

Thinking About Data

Ryan Greenup

May 2, 2020

Contents

TODO Overheads	1
.1 TODO Install Emacs Application Framework	1
TODO Install a live preview for equations in org-mode	1
TODO Experiment with using Bookdown to Merge all RMD Files	1
Unit Information	ATTACH 1
Deriving the Normal Distribution	2
Power Series	SERIES 2
.1 Example	2
.2 Representing a function as a Power Series	2
.3 Calculus Rules and Series	3
.4 Taylor Series	3
Modelling Normal Distribution	4
.1 what is the y -axis in a Density curve? GGPlot2:ATTACH	4
.2 Defining the Normal Distribution	8
.3 Modelling only distance from the mean	9
.4 Incorporating Proportional to Frequency	10
.5 Putting both Conditions together	10
Understanding the p-value	16
False Positive Rate	17
False Discovery Rate	17
Measuring Probability	17
Comparing α and the p-value	18
Wikipedia Links	18
Calculating Power	19
Example	19
.1 Problem	19
.2 Solution	19
.3 Step 1; Find the Critical Sample Mean (\bar{x}_{crit})	20
.4 Step 2: Find the Difference between the Critical and True Means as a Z-Value (prob of Type II)	20
.5 Step 3: State the value of β	21

.6	Step 4: State the Power Value	21
Confidence Intervals		21
	Khan Academy	21
	The Confidence Interval is not the probability ATTACH	21
.1	The Reasoning	22
	Classical Confidence Interval	22
Weekly Material		23
(3);	Comparison of Population Samples WK3	23
.1	Lecture	23
.2	Tutorial	25
DONE (4);	Using Student's t-Distribution WK4	25
DONE (5)	Discrete Distributions (Mapping Disease) WK5	26
.1	DONE Quizz for t Distribution Material	26
.2	Lecture	26
.3	Tutorial	29
DONE (6)	Paired t-test (Observation or Experiment) WK6	29
TODO (7)	Corellation (Do Taller People Earn More) WK7	29
.1	DONE Lecture	29
TODO (8)	No Really do they Earn More? WK8	32
.1	TODO Lecture	32
.2	TODO Tutorial	32
wk9	Break WK9	32
TODO (9)	Do redheads have a lower pain threshold? WK10	32
TODO (10)	What is Normal? WK11	32
TODO (11)	Normality as opposed to deviant etc. WK12	32
TODO (12)	When it all goes Wrong WK13	32
TODO (13)	Exam Prep WK14	32
	Symlinks to Material	33
.1	Lecture	33
.2	Tutorial	33
Central Limit Theorem		34
<ul style="list-style-type: none"> ▪ PDF Version ▪ HTML Version <ul style="list-style-type: none"> – OnLine HTML 		

TODO Overheads

TODO Install Emacs Application Framework

This is going to be necessary to deal with not just equations but links, tables and other quirks

[Install it from here](#)

The reason for this is that generating latex preview fragments is just far too slow to be useful in any meaningful fashion.

TODO Install a live preview for equations in org-mode

[Here is one example](#) but there was a better one I was using

TODO Experiment with using Bookdown to Merge all RMD Files

[Using Bookdown Package](#)

Unit Information

ATTACH

- Learning Guide
 - [Zoom Tutorial](#)
 - [Zoom Lecture](#)

Deriving the Normal Distribution

Power Series

Series

A function f :

$$f(z) = \sum_{i=0}^{\infty} [C_n (z - a)^n], \quad \exists z \in \mathbb{C}$$

Is a [Power Series](#) a and will either:

- Converge only for $x = a$,
- converge $\forall x$
- converge in the circle $|z - a| < R$

Example

Take some function equal to the following power series:

$$f(x) = \sum_{n=0}^{\infty} [n! \cdot x^n]$$

Because the terms inside the power series has a factorial the only test that will work is the limit ratio test so we use that to evaluate convergence. ¹

let $a_n = n! \cdot x^n$:

¹Refer to [Solving Series Strategy](#)

$$\begin{aligned}\frac{\lim_{n \rightarrow \infty} |a_{n+1}|}{\lim_{n \rightarrow \infty} |a_n|} &= \lim_{n \rightarrow \infty} \left| \frac{(n+1)! \cdot x^n \cdot x}{n! \cdot x^n} \right| \\ &= (n+1) \cdot |x| \\ &= 0 \iff x = 0\end{aligned}$$

\therefore The power series converges if and only $x = 0$.

Representing a function as a Power Series

Ordinary functions can be represented as power series, this can be useful to deal with integrals that don't have an elementary anti-derivative.

1. Geometric Series First take the Series:

$$\begin{aligned}S_n &= \sum_{k=0}^n r^k \\ &= 1 + r + r^2 + r^3 \dots + r^{n-1} + r^n \\ \implies r \cdot S_n &= r + r^2 + r^3 + r^4 \dots + r^n + r^{n+1} \\ \implies S_n - r \cdot S_n &= 1 + r^{n+1} \\ \implies S_n &= \frac{1 + r^{n+1}}{1 - r}\end{aligned}$$

So now consider the geometric series:

$$\begin{aligned}\sum_{k=0}^{\infty} [x^k] &= \lim_{n \rightarrow \infty} \left[\sum_{k=0}^n x^k \right] \\ &= \lim_{n \rightarrow \infty} \left[\frac{1 + x^{n+1}}{1 - x} \right] \\ &= \frac{1 + \lim_{n \rightarrow \infty} [x^{n+1}]}{1 - x} \\ &= \frac{1 + 0}{1 - x} \\ &= \frac{1}{1 - x}\end{aligned}$$

2. Using The Geometric Series to Create a Power Series Take for example the function:

$$g(x) = \frac{1}{1 + x^2}$$

This could be represented as a power series by observing that:

$$\frac{1}{1 - \#_1} = \sum_{n=0}^{\infty} [\#_1^n]$$

And then simply putting in the value of $\#_1 = (-x^2)$:

$$\frac{1}{1 - (-x^2)} = \sum_{n=0}^{\infty} [(-x^2)^n]$$

Calculus Rules and Series

The laws of differentiation allow the following relationships:

1. Differentiation

$$\frac{d}{dx} \left(\sum_{n=1}^{\infty} c_n (z - a)^n \right) = \sum_{n=1}^{\infty} \left[\frac{d}{dx} (c_n (z - a)^n) \right]$$

2. Integration

$$\int \left(\sum_{n=1}^{\infty} c_n (z - a)^n \right) dx = \sum_{n=1}^{\infty} [c_n (z - a)^n]$$

Taylor Series

This is the important one, the idea being that you can use this to easily represent any function as an infinite series:

Consider the pattern formed by taking derivatives of $f(z) = \sum_{n=1}^{\infty} c_n (z - a)^n$:

$$f(z) = c_0 + c_1 (z - a) + c_2 (z - a)^2 + c_3 (z - a)^3 + \dots$$

$$\implies f(a) = c_0$$

$$f'(z) = c_1 + 2c_2 (z - a) + 3c_3 (z - a)^2 + 4c_4 (z - a)^3$$

$$\implies f'(a) = c_1$$

$$f''(z) = 2c_2 + 3 \times 2 \times c_3 (z - a) + 4 \times 3c_4 (z - a)^2 + \dots$$

$$\implies f''(a) = 2 \cdot c_2$$

$$f'''(z) = 3 \times 2 \times 1 \cdot c_3 + 4 \times 3 \times 2c_4 (z - a) + \dots$$

$$\implies f'''(a) = 3!c_3$$

Following this pattern forward:

$$f^{(n)}(a) = n! \cdot c_n$$

$$\implies c_n = \frac{f^{(n)}(a)}{n!}$$

Hence, if there exists a power series to represent the function f , then it must be:

$$f(z) = \sum_{n=0}^{\infty} \left[\frac{f^{(n)}(a)}{n!} (z - a)^n \right]$$

If the power series is centred around 0, it is then called a *McLaurin Series*.

1. Power Series Expansion of e

$$\begin{aligned} f(z) = e^z &= \sum_{n=0}^{\infty} \left[\frac{f^{(n)}(0)}{n!} \cdot x^n \right] \\ &= \sum_{n=0}^{\infty} \left[\frac{e^0}{n!} x^n \right] \\ &= \sum_{n=0}^{\infty} \left[\frac{x^n}{n!} \right] \end{aligned}$$

Modelling Normal Distribution

The Normal Distribution is a probability density function that is essentially modelled after observation.²

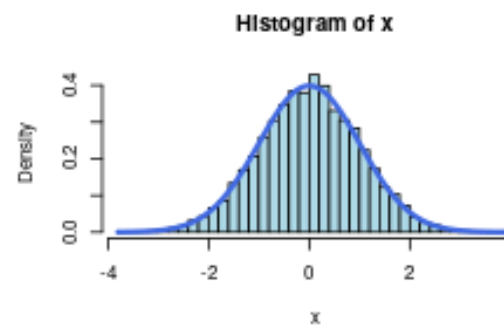
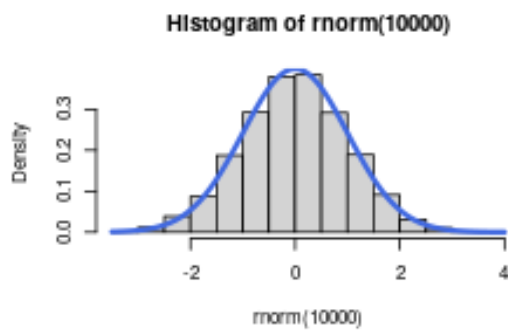
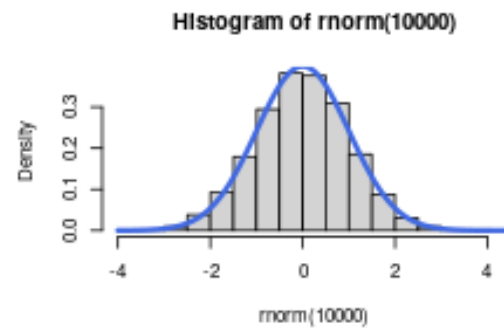
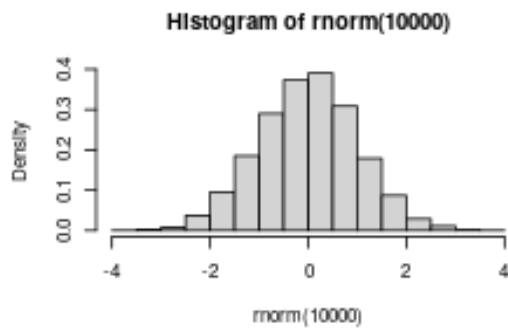
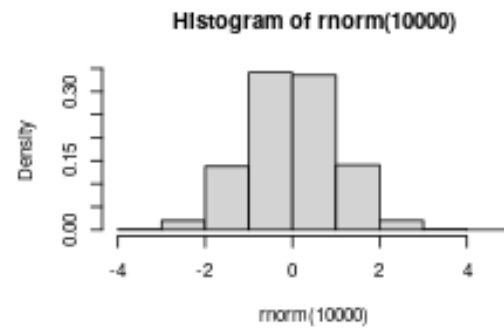
what is the y -axis in a Density curve?

ggplot2:ATTACH

Consider a histogram of some continuous normally distributed data:

```
1      # layout(mat = matrix(1:6, nrow = 3))
2      layout(matrix(1:6, 3, 2, byrow = TRUE))
3
4
5      x <- rnorm(10000, mean = 0, sd = 1)
6      sd(x)
7      hist(rnorm(10000), breaks = 5, freq = FALSE)
8      ## curve(dnorm(x, 0, 1), add = TRUE, lwd = 3, col = "royalblue")
9
10     hist(rnorm(10000), breaks = 10, freq = FALSE)
11     ## curve(dnorm(x, 0, 1), add = TRUE, lwd = 3, col = "royalblue")
12
13     hist(rnorm(10000), breaks = 15, freq = FALSE)
14     ## curve(dnorm(x, 0, 1), add = TRUE, lwd = 3, col = "royalblue")
15
16     hist(rnorm(10000), breaks = 20, freq = FALSE)
17     curve(dnorm(x, 0, 1), add = TRUE, lwd = 3, col = "royalblue")
18
19     hist(rnorm(10000), breaks = 25, freq = FALSE)
20     curve(dnorm(x, 0, 1), add = TRUE, lwd = 3, col = "royalblue")
21
22     hist(x, breaks = 30, freq = FALSE, col = "lightblue")
23     curve(dnorm(x, 0, 1), add = TRUE, lwd = 3, col = "royalblue")
```

