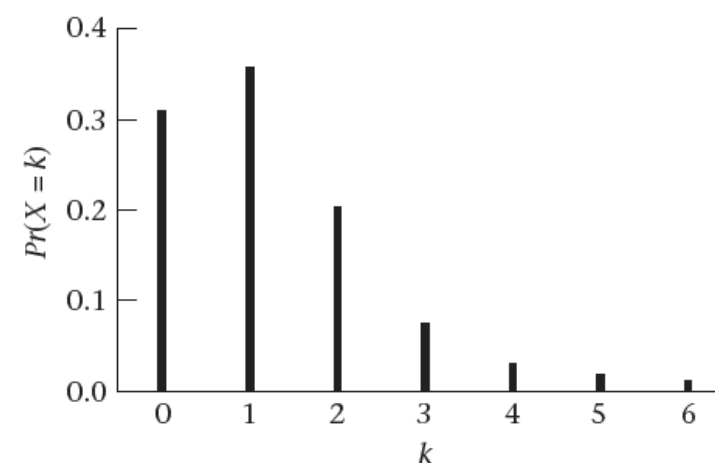
## Normal Approximation to the Poisson

Now, let's see what conditions must hold for the normal distribution to approximate the Poisson. Recall that the Poisson distribution is defined by one parameter, $\lambda$, and that $\mu = \lambda$ and $\sigma^2 = \lambda$ (the mean and the variance are equal). We know that Poisson distributions can be very right-skewed, with lots of probability mass on low values and little probability mass on large values.

So, it will take a large value of $\lambda$ to produce a symmetric distribution. $\lambda \geq 10$ *generally* achieves this.

**Figure 4.5**  Distribution of the number of deaths attributable to typhoid fever over various time intervals
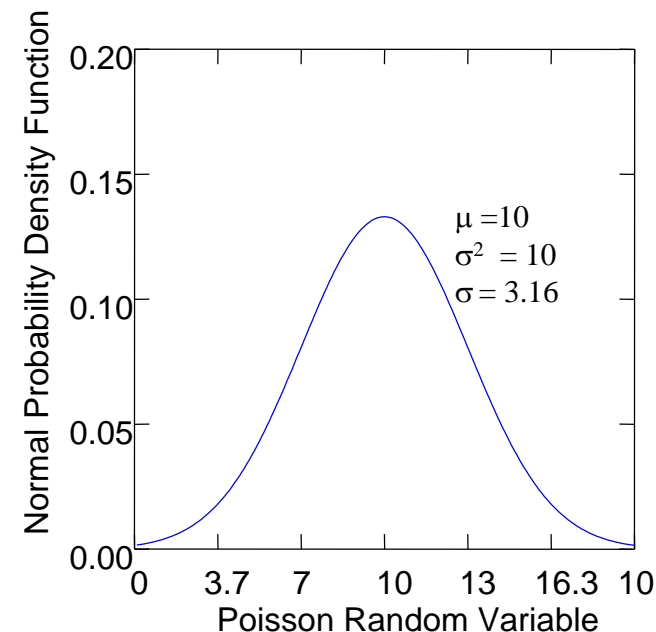


(a) 6 months



(b) 3 months

*Rosner, Fundamentals of Biostatistics, 7th edition, 2010*

We apply the same kind of continuity correction for finding Poisson probabilities based on the normal approximation as we do for the binomial.

Let's say we want to find P(X = 10) – where X is the number of bacteria colonies growing on a culture dish. Assume $\lambda$ = 10. Using the continuity correction we have

P( ___ $\leq$ X $\leq$ ___ ) =



What is the exact Poisson probability for this example, and how does it compare to the approximation? *(Is the accuracy within +/- 0.001?)*

```
pnorm(0.16)-pnorm(-0.16)
[1] 0.1271189

dpois(10,10)
[1] 0.12511
```

What about P(X ≥ 15)?

With the continuity correction, P(X ≥ ___ ) = P(Z > _____ ) = P(Z > ____ ) = 0.077

What is the exact Poisson probability for this example, and how does it compare to the approximation? *(Is the accuracy within +/- 0.001?)*

```
1 - pnorm(_____)
[1] 0.07736446

1 - ppois(14, 10)
[1] 0.08345847
```

If λ = 1000, is the normal approximation's accuracy within +/- 0.001 for P(975 ≤ X ≤ 1050)?

```
pnorm(_____)-pnorm(_____)
[1] 0.7348499

ppois(1050,1000) - ppois(974,1000)
[1] 0.7334318
```

## C) Assessing normality in a sample

Things to look for regarding normality of an empirical (sample) distribution:

Does the _mean = median = mode_? This will be true for any symmetric distribution. Is that enough to be called normal? ... no!

