

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

L15: Normal and Binomial distributions

February 28, 2020

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Today's objectives

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

- ▶ Normal quantile plots
- ▶ Introduce the binomial distribution
- ▶ What kinds of outcomes follow a binomial
- ▶ Understanding the probability space for a binomial
- ▶ What is the theoretical distribution for binomials

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

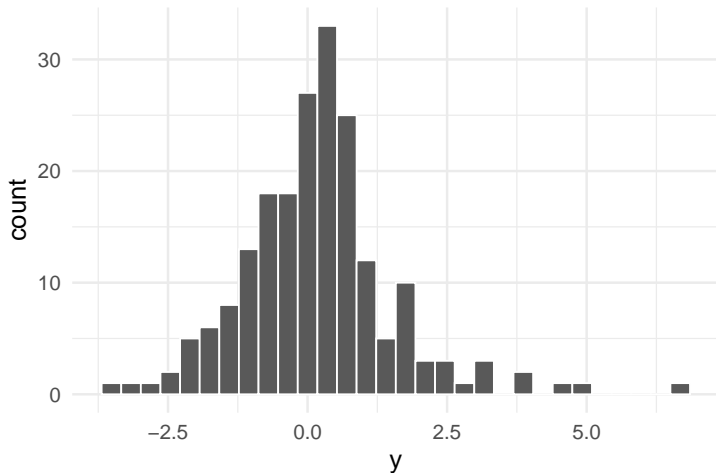
The Normal quantile plot (a.k.a the Q-Q plot)

The Normal quantile plot (a.k.a the Q-Q plot)

- ▶ The purpose of making a Q-Q plot is to examine the Normality of a distribution of a variable.
- ▶ If you want to know whether variable is Normally distributed you could examine its histogram to see if it is unimodal and symmetric. However, it is still sometimes hard to say if it is truly Normal. To do so you can use a Q-Q plot.

Are these data Normally distributed?

- The data is unimodal and symmetric, but is its distribution Normal?



The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Making a QQ plot step by step

1. First, arrange the variable in ascending order. Calculate the percentile for each measurement. For example if you had ten measurements in ascending order, the first measurement is at the 10th percentile because 10% of the data is at or below its value. The second measurement is at the 20% because 20% of the data is at or below its value, and so forth.
2. Then, for each of the percentiles you calculated, use that percentile to calculate the corresponding x-value of the Normal distribution that occurs at that percentile. For example, at $x = -1$ at the 16th percentile of the $N(0, 1)$ distribution.
3. Make a scatter plot of the calculated x-values on the x-axis and the original variable values on the y-axis.
4. The closer the data is to a straight line, the more closely it approximates a Normal distribution.

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Making a QQ plot step by step

#1. calculate the percentile :

```
example_data <- example_data %>% arrange(y) %>%  
  mutate(quantile = row_number()/n())
```

2. then calculate the x-value at each percentile from the previous step

3. this x-value can be called a z-score because it is from the standard Normal

```
example_data <- example_data %>%  
  mutate(z_score = qnorm(quantile, mean = 0, sd = 1))  
head(example_data)
```

##	y	quantile	z_score
## 1	-3.450179	0.005	-2.575829
## 2	-3.239533	0.010	-2.326348
## 3	-2.671305	0.015	-2.170090
## 4	-2.496146	0.020	-2.053749
## 5	-2.479696	0.025	-1.959964
## 6	-2.252152	0.030	-1.880794