²The Normal Distribution



(Or in ggplot2) as described in listing 1 and shown in figure ??

```

1   library(tidyverse)
2   library(gridExtra)
3   x <- rnorm(10000)
4   x <- tibble::enframe(x)
5   head(x)
6   PlotList <- list()
7   for (i in seq(from = 5, to = 30, by = 5)) {
8     PlotList[[i/5]] <- ggplot(data = x, mapping = aes(x = value)) +
9       geom_histogram(aes(y = ..density..), col = "royalblue", fill =
10         ↪ "lightblue", bins = i) +
11       stat_function(fun = dnorm, args = list(mean = 0, sd = 1))+
12       theme_classic()
13   }
14   # arrangeGrob(grobs = PlotList, layout_matrix = matrix(1:6, nrow =
15     ↪ 3))
16   grid.arrange(grobs = PlotList, layout_matrix = matrix(1:6, nrow = 3))

```


there is no package called 'tidyverse'

there is no package called 'tidyverse'

Figure 1: Histograms Generated in ggplot2

Observe that the outline of the frequencies can be made arbitrarily close to a curve given that the bin-width is made sufficiently small. This curve, known as the probability density function, represents the frequency of observation around that value, or more accurately the area beneath the curve around that point on the x -axis will be the probability of observing values within that corresponding interval.

Strictly speaking the curve is the rate of change of the probability at that point as well.

Defining the Normal Distribution

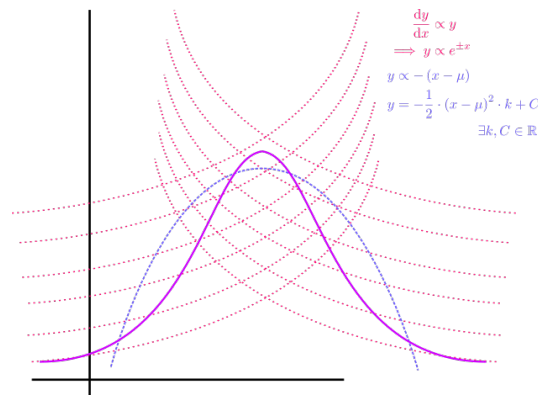
Data are said to be normally distributed if, the plot of the frequency density curve is such that:

- The rate of change is proportional to:
 - The distance of the score from the mean
 - * $\frac{d}{dx}(f) \propto -(x - \mu)$
 - The frequencies themselves.
 - * $\frac{d}{dx} \propto f$

If the Normal Distribution was only proportional to the distance from the mean (i.e. $x \propto -(x - \mu)$) the model would be a parabola that dips below zero, as shown in .3, so it is necessary to provide the restriction that the rate of change is also proportional to the frequency (i.e. $y \propto y$).

let f be the frequency of observation around x , following these rules the plot would come to look something like figure ??:

Bell Curve



Modelling only distance from the mean

If we presumed the frequency (which we will call f on the y -axis) was proportional only to the distance from the mean the model would be a parabola:

$$\begin{aligned} \frac{df}{dx} &\propto -(x - \mu) \\ \frac{df}{dx} &= -k(x - \mu), \quad \exists k \in \mathbb{R} \\ \int \frac{df}{dx} dx &= - \int (x - \mu) dx \end{aligned}$$

Using integration by substitution:

$$\begin{aligned} \text{let: } v &= x - \mu \\ \implies \frac{dv}{dx} &= 1 \\ \implies dv &= dx \end{aligned}$$

and hence

$$\begin{aligned}
\int \frac{df}{dx} dx &= - \int (x - \mu) dx \\
\Rightarrow \int dp &= - \int v dv \\
p &= -\frac{1}{2} v^2 \cdot k + C \\
p &= -\frac{1}{2} (x - \mu)^2 \cdot k + C
\end{aligned}$$

Clearly the problem with this model is that it allows for probabilities less than zero, hence the model needs to be refined to:

- incorporate a slower rate of change for smaller values of f (approaching 0)
- incorporate a faster rate of change for larger values of f
 - offset by the condition that $\frac{df}{dx} \propto -(x - \mu)$

Incorporating Proportional to Frequency

In order to make the curve bevel out for smaller values of f it is sufficient to implement the condition that $\frac{df}{dx} \propto f$:

$$\begin{aligned}
\frac{df}{dx} &\propto f \\
\int \frac{1}{f} \cdot \frac{df}{dx} dx &= k \cdot \int dx \\
\ln |f| &= k \cdot x \\
f &= C \cdot e^{\pm x} \\
f &\propto e^{\pm x}
\end{aligned}$$

Putting both Conditions together

So in order to model the bell-curve we need:

$$\begin{aligned}
f &\propto f \wedge f \propto -(x - \mu) \\
\Rightarrow \frac{df}{dx} &\propto -f(x - \mu) \\
\int \frac{1}{f} df &= -k \cdot \int (x - \mu) dx \\
\ln |f| &= -k \int (x - \mu) dx
\end{aligned}$$

because $f > 0$ by definition, the absolute value operators may be dispensed with:

$$\begin{aligned}
\ln(f) &= -k \cdot \frac{1}{2} (x - \mu)^2 + C \\
f &\propto e^{\frac{(x - \mu)^2}{2}}
\end{aligned}$$

Now that the function has been solved it is necessary to apply the IC's in order to further simplify it.