Other measures have to do with how skewed the distribution is. As we've talked about previously, a normal distribution has a _skewness_ value of 0. A skewness value greater than 0 indicates a distribution is skewed to the right (+), less than 0 means skewed to the left (-).

Normal distributions have only so much probability mass in the middle and only so much out in the extremes or tails of the distribution. The measure called _kurtosis_ measures (mostly, according to Peter Westfall's work) how light or heavy tailed a distribution is. A true normal distribution has a kurtosis value of 3.

_Plots_ are probably the best measures of normality - our eyes are really pretty good at pattern recognition. Be sure to use yours every time you approach a set of data!

A *Stem-and-Leaf* plot looks like a histogram turned on its sides - when we smooth the histogram with our eyes do we see symmetry? Do we see a too heavy-tailed or a too light-tailed distribution?

```
set.seed(515)

stemdat <- rnorm(50, mean=10)

stem(stemdat)

   The decimal point is at the |

    7 |
    8 | 002
    8 | 669
    9 | 1123
    9 | 566777888889
   10 | 00123334
   10 | 5555677899
   11 | 01134
   11 | 58
   12 | 12
   12 | 5
```

```
set.seed(515)

stemdat2 <- rexp(50, rate=1)

stem(stemdat2)

   The decimal point is at the |

    0 | 0001112222233344444
    0 | 555666777777888
    1 | 012223
    1 | 5799
    2 | 04
    2 |
    3 | 4
    3 | 9
    4 | 0
    4 | 8
```
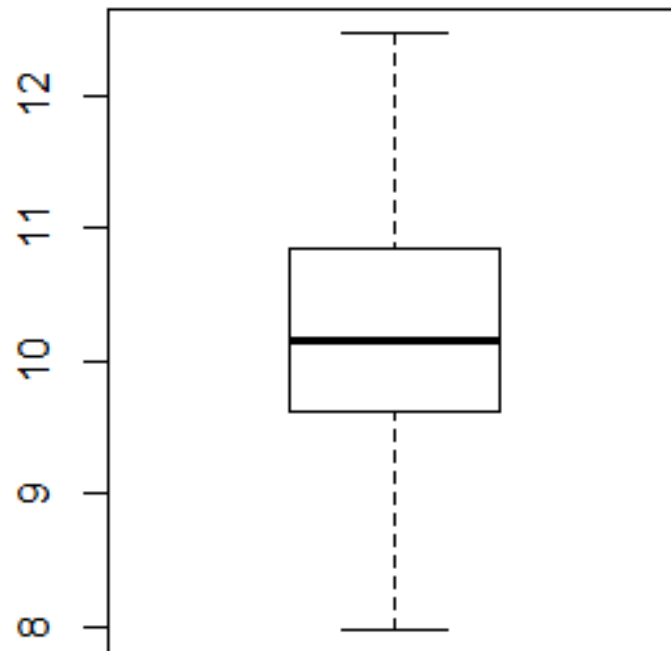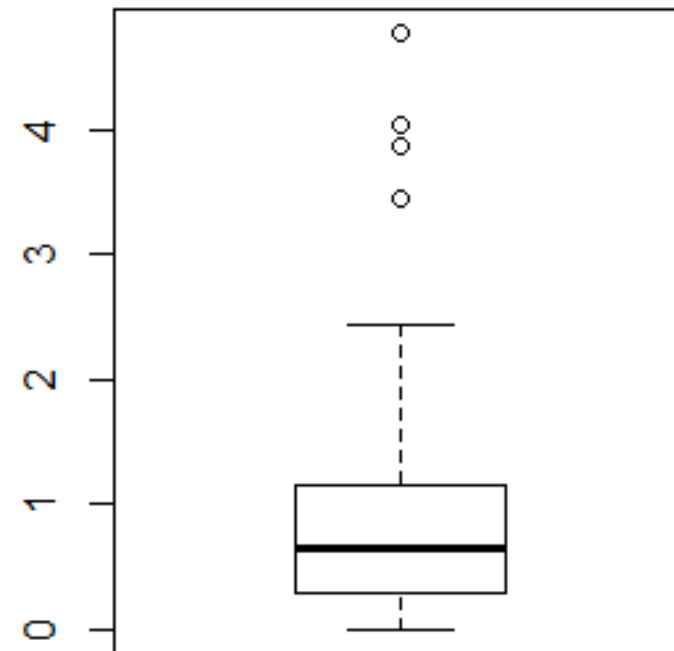
The _box plot_ is also useful to us - the distances of the hinges to the median should be the same for a symmetric distribution. (R function: boxplot.)