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

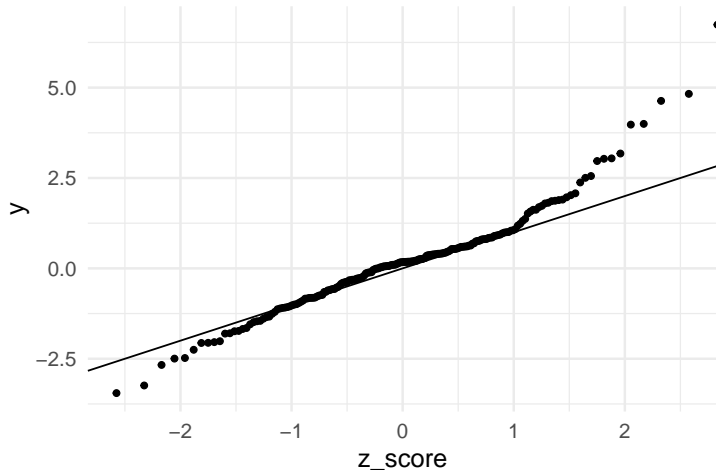
Sampling distribution of
binomial

Example trial of size 2

Example trial of size 10

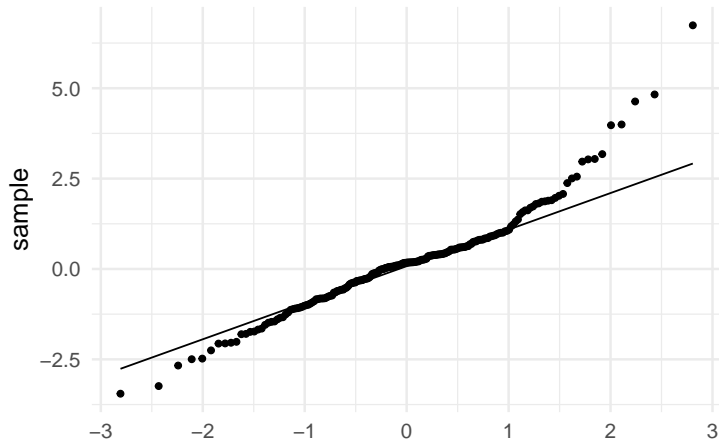
Look at the QQ plot for these data

- Notice that the data overlays the 45 degree line in the middle but not in the tails of the distribution. This sort of pattern shows that these data are “wider” (have larger standard deviation) than a Normally distributed variable.



Easy way to make a qqplot() where R does all the calculating for you

```
ggplot(example_data, aes(sample = y)) + stat_qq() + stat_qq_line()  
theme_minimal(base_size = 15)
```



Another example: Seed Data

```
library(readr)
seed_data <- read_csv("Ch04_seed-data")
```

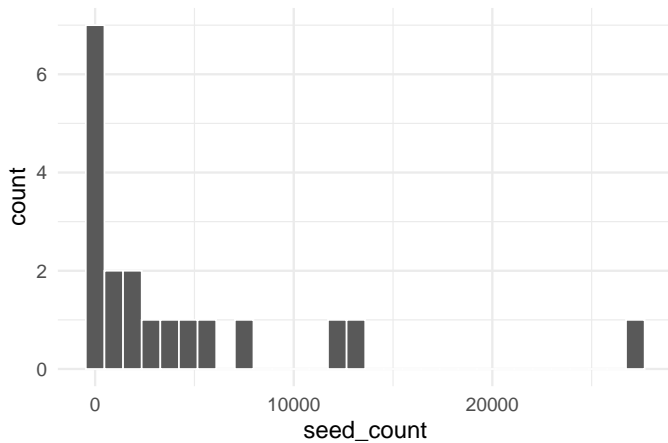
```
## Parsed with column specification:
## cols(
##   species = col_character(),
##   seed_count = col_double(),
##   seed_weight = col_double()
## )
```

```
head(seed_data)
```

```
## # A tibble: 6 x 3
##   species      seed_count seed_weight
##   <chr>          <dbl>       <dbl>
## 1 Paper birch      27239         0.6
## 2 Yellow birch    12158         1.6
```

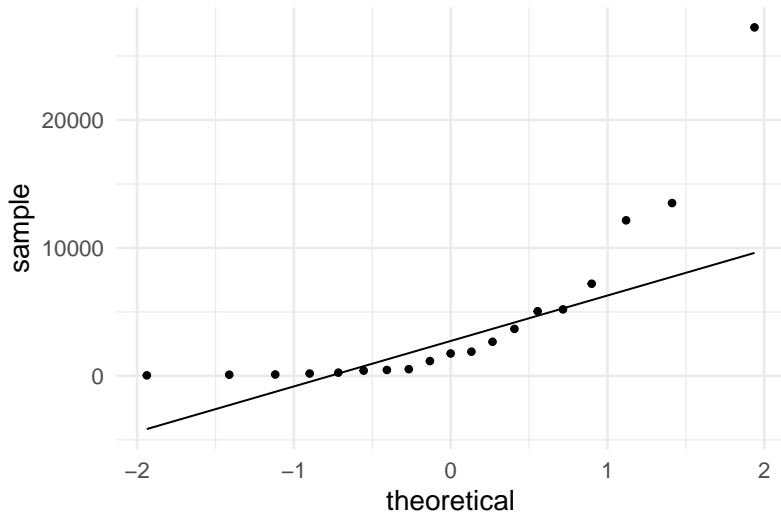
Another example

Check out its distribution. It definitely does not look normal:



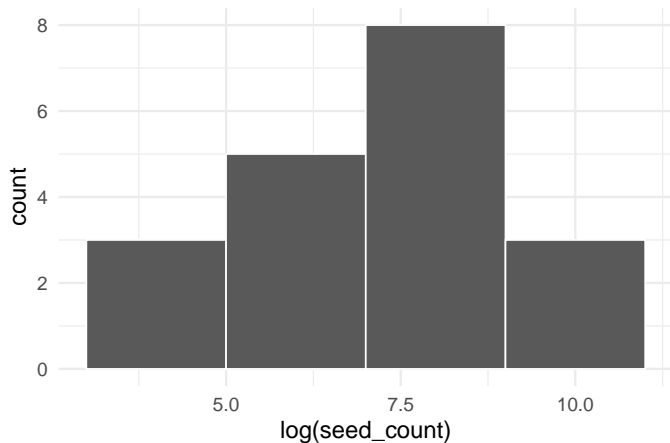
Another example

And look at its QQ plot. Does the data appear to follow a Normal distribution?



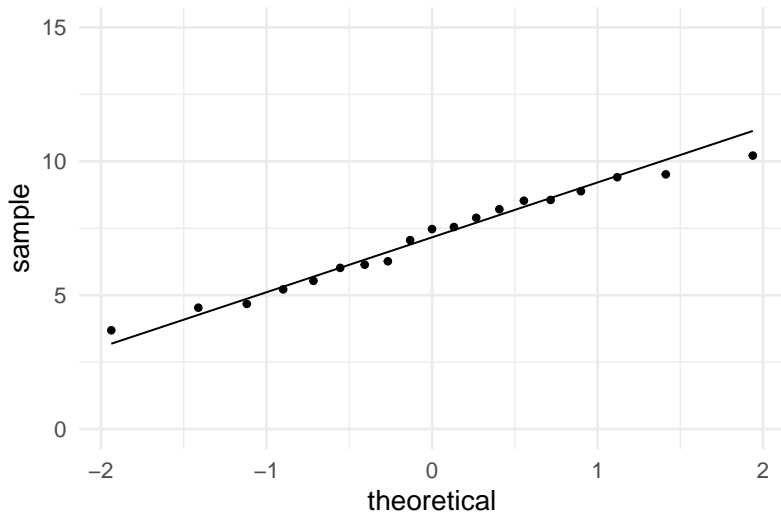
Another example (logged)

You might remember that we took the log of seed_count before we used it in regression. The log values look like this:



Another example (logged)

How does the QQ plot look for the logged variable?



The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

- ▶ Review the QQ plots from the book on page 290-292 of B&M Edition 4
- ▶ Try and gain intuition about when a variable does not appear to fit a Normal distribution
 - ▶ Was the distribution skewed?
 - ▶ Was there an outlier?
- ▶ For each scenario how do these deviations from Normality affect the QQ plot?

Other types of outcomes

We've now seen how we can use a normal probability distribution to help us evaluate continuous variables.

What about expected values(probabilities) for all the outcomes that have only 2 possibilities

The binomial setting and binomial distributions

- ▶ An elementary school administers eye exams to 800 students. How many students have perfect vision?
- ▶ A new treatment for pancreatic cancer is tried on 250 patients. How many survive for five years?
- ▶ You plant 10 dogwood trees. How many live through the winter?

What are the common threads to each of these questions?