1. IC, Probability Adds to 1 The area bound by the curve must be 1 because it represents probability, hence:

$$1 = \int_{-\infty}^{\infty} f df$$

$$1 = -C \int_{-\infty}^{\infty} e^{\frac{k}{2}(x-\mu)^2} df$$

Using integration by substitution:

$$\text{let: } u^2 = \frac{k}{2} (x - \mu)^2$$

$$u = \sqrt{\frac{k}{2}} (x - \mu)$$

$$\frac{du}{dx} = \sqrt{\frac{k}{2}}$$

hence:

$$1 = -C \int_{-\infty}^{\infty} e^{\frac{k}{2}(x-\mu)^2}$$

$$1 = \sqrt{\frac{2}{k}} \cdot C \int_{-\infty}^{\infty} e^{-u^2} du$$

$$1^2 = \left(\sqrt{\frac{2}{k}} \cdot C \int_{-\infty}^{\infty} e^{-u^2} du \right)^2$$

$$1^2 = \left(\sqrt{\frac{2}{k}} \cdot C \int_{-\infty}^{\infty} e^{-u^2} du \right) \times \left(\sqrt{\frac{2}{k}} \cdot C \int_{-\infty}^{\infty} e^{-u^2} du \right)$$

Because this is a definite integral u is merely a dummy variable and instead we can make the substitution of x and y for clarity sake.

$$1^2 = \left(\sqrt{\frac{2}{k}} \cdot C \int_{-\infty}^{\infty} e^{-x^2} dx \right) \times \left(\sqrt{\frac{2}{k}} \cdot C \int_{-\infty}^{\infty} e^{-y^2} dy \right)$$

Now presume that the definite integral is equal to some real constant $\beta \in \mathbb{R}$:

$$1 = \frac{2}{k} \cdot C^2 \int_{-\infty}^{\infty} e^{-y^2} dy \times \beta$$

$$= \frac{2}{k} \cdot C^2 \int_{-\infty}^{\infty} \beta \cdot e^{-y^2} dy$$

$$= \frac{2}{k} \cdot C^2 \cdot \int_{-\infty}^{\infty} \left(\int_{-\infty}^{\infty} e^{-x^2} dx \right) e^{-y^2} dy$$

$$= \frac{2}{k} \cdot C^2 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-(x^2+y^2)} dx dy$$

This integral will be easier to evaluate in polar co-ordinates, a double integral may be evaluated in polar co-ordinates using the following relationship: ³

$$\iint_D f(x, y) dA = \int_{\alpha}^{\beta} \int_{h_1(\phi)}^{h_2(\phi)} f(r \cdot \cos(\phi), r \cdot \sin(\phi)) dr d\phi$$

hence this simplifies to:

$$1 = \frac{2}{k} c^2 \int_0^{2\pi} \int_0^r r \cdot e^{(r \cdot \cos \theta)^2 + (r \cdot \sin \theta)^2} dr d\theta$$

$$1 = \frac{2}{k} c^2 \int_0^{2\pi} \int_0^r r \cdot e^{r^2} dr d\theta$$

Because the integrand is of the form $f'(x) \times g(f(x))$ we may use integration by substitution:

$$\text{let: } u = -r^2$$

$$\frac{du}{dr} = -2r$$

$$dr = -\frac{1}{2r} du$$

and hence:

$$1 = \frac{2}{k} c^2 \int_0^{2\pi} \int_0^r r \cdot e^{r^2} dr d\theta$$

$$\Rightarrow 1 = -\frac{2}{k} c^2 \int_0^{2\pi} \int_0^{\infty} r \cdot e^{r^2} dr d\theta$$

$$1 = \frac{2}{k} c^2 \int_0^{2\pi} \int_0^r r \cdot e^{r^2} dr d\theta$$

$$\Rightarrow 1 = -\frac{2}{k} c^2 \int_0^{2\pi} \int_0^{\infty} -\frac{1}{2} e^{-u} du d\theta$$

$$= \frac{1}{k} c^2 \int_0^{2\pi} \int_0^{\infty} e^{-u} du d\theta$$

$$= \frac{1}{k} c^2 \int_0^{2\pi} [-e^{-u}]_0^{\infty} d\theta$$

$$1 = \frac{1}{k} c^2 2\pi$$

$$\Rightarrow C^2 = \frac{k}{2\pi}$$

So from before:

$$f = -C \cdot e^{k \cdot \frac{(x-\mu)^2}{2}}$$

$$= -\sqrt{\frac{k}{2\pi}} \cdot e^{k \cdot \frac{(x-\mu)^2}{2}}$$

³Calculus III - Double Integrals in Polar Coordinates

so now we simply need to apply the next initial condition.

2. IC, Mean Value and Standard Deviation

- (a) Definitions The definition of the expected value, where $f(x)$ is a probability function is: ⁴

$$\mu = E(x) = \int_a^b x \cdot f(x) dx$$

That is, roughly, the sum of the expected proportion of occurrence.

The definition of the variance is:

$$V(x) = \int_a^b (x - \mu)^2 f(x) dx$$

which can be roughly interpreted as the sum of the proportion of squared distance units from the mean. The standard deviation is $\sigma = \sqrt{V(x)}$.

- (b) Expected Value of the Normal Distribution The expected value of the normal distribution is μ , this can be shown rigorously:

$$\begin{aligned} \text{let: } v &= x - \mu \\ \implies dv &= dx \end{aligned}$$

Observe that the limits of integration will also remain as $\pm\infty$ following the substitution:

$$\begin{aligned} E(v) &= \int_{-\infty}^{\infty} v \times f(v) dv \\ &= k \cdot \int_{-\infty}^{\infty} v \cdot e^{v^2} dv \\ &= \frac{1}{2} \left[e^{x^2} \right]_{-\infty}^{\infty} \\ &= \frac{1}{2} \lim_{b \rightarrow \infty} \left[\left[e^{x^2} \right]_{-b}^b \right] \\ &= \frac{1}{2} \lim_{b \rightarrow \infty} \left[e^{b^2} - e^{(-b)^2} \right] \\ &= \lim_{b \rightarrow \infty} [0] \times \frac{1}{2} \\ &= \frac{1}{2} \times 0 \\ &= 0 \end{aligned}$$

Hence the Expected value of the standard normal distribution is $0 = x - \mu$ and so $E(x) = \mu$.