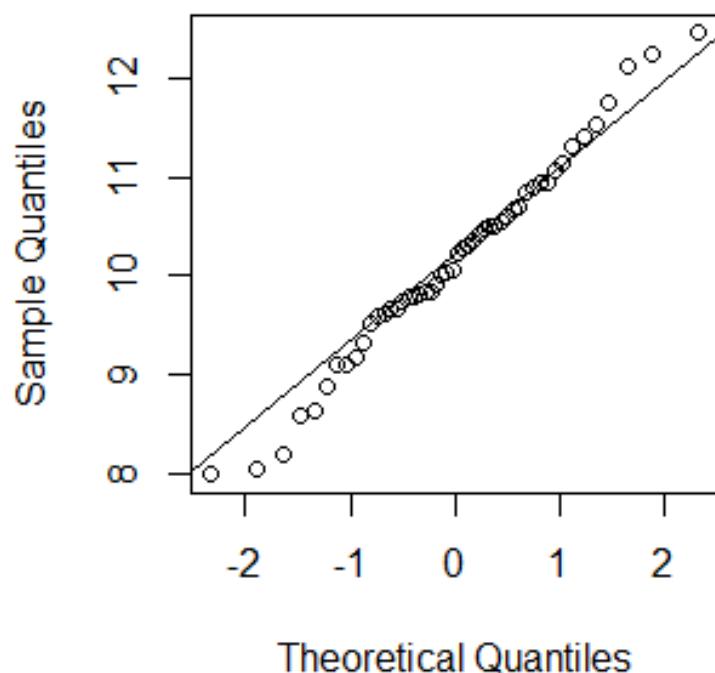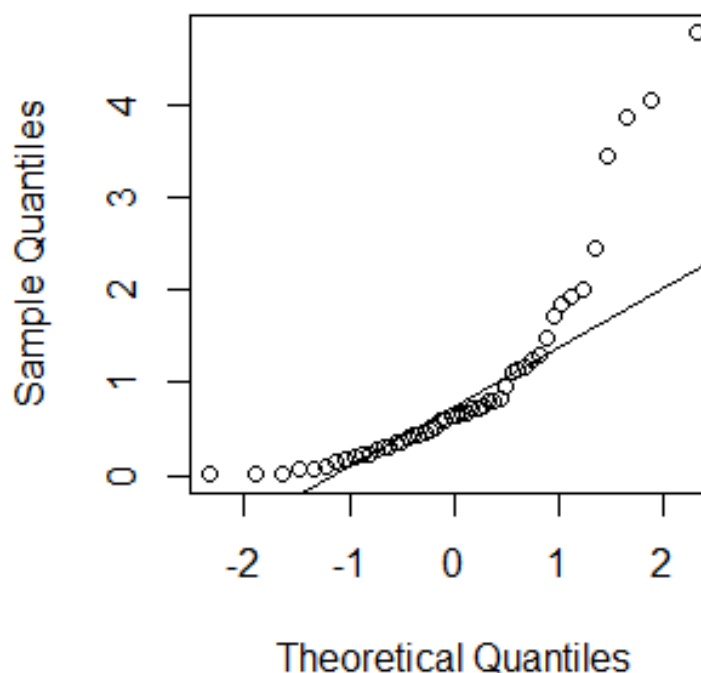Another plot that we can look at is the _Normal probability plot_. This is also known as a Q-Q plot and its shows actual (approximate) percentiles from the data versus expected percentiles (expressed as standard deviation units) of a normal distribution. The expected percentiles are based on those of a normal distribution with $\mu = \bar{X}$ and $\sigma^2 = s^2$. Deviations from a straight line are informative. The appendix includes a sheet with some patterns that you might see and what they suggest. (R functions: `qqnorm` and `qqline`.)