L15: Normal and
Binomial
distributions

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

- ▶ Something is done n number of times.
- ▶ The outcome of interest for each question is categorical (binary - two levels)

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Binomial Distributions

Binomial Distributions

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10



Figure 1: Bernoulli

- ▶ Bernoulli Random Variable: The variable must assume one of two possible mutually exclusive outcomes
- ▶ Each trial of the BRV results in either a success or failure of the event happening
- ▶ Derived from the experiment: counting the number of occurrences of an event in n independent trials
- ▶ Random Variable: X = number of times the event happens in the fixed number of trials (n)
- ▶ Parameters
 - ▶ n = number of trials
 - ▶ p = probability of success (event happening)

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

- ▶ The n observations are independent. Knowing the result of one observation does not change the probabilities assigned to other observations
- ▶ Each observations is either a “success” or a “failure” (usually noted with 0 or 1). These terms are used for convenience.
- ▶ The probability of success, call it p is the same for each observation.

$$P(0 \cup 1) = P(0) + P(1) = 1$$

Example 1

A researcher has access to 40 men and 40 women and selects 10 of them at random to participate in an experiment. The number of women selected can be represented by X . Is X binomially distributed?

- Read the question carefully. What is the probability of selecting a woman when there are 40 individuals. If a woman is chosen, what is the probability of selecting a woman the second time?

Example 2

A pharmaceutical company inspects a simple random sample of 10 empty plastic containers from a shipments of 10,000. They are examined for traces of benzene. Suppose that 10% of the containers in the shipment contain benzene. Let X represent the number of containers contaminated with benzene. Is X binomially distributed?

- ▶ Issue: Each time you sample one bottle, it affects the change that the next bottle will be contaminated. However given that the population is size 10,000 and the sample size is 20, the effect of one sample's success status on the next bottle's success status is negligible.
- ▶ Here the distribution of X is *approximately* Binomial:

$$X \sim \text{Binom}(10000, 0.10)$$

where \sim is read as “approximately distributed as”.

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

If X has the binomial distribution with n observations and probability p of success on each observation, the possible values of X are 0, 1, 2, ..., n . If k is any one of these values,

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

- ▶ Read $\binom{n}{k}$ as “ n choose k ”. It counts the number of ways in which k successes can be arranged among n observations.
 - ▶ The binomial probability is this count multiplied by the probability of any one specific arrangement of the k successes.

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

**Sampling distribution of
binomial**

Example Trial of 2

Example trial of size 10

Sampling distribution of binomial

Definition: sampling distribution

A **sampling distribution** is shown as the distribution (with a histogram) of a **sample statistic** after taking many samples.

The distribution of the number of successes across many samples is called the **sampling distribution** for X . with mean $\#$ successes denoted by \bar{x}

The distribution of the proportion of successes across many samples is called the **sampling distribution** for p with mean proportion successes denoted by \hat{p}

Binomial approximation when N is much larger than n

Choose a simple random sample of size n from a population with proportion p of successes. When the population size (N) is much larger than the sample, the count X of successes in the sample has approximately the binomial distribution with parameters n and p .

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Example Trial of 2

In 1987, 29% of the adults in the United States smoked cigarettes, cigars, or pipes.

Let Y be a random variable that represents smoking status.

- ▶ $Y = 1$, an adult is currently a smoker
- ▶ $Y = 0$, an adult is not a current smoker

The two values of Smoking status are mutually exclusive and exhaustive.

What is the probability a randomly selected person is a smoker? $P(Y=1)$

What is the probability a randomly selected person is a non-smoker? $P(Y=0)$

Suppose that we randomly select two individuals from the population of adults in the United States.

The random variable X represents the number of persons in the pair who are current smokers.

First Person(Y_1)	Second Person (Y_2)	Probability	Number of Smokers (X)
0	0		0
1	0		1
0	1		1
1	1		2

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Smoking status

Suppose that we randomly select two individuals from the population of adults in the United States.

The random variable X represents the number of persons in the pair who are current smokers.

First Person (Y_1)	Second Person (Y_2)	Probability	Number of Smokers (X)
0	0	$(1 - p) \times (1 - p)$	0
1	0	$p \times (1 - p)$	1
0	1	$(1 - p) \times p$	1
1	1	$p \times p$	2

Remember that we can multiply to get the probabilities here because the events are independent

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Smoking status

Recall that for one trial, $P(\text{event}) + P(\text{no event}) = 1$

So $P(\text{smoker}) + P(\text{not smoker}) = 1$

We know $P(\text{smoker}) = 29\%$ so :

$1 - 0.29 = 0.71$ (probability of non smoker)

Calculate by hand using a table

If 29% of US adults smoked, $p=0.29$, what are the values of these probabilities?