- (c) Variance of the Normal Distribution Now that the expected value has been confirmed, consider the variance of the distribution:

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 \times f(x) dx$$

⁴Expected Value and Variance

Now observe that $(x - \mu)$ appears as an exponential and as a factor if this is redefined as $w = x - \mu \implies dx = dw$ we have:

$$\sigma^2 = \sqrt{\frac{k}{2}} \int_{-\infty}^{\infty} w^2 e^{-\frac{k}{2}w^2} dw$$

Now the integrand is of the form $f(x) \times g(x)$ meaning that the only strategy to potentially deal with it is integration by parts:

$$\int u dv = u \cdot v - \int v du$$

where:

- u is a function that simplifies with differentiation
- dv is something that can be integrated

$$\begin{aligned} u &= w & dv &= w \cdot e^{-\frac{k}{2}w^2} dw \\ \implies du &= dw & \implies v &= \int w \cdot e^{-\frac{k}{2}w^2} dw \\ & & \implies v &= \frac{1}{k} e^{-\frac{k}{2}w^2} \end{aligned}$$

Hence the value of the variance may be solved:

Now that the expected value has been confirmed, consider the variance of the distribution:

$$\begin{aligned} \sigma^2 &= \int_{-\infty}^{\infty} (x - \mu)^2 \times f(x) dx \\ &= \int_{-\infty}^{\infty} (x - \mu)^2 \times \left(\sqrt{\frac{k}{2\pi}} e^{-\frac{k}{2}(x-\mu)^2} \right) dx \\ &= \sqrt{\frac{k}{2\pi}} \int_{-\infty}^{\infty} (x - \mu)^2 \times \left(e^{-\frac{k}{2}(x-\mu)^2} \right) dx \end{aligned}$$

Now observe that $(x - \mu)$ appears as an exponential and as a factor if this is redefined as $w = x - \mu \implies dx = dw$ we have:

$$\sigma^2 = \sqrt{\frac{k}{2}} \int_{-\infty}^{\infty} w^2 e^{-\frac{k}{2}w^2} dw$$

Now the integrand is of the form $f(x) \times g(x)$ meaning that the only strategy to potentially deal with it is integration by parts:

$$\int u dv = u \cdot v - \int v du$$

where:

- u is a function that simplifies with differentiation
- dv is something that can be integrated

$$\begin{aligned} u &= w & dv &= w \cdot e^{-\frac{k}{2}w^2} dw \\ \implies du &= dw & \implies v &= \int w \cdot e^{-\frac{k}{2}w^2} dw \\ & & \implies v &= \frac{1}{k} e^{-\frac{k}{2}w^2} \end{aligned}$$

Hence the value of the variance may be solved:

$$\begin{aligned}
\sigma^2 &= \sqrt{\frac{k}{2\pi}} \int_{-\infty}^{\infty} w^2 e^{-\frac{k}{2}w^2} dw \\
&= \sqrt{\frac{k}{2\pi}} \left[u \cdot v - \int v du \right]_{-\infty}^{\infty} \\
&= \sqrt{\frac{k}{2\pi}} \left(\left[\frac{-w}{k} \cdot e^{-\frac{k}{2}w^2} \right]_{-\infty}^{\infty} - \frac{1}{k} \int_{-\infty}^{\infty} e^{\frac{k}{2}w^2} dw \right) \\
&= \sqrt{\frac{k}{2\pi}} \left[\frac{-w}{k} \cdot e^{-\frac{k}{2}w^2} \right]_{-\infty}^{\infty} - \frac{1}{k} \left(\sqrt{\frac{k}{2\pi}} \int_{-\infty}^{\infty} e^{\frac{k}{2}w^2} dw \right)
\end{aligned}$$

The left term evaluates to zero and the right term is the area beneath the bell curve with mean value 0 and so evaluates to 1:

$$\begin{aligned}
\sigma^2 &= 0 - \frac{1}{k} \\
\Rightarrow k &= \frac{1}{\sigma^2}
\end{aligned}$$

So the function for the density curve can be simplified:

$$\begin{aligned}
&= -\sqrt{\frac{k}{2\pi}} \cdot e^{k \cdot \frac{(x-\mu)^2}{2}} \\
&= \sqrt{\frac{1}{2\pi\sigma^2}} \cdot e^{\frac{1}{2} \cdot \frac{(x-\mu)^2}{\sigma^2}}
\end{aligned}$$

now let $z = \frac{x-\mu}{\sigma} \Rightarrow dz = \frac{dx}{\sigma}$, this then simplifies to:

$$f(x) = \sqrt{\frac{1}{2\pi}} \cdot e^{-\frac{1}{2}z^2}$$

Now using the power series identity from BEFORE :

$$e^{-\frac{1}{2}z^2} = \sum_{n=0}^{\infty} \left[\frac{\left(-\frac{1}{2}z^2\right)^n}{n!} \right]$$

We can solve the integral of $f(x)$ (which has no elementary integral).

$$\begin{aligned}
f(x) &= \sqrt{\frac{1}{2\pi}} \cdot \sum_{n=0}^{\infty} \left[\frac{\left(-\frac{1}{2}z^2\right)^n}{n!} \right] \\
\int f(x) dx &= \frac{1}{\sqrt{2\pi}} \int \sum_{n=0}^{\infty} \left[\frac{\left(-\frac{1}{2}z^2\right)^n}{n!} \right] dz \\
&= \frac{1}{\sqrt{2\pi}} \cdot \sum_{n=0}^{\infty} \left[\int \frac{(-1)^{-1} z^{2n}}{2^n \cdot n!} dz \right] \\
&= \frac{1}{\sqrt{2\pi}} \cdot \sum_{n=0}^{\infty} \left[\frac{(-1)^n \cdot z^{2n+1}}{2^n (2n+1) n!} \right]
\end{aligned}$$

Although this is a power series it still gives a method to solve the area beneath the curve of the density function of the normal distribution.