stemdat (Normal) — Sample Quantiles vs Theoretical Quantiles

stemdat2 (Exponential) — Sample Quantiles vs Theoretical Quantiles

Finally, there are a number of _statistical tests_ that look for deviations from normality. Attached to these tests is a probability that you'd see a value of that statistic or one bigger (in most cases; or one smaller for one of the tests**) if the data are really, really normal, i.e. a _p-value_ of this test.
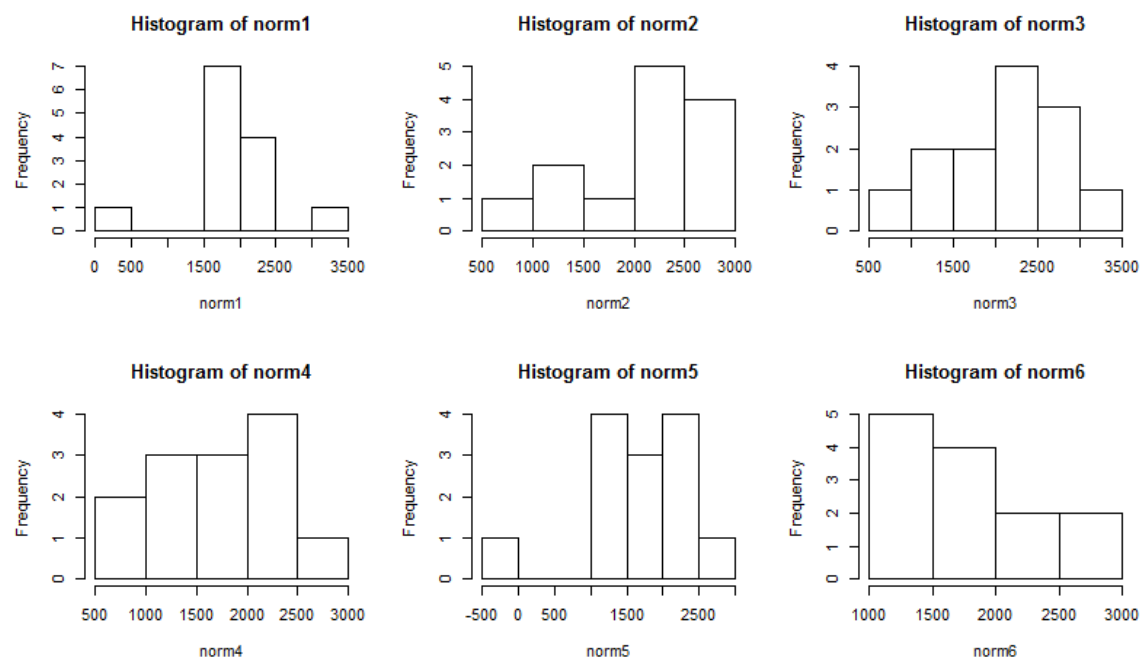
It's best _not_ to weight the results of these tests too heavily - small sample sizes can lead to a decision of normality even when the data _are not_ normal and large sample sizes can lead to a decision of non-normality even when the data _are_ very close to normal.

| Normality Test (Statistic) | Function in R (Package) | stemdat Statistic; p-value | stemdat2 Statistic/p-value |
|---|---|---|---|
| Shapiro-Wilk (W)** | `shapiro.test` (stats package) | W=0.987; p=0.846 | W=0.753; p<0.001 |
| Anderson-Darling (D) | `ad.test` (nortest package) | D=0.073; p=0.730 | D=0.230; p<0.001 |
| Cramer-von Mises (W) | `cvm.test` (nortest package) | W=0.028; p=0.871 | W=0.711; p<0.001 |
| Lilliefors/Kolmogorov-Smirnov (A) | `lillie.test` (nortest package) | A=0.192; p=0.892 | A=4.087; p<0.001 |

Notice the difference in statistic values between our normal and exponential data sets, and how Shapiro-Wilk is significant for smaller W statistics.

One of the best ways to assess normality is through simulation. Generating a series of histograms from simulated data with a given mean and variance can remind us of the variety of patterns that are possible even when the underlying distribution is truly normal.

Based on a data set of dieters, if $\bar{x}$ = 1844 kcal and s = 638 kcal, here are six samples of size 13 randomly generated from a normal distribution with μ = 1844, σ = 638.



```r
# Set up for 2x3 matrix of histograms
par(mfrow=c(2,3))

# Generate six random samples, histograms
set.seed(125)

norm1 <- rnorm(13, mean = 1844, sd = 638)
hist(norm1)

norm2 <- rnorm(13, mean = 1844, sd = 638)
hist(norm2)

norm3 <- rnorm(13, mean = 1844, sd = 638)
hist(norm3)

norm4 <- rnorm(13, mean = 1844, sd = 638)
hist(norm4)

norm5 <- rnorm(13, mean = 1844, sd = 638)
hist(norm5)

norm6 <- rnorm(13, mean = 1844, sd = 638)
hist(norm6)
```

To summarize, it's best to take a _holistic_ approach to assessing normality. Look at all of the available evidence and gather a full impression.

So why the heavy emphasis on the normal distribution? Along with the sentiments expressed by Youden, many of the statistical tests used in inferential statistics have their basis in the normal distribution.

<div align="center">

THE

NORMAL

LAW OF ERROR

STANDS  OUT  IN THE

EXPERIENCE OF MANKIND

AS  ONE  OF  THE   BROADEST

GENERALIZATIONS OF NATURAL

PHILOSOPHY  -   IT SERVES AS  THE

GUIDING INSTRUMENT IN RESEARCHES

IN THE PHYSICAL AND SOCIAL SCIENCES AND

IN MEDICINE  AGRICULTURE  AND  ENGINEERING.

IT IS AN INDISPENSABLE TOOL FOR THE ANALYSIS AND THE

INTERPRETATION OF THE DATA OBTAINED BY OBSERVATION AND EXPERIMENT

W.J. YOUDEN

</div>