The random variable X represents the number of persons in the pair who are current smokers.

First Person(Y_1)	Second Person (Y_2)	Probability	Number of Smokers (X)
0	0	$(.71) \times (.71) = .5041$	0
1	0	$.29 \times (.71) = .2059$	1
0	1	$(.71) \times .29 = .2059$	1
1	1	$.29 \times .29 = .0841$	2

Notice here that the sum of column 3 is $= 1$

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

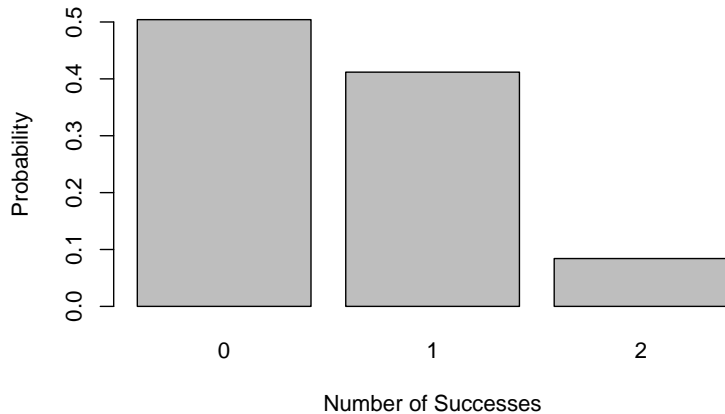
Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Probability distribution for 2 selected individuals

Probability Distribution



The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

In the binomial distribution the sum of all probabilities of potential outcomes equals 100% (1.0)

If you have a certain number of events with probability of success (p),

What is probability that X is occurs at least once $P(X \geq 1)$?

$$P(X \geq 1) = 1 - P(X = 0)$$

If you have a certain number of events with probability of success (p),

what is probability that X occurs fewer than twice $P(X < 2)$?

$$P(X < 2) = P(X = 1) + P(X = 0)$$

Binomial Probability Distributions: for n trials

What if we are interested in the expected outcomes if select a larger group of individuals? It starts to get cumbersome to write out that table by hand. The general expression of the probability distribution of a binomial random variable X where x is the number of successes in a sample of size n .

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

where $n = 1, 2, 3, \dots$ and $x = 0, 1, \dots, n$.

Binomial Combinations: How many combinations of n people give x successes?

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

$$\binom{2}{0} = \frac{2!}{0!(2-0)!} = 1$$

$$\binom{2}{1} = \frac{2!}{1!(2-1)!} = 2$$

****remember that $0! = 1$**

So for 1 success in 2 individuals ($n=2$, $x=1$) where 29% are smokers ($p=0.29$)

$$P(X = x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

$$P(X = 1) = \binom{2}{1} 0.29^1 (1 - 0.29)^{2-1} = 0.4118$$

Binomial Probability Distributions

We can use the formulas as shown to calculate the probability of a given number of successes from a binomial by hand

Or we could use R with 'dbinom(#successes,size,probability of success)'

This function calculates the probability of observing x successes when $X \sim \text{Binom}(n, p)$

```
dbinom(1,size=2,prob=0.29)
```

```
## [1] 0.4118
```

let's look further at sampling distributions using our container example. . .

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Example trial of size 10

Sampling distribution of a count in R

First, set up a large population of size 10,000 where 10% of the containers are contaminated by benzene. We call benzene a “success” since it is coded as 1. We can see that 10% of the containers are contaminated and 1000 bottles are “successes”

We simulate these data:

```
container.id <- 1:10000  
benzene <- c(rep(0, 9000), rep(1, 1000))  
pop_data <- data.frame(container.id, benzene)
```

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Sampling distribution of a count in R

```
# Calculate the population number of bottles contaminated by benzene and the  
# population mean proportion
```

```
pop_stats <- pop_data %>% summarize(pop_num_successes = sum(benzene),  
                                     pop_mean = mean(benzene))
```

```
pop_stats
```

```
##   pop_num_successes pop_mean  
## 1                1000     0.1
```

Sampling distribution of a count in R

Take a sample of size 10 from the population. Note that 10 is much smaller than 10,000.