Understanding the p-value

Let's say that I'm given 100 vials of medication and in reality only 10 of them are actually effective.

POS	POS	POS	POS	POS	POS	POS	POS	POS	POS
:-:	:-:	:-:	:-:	:-:	:-:	:-:	:-:	:-:	:-:
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG
NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG	NEG

We don't know which ones are effective so It is necessary for the effective medications to be detected by experiment. Let:

- the p-value be 9% for detecting a significant effect
- assume the statistical power is 70%

So this means that the corresponding errors are:

1. Of the 90 Negative Drugs, $\alpha \times 90 \approx 8$ will be identified as Positive (False Positive) a. This means 72 will be correctly identified as negative. (TN)
2. Of the 10 Good drugs $\beta \times 10 = 3$ will be labelled as negative (False Negative) b. This means 8 will be correctly identified as positive (True Positive)

These results can be summarised as:

	Really Negative	Really Positive
Predicted Negative	TNR; $(1 - \alpha)$	FNR; $\beta \times 10 = 3$
Predicted Positive	FPR; $FPR = \alpha \times 90 \approx 8$	TPR $(1 - \beta)$

And a table visualising the results:

TP	TP	TP	TP	TP	TP	TP	FN	FN	FN
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
FP	FP	FP	FP	FP	FP	FP	FP	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN
TN	TN	TN	TN	TN	TN	TN	TN	TN	TN

So looking at this table, it should be clear that:

- If the null hypothesis had been true, the probability of a False Positive would indeed have been $\frac{8}{90} \approx 0.09$
- The probability of incorrectly rejecting the null hypothesis though is the number of FP from anything identified as positive $\frac{7}{7+6} \approx 0.5$

False Positive Rate

The False Positive Rate is expected to be α it is:

$$\begin{aligned}
 E(\text{FPR}) &= \alpha; \\
 \text{FPR} &= \frac{FP}{N} \\
 &= \frac{FP}{FN + TP} \\
 &= \frac{8}{8 + 72} \\
 &= 9\%
 \end{aligned}$$

False Discovery Rate

The False discovery Rate is the proportion of observations considered as positive (or significant) that are False Positives. If you took all the results you considered as positive and pulled one out, the probability that one was a false positive (and you were committing a type I error) would be the FDR and could be much higher than the FPR.

Measuring Probability

In setting α as 9% I've said that 'if the null hypothesis was true and every vial was negative, 9% of them would be false positives', this means that in practice 9% of the negative vials would be detected as false

positives (I wouldn't count the positives because my α assumption was made under the assumption that everything was negative, hence 9% of the negative vials will be false positives).

So this measures the probability of rejecting the null hypothesis if it were true.

It does not measure the probability of rejecting the null hypothesis but then being mistaken, because to reject the null hypothesis it is necessary to consider observations that are considered positive (whether or not they actually are), the number of those that are False Positive would represent the probability of committing a type 1 error in that experiment

So the p -value measures the probability of committing a type I error under the assumption that the null hypothesis is true.

The FDR represents the actual probability of committing a type I error when taking multiple comparisons.

Comparing α and the p -value

The distinction between α and p -value is essentially that the α value is set as a significance standard and the p -value represents the probability of getting a test-statistic \geq the observed value

The α value is the probability of

Rejecting the null hypothesis under the assumption that the null hypothesis is true.

This will be the False Positive Rate:

The proportion of Negative Observations misclassified as Positive will be the False Positive Rate.

Be careful though because this is not necessarily the *probability of incorrectly rejecting the null hypothesis* there is also the the $FDR = \frac{FP}{TP+FP}$:

The proportion of observations classified as positive that are false positives, this estimates the probability of rejecting the null hypothesis and being wrong. (whereas the α value is the probability of rejecting the null hypothesis under the assumption it was true this is different from the probability of rejecting H_0 and being wrong, which is the FDR).

The p -value is the corresponding probability of the test statistic that was returned, so they mean essentially the same thing, but the α value is set before hand and the p -value is set after the fact:

The p -value is the probability, under the assumption that there is no true effect or no true difference, of collecting data that shows a difference equal to or more extreme than what was actually observed.

Wikipedia Links

Helpful Wikipedia Links

- [False Positive Rate](#)
- [False Discovery Rate](#)
- [Sensitivity and Specificity](#)

- ROC Curve
 - This has all the TP FP calculations
- Type I and Type II Errors
 - This has the useful Tables and SVG Density Curve

Calculating Power

Statistical Power is the probability of rejecting the null hypothesis assuming that the null hypothesis is false (True Positive).

Complementary to the *False Positive Rate* and *False Detection Rate*, the power is distinct from the probability of correctly rejecting the null hypothesis, which is the probability of selecting a True Positive from all observations determined to be positive (the Positive Predictive Value or the [Precision](#)):

$$PPV = \frac{TP}{TP + FP}$$

$$FDR = \frac{FP}{TP + FN}$$

$$\alpha = \frac{FP}{N} = \frac{TP}{TN + FP}$$

$$\beta = \frac{TP}{P} = \frac{TP}{TP + FN}$$

Example

Problem

An ISP stated that users average 10 hours a week of internet usage, it is already known that the standard deviation of this population is 5.2 hours. A sample of $n = 100$ was taken to verify this claim with an average of \bar{x} .

A worldwide census determined that the average is in fact 12 hours a week not 10.

Solution

1. Hypotheses

- (a) H_0 : **The Null Hypothesis** that the average internet usage is 10 hours per week
- (b) H_a : **The Alternative Hypothesis** that the average internet usage exceeds 10 hours a week

2. Data

Value	Description
$n = 100$	The Sample Size
$\sigma = 5.2$	The Standard Deviation of internet usage of the population
$\mu = 10$	The alleged average internet usage.
$\bar{x} = 11$	The average of the sample
$\mu_{True} = 12$	The actual average from the population
$\alpha = 0.05$	The probability of a type 1 error at which the null hypothesis is rejected
$\beta = ??$	The probability of a type 2 error

Step 1: Find the Critical Sample Mean (\bar{x}_{crit})

The Central Limit Theorem provides that the mean value of a sample that is:

- sufficiently large, or
- drawn from a normally distributed population

will be normally distributed, so if we took various samples of a population and recorded all the sample means in a set \bar{X} we would have: will be normally distributed, so if we took various samples of a population and recorded all the sample means in a set \bar{X} we would have:

$$\bar{X} \sim \mathcal{N}\left(\mu, \left(\frac{\sigma}{\sqrt{n}}\right)\right)$$

And hence we may conclude that:

$$\begin{aligned}
 Z &= \frac{\bar{x} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)} \\
 \implies \bar{x}_{crit} &= \mu + z_{\alpha} \cdot \left(\frac{\sigma}{\sqrt{n}}\right) \\
 \bar{x}_{crit} &= \mu + z_{0.05} \cdot \left(\frac{\sigma}{\sqrt{n}}\right) \\
 \bar{x}_{crit} &= \mu + 1.645 \cdot \left(\frac{5.2}{\sqrt{100}}\right) \\
 &= 10.8554
 \end{aligned}$$

Thus H_0 is rejected for a sample mean of 10.86 hours per week at a confidence level of $\alpha = 0.05$.

Step 2: Find the Difference between the Critical and True Means as a Z-Value (prob of Type II)

The probability of accepting the null hypothesis assuming that it is false, is the probability of getting a value less than the critical value given that the mean value is actually 12:

$$\begin{aligned}
 z &= \frac{\bar{x}_{crit} - \mu_{true}}{\left(\frac{\sigma}{\sqrt{n}}\right)} \\
 &= \frac{10.86 - 12}{\frac{5.2}{10}} = -2.2
 \end{aligned}$$

Step 3: State the value of β

$$\begin{aligned}\beta &= P(\text{Type II Error}) \\ &= P(H_0 \text{ is not rejected} \mid H_0 \text{ is false}) \\ &= P(\mu_{\bar{X}_{\text{crit}}} < \bar{x}_{\text{crit}} \mid \mu = 12) \\ &= 0.014\end{aligned}$$

Step 4: State the Power Value

$$\begin{aligned}\text{Power} &= (H_0 \text{ is not rejected} \mid H_0 \text{ is false}) \\ &= P(\mu_{\bar{X}_{\text{crit}}} < \bar{x}_{\text{crit}}) \\ &= 1 - \beta \\ &= 1 - 0.014 \\ &= 98.6\%\end{aligned}$$

Confidence Intervals

Khan Academy

According to [Khan Academy](#):

This means that for any sample drawn from the population, the true population value would be found within this interval for 0.95 of those samples

The Confidence Interval is not the probability

ATTACH

adapted from:

- <https://qr.ae/pNrV1x>
- <https://qr.ae/pNrV6y>

PDF Version

I assume that the motivation for this question is that most statistics books emphasize the fact that, once you have taken a sample and constructed the confidence interval (CI), there is no longer any randomness left in a CI statement (except for the Bayesian point of view which thinks of μ as being a random variable).

That is, when reporting a CI: I am 95% confident that the mean is between 25.1 and 32.6 is correct. There is a 95% probability that the mean is between 25.1 and 32.6 is WRONG. Either μ is in that interval or not; there is no probability associated with it.

The Reasoning

Suppose that somewhere on the wall is an invisible bullseye a special point (call it μ) which only I can see. I'm going to throw a dart at μ . Based on long observation, you know that when I throw a dart at something, 95% of the time, my dart will hit within 6 inches of what I was aiming at. (The other 5% of the time, I miss by more than 6 inches.) When you see where that dart lands, you will draw a circle around it with a radius of 6 inches.

It is correct to say:

The probability that μ will be in that circle is 95%.

The reason that is correct is, I have not yet thrown the dart, so the location of the circle is random, and in 95% of repetitions of this dart-throwing, circle-drawing routine, μ will be in the resulting circle. Now if I actually take aim, throw my dart, and it hits

right here \implies .

It is no longer correct to talk about probabilities. You can be pretty sure that μ is within that circle. To be specific, pretty sure = 95% confident. But you cannot say that the probability that μ is in that circle is 95%, because μ is not random. This throw might have been one of the 5% of throws that miss μ by more than 6 inches.

Lets assume we want a 95% CI for μ from a normal population with a known standard deviation σ , so the margin of error is:

$$M = 1.96 \frac{\sigma}{\sqrt{n}}$$

Then \bar{X} is the dart we are throwing at μ .

Before you take the sample and compute the mean, you have:

$$P(\bar{X} - M < \mu < \bar{X} + M) = 95\%$$

This is correct because \bar{X} is a random variable. However, once you compute the mean \bar{x} (lowercase x meaning it is now just a number, not a random quantity), the inequality:

$$\bar{x} - M < \mu < \bar{x} + M$$

is either true or false; the dart has landed at \bar{x} , and we don't know if this was one of the throws that is within M of μ .

Classical Confidence Interval

A classical confidence interval contains all values for which the data do not reject the null hypothesis that the parameter is equal to that value.

This does not necessarily tell you anything regarding the probability that the parameter is in the interval.

If however you intended to take a sample of the data and draw a 92% confidence interval, there would be a 92% probability of the population mean being within that interval, if however you drew that sample and created that interval the the probability of the invisible point μ being within that interval can't really be known because we just don't know where it is relative to the dart (i.e. how well the sample reflects the population).

Weekly Material

(3); Comparison of Population Samples

wk3

Lecture

1. Boxplots The delimiting marks in box plots correspond to the median and interquartile range (which is basically the median of all data below the median):

```
1 library("ggplot2")
2 ggplot(iris, aes(x = Sepal.Width, y = Sepal.Length, color =
  ↳ Species)) +
3 geom_boxplot() +
4 theme_bw()
```

there is no package called 'ggplot2'

Figure 2: Bar Plots Generated in ggplot2

2. Lecture Announcements Everything is online now. We'll be using *Zoom* a lot.
 - (a) **DONE** Finish Quiz 1
 - (b) **TODO** Finish Quiz 1 30 minutes to finish it, Test your computer First.
 - (c) Post pacman on the mailing list.
3. Naming Variables Attribute ... Data Base
4. **DONE** Review Chi Distribution,
 - (a) Is it in VNote?
 - (b) Should I put it in org-mode?

5. **DONE** Fix YAML Headers in `rmd` to play ball with Notable

(a) **DONE** Post this use-case to Reddit

(b) **DONE** Fix YAMLTags and TagFilter and Post to Reddit BASH [The Bash Script is here](#)

i. **TODO** Post an easy way to use this to the mailing list [Here](#) is a start that I've made

ii. Should I have all the lists and shit in `/tmp`

Pros	Cons
Less Mess	Harder to Directly watch what's happening
Easier to manage ⁵	

should I use `/tmp` or `/tmp/Notes` or somting?

iii. **DONE** Is there an easy way to pass the md off to vnote Or should I just use the `ln -s ~/Notes/DataSci /tmp/notes` trick?

Follow the instructions [here](#), it has to be done manually and then symlinked SCHEDULED: <2020-03-21 Sat>

iv. Should I have all the lists and shit in `/tmp`

Pros	Cons
Less Mess	Harder to Directly watch what's happening
Easier to manage ⁵	

should I use `/tmp` or `/tmp/Notes` or somting?

(c) **DONE** Is there a way to fix the Text Size of Code in emacs when I zoom out? Yeah just disable `M-x mixed-pitch-mode`

6. Calculating mean

```
1 library(tidyverse)
2 bwt <- c(3429, 3229, 3657, 3514, 3086, 3886)
3 (bwt <- sort(bwt))
4 mean(bwt)
5 mean(c(3429, 3514))
6 median(bwt)
7 max(bwt)-min(bwt)
```

The mean value is nice in that it has good mathematical properties, so for predictions and classifications (like gradient descent), if the model contains the mean the model will be smooth and the mean will lead to a well behaved model with respect to the derivative.

The Median value, however is more immune to large outliers, for example:

⁵By which I mean I'm not sure if the directory that the `00tagmatch` directory `00taglist` file will be made in are the wd of bash, or, if they are the location of `~/Notes`

I'm also not sure how that will be influenced by looking for `#tags` in the `~/Notes/DataSci` Directory

```

1 library(tidyverse)
2 x <- c(rnorm(10), 9) * 10 %>% round(1)
3 mean(x); median(x)

```

7. Calculating Range

```

1 range(bwt)
2 bwt %>% range %>% diff

```

8. Calculating Variance

```

1 (var <- (bwt-mean(bwt))^2 %>% mean)
2 var(bwt)
3 (sd <- (bwt-mean(bwt))^2 %>% %>% sqrt) # Not using n-1 !!
4 (sd <- sqrt(sum((bwt-mean(bwt))^2)/(length(bwt) -1)))
5 sd(bwt)
6 mean(sum((bwt-mean(bwt))^2))

```

9. InterQuartile Data

Tutorial

The tutorial work is located at ~/Notes/DataSci/ThinkingAboutData and linked here:

- PDF
- [HTML](#)
- [MD](#)
- [RMD](#)

DONE (4); Using Student's t-Distribution

wk4

The tutorial work is located at ~/Notes/DataSci/ThinkingAboutData and linked here:

- PDF
- [HTML](#)
- [MD](#)
- [RMD](#)

DONE (5) Discrete Distributions (Mapping Disease)

wk5

- Lecture
- Practical
 - [RMD File](#)

DONE Quizz for t Distribution Material

[The Quiz has been Released](#)

- [03 Tutorial](#)

[Tutorial 3](#)

- [04 Tutorial](#)
- Lecture 03
- Lecture 04

Lecture

The Poisson Model is the Binomial Model stretch towards its limits.

1. Combinatorics

The [Counting Formulas](#) are:

	selection	ordered	unordered
With Repetition		n^m	$\binom{m+n-1}{n}$
Without Repetition		$n_{(m)}$	$\binom{n}{m}$

Where:

- $\binom{n}{m} = \frac{n!}{m!(n-m)!}$
- $n_{(m)} = \frac{n!}{(n-m)!}$
- $n! = n \times (n-1) \times (n-2) \times \dots \times 1$

2. Binomial Distribution

A Binomial experiment requires the following conditions:

- We have n independent events.
- Each event has the same probability p of success.
- We are interested in the number of successes from the n trials (referred to as size in **R**).
- The probability of k successes from the n trials is a Binomial distribution with

probabilities:

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

- (a) Problem Hard drives have an annual failure rate of 10%, what is the probability of 2 hard drives failing after 3 years?

This means that:

- The number of repetitions is 3 ($n = \text{size} = 3$)
- The statistic we are interested in is 2 ($k = x = 2$)

```
1      dbinom(x = 2, size = 3, prob = 0.1)
2
3      ## For Hard Drives
4          ## k/x....is the number of years
5          ## n/ ....size is the number of failures
6          ## p ....is the probability of failure
7
8          ## choose(n,k)*p^k*(1-p)^(n-k)
9          ## dbinom(x = 2, size = 4, prob = p)
```

So in this case there would only a 2% chance.

3. Poisson An interesting thing with the poisson distribution is that the mean value and the variance are both equal.

The expected value is the limit that the mean value would approach if the sample was made arbitrarily large, the value is denoted λ

Poisson is French for fish so sometimes people call the distribution the *fishribution*.

4. Binomial and Poisson

The Poisson distribution is derived from the Binomial distribution.

If $(1 - p)$ is close to 1 then $np(1 - p) \approx np$, so for very large sample sizes we have the expected value equal to the variance, and a *Poisson* distribution.

This has something to do with widely increasing the number of trials, like say if we had an infinite number of trials with the probability of success in a given hour as 30%.

For a binomial