- ▶ How many contaminated bottles are we expecting in the sample?
- ▶ Given that we sample 10, what is the full range of possible values we could see for X , the number of successes and p the proportion of successes?
- ▶ Which values are most likely?

```
# first sample
set.seed(1)
sample_data <- pop_data %>% sample_n(10)
sample_data %>% summarize(sample_num_successes = sum(benzene),
                           sample_mean = mean(benzene))
```

```
##      sample_num_successes sample_mean
## 1                        2          0.2
```

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Sampling distribution of a count in R

We only took one sample, and got 2 successes for a sample mean of 20%. Is that usual or unusual?

To see what is most likely, we need to imagine repeatedly taking samples of size 10 from the population and calculating the sample number of successes and proportion of successes for each sample.

For the next few slides, we focus on the sampling distribution for X .

Sampling distribution of a count in R

The embedded code takes 1000 samples each of size 10.

It then calculates the mean sample proportion and number of successes for each sample and stores all the results in a data frame.

You don't need to know how the code works.

Sampling distribution of a count in R

Here are the first rows of the data frame we made on the previous slide. Each row represents an independent sample from the population.

```
head(many.sample.stats)
```

##	sample_proportion	sample_num_successes	sample.id
## 1	0.1	1	1
## 2	0.1	1	2
## 3	0.0	0	3
## 4	0.0	0	4
## 5	0.0	0	5
## 6	0.1	1	6

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Sampling distribution of a count in R

We want to know: Of the 1000 samples, what percent observed 0 contaminated bottles? What percent observed 1 contaminated bottle? And so on. We can use dplyr functions to calculate this and plot the results in a histogram.

```
aggregated.stats <- many.sample.stats %>%  
  group_by(sample_num_successes) %>%  
  summarize(percent = n()/1000)
```

Sampling distribution of a count in R

```
aggregated.stats
```

```
## # A tibble: 6 x 2
##   sample_num successes percent
##           <dbl>     <dbl>
## 1             0     0.32
## 2             1     0.419
## 3             2     0.179
## 4             3     0.067
## 5             4     0.013
## 6             5     0.002
```

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

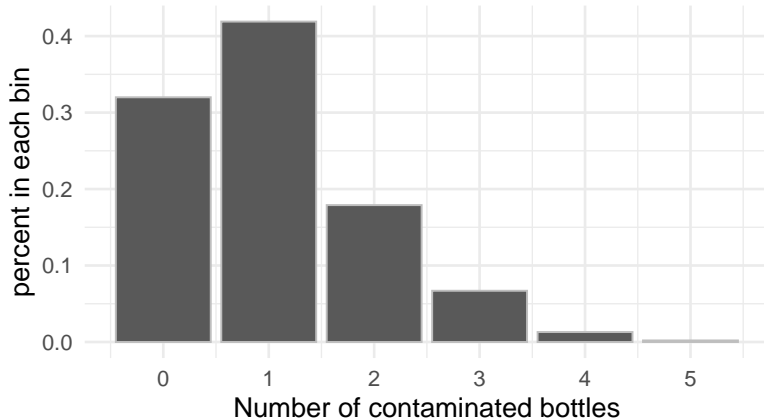
Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Sampling distribution of a count in R

Histogram of the number of
successes observed across 1000 samples



The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Sampling distribution of a count in R

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

As we will see in a moment, this histogram *approximates* the shape of the binomial distribution with $n = 10$ and $p = 0.1$. Observing one success is the most likely outcome. Why is that?

Worked probabilities, $x = 0$

We sampled $n=10$ bottles where the probability of success on any one pick is 10%.

- ▶ What is the chance of observing zero contaminated bottles?
- ▶ This means the first bottle is not contaminated and the second bottle is not contaminated, and ... and the tenth bottle is not contaminated

$$P(X_1 = 0 \text{ and } X_2 = 0 \text{ and } \dots \text{ and } X_{10} = 0)$$

$= P(X_1 = 0) \times P(X_2 = 0) \times \dots \times P(X_{10} = 0)$, using the multiplication rule for independent events

$$= (0.90)^{10}$$

$$= 0.3486784 = 34.9\%$$

Worked probabilities, $x = 1$

- ▶ What is the chance of observing exactly one contaminated bottle?
- ▶ Suppose that the first bottle was contaminated, then all the rest had to be not contaminated. What is the probability of observing this specific sequence of events?

$$\begin{aligned} &P(X_1 = 1 \text{ and } X_2 = 0 \text{ and } X_3 = 0 \text{ and...and } X_{10} = 0) \\ &= P(X_1 = 1) \times P(X_2 = 0) \times P(X_3 = 0) \dots \times P(X_{10} = 0) \\ &= (0.1)^1 (0.90)^9 \\ &= 0.03874205 = 3.87\% \end{aligned}$$

But we're not done. This is only one specific way of observing exactly one contaminated bottle. What is another way? How many ways are there to observed exactly one contaminated bottle when there are ten bottles?

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Worked probabilities, $x = 1$

There are ten ways to observe exactly one contaminated bottle:

- ▶ 1, 0, 0, 0, 0, 0, 0, 0, 0, 0
- ▶ 0, 1, 0, 0, 0, 0, 0, 0, 0, 0
- ▶ 0, 0, 1, 0, 0, 0, 0, 0, 0, 0
- ▶ 0, 0, 0, 1, 0, 0, 0, 0, 0, 0
- ▶ 0, 0, 0, 0, 1, 0, 0, 0, 0, 0
- ▶ ...
- ▶ 0, 0, 0, 0, 0, 0, 0, 0, 0, 1

Worked probabilities, $x = 1$

Remember

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

$$\binom{10}{1} = \frac{10!}{1!(10-1)!} = 10$$

Worked probabilities, $x = 1$

Each of these ten ways has the same probability of occurring.

$P(\text{observed exactly 1 contaminated bottle}) =$

$P(\text{1st bottle is contaminated, and rest are not OR 2nd bottle is contaminated, and rest are not})$

$= (0.1)^1(0.9)^9 + (0.1)^1(0.9)^9 + \dots + (0.1)^1(0.9)^9$, using the addition rule for disjoint events

$$= 10 \times (0.1)^1(0.9)^9$$

$$= 0.3874205 = 38.7\%$$

Worked probabilities, $x = 1$

We can check our calculations using the `dbinom()` function in R.

```
dbinom(x = 1, size = 10, prob = 0.1)
```

```
## [1] 0.3874205
```

This is exactly the answer we obtained.

Worked probabilities, $x = 2$

What is chance of observing exactly two contaminated bottles?

Following the same line of thinking, suppose that the first two bottles were contaminated. The chance of this happening is:

$$(0.1)^2(0.9)^8 = 0.004303672$$

But how many ways are there to observe exactly two contaminated bottles?

Worked probabilities, $x = 2$

Remember

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

$$\binom{10}{2} = \frac{10!}{2!(10-2)!} = 45$$

You could write out all the possibilities like last time, but there are a lot more!

Worked probabilities, $x = 2$

Note: we can get our calculators or R to perform this calculation for us. On our calculator, we need the button $\binom{n}{k}$, pronounced “n choose k”, and asks how many ways are there to have k successes when there are n individuals? In R we need the function `choose(n, k)`

```
choose(10, 2)
```

```
## [1] 45
```

There are 45 ways to observe exactly two contaminated bottles when you have ten bottles observed.

Make sure you can also perform this calculation on your calculator!

The Normal quantile plot
(a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of
binomial

Example Trial of 2

Example trial of size 10

Worked probabilities, $x = 2$

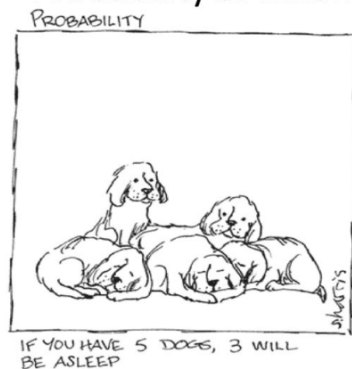
To get the probability of observing exactly 2 contaminated bottles, because all of the possible combinations are equally possible, we can multiply 45 by the probability of observing the first two bottles as being contaminated:

$$45 \times (0.1)^2(0.9)^8 = 0.1937102 = 19.4\%$$

Check using R:

```
#fill in during class
```

Comic Relief



L15: Normal and Binomial distributions

The Normal quantile plot (a.k.a the Q-Q plot)

Binomial Distributions

Sampling distribution of binomial

Example Trial of 2

Example trial of size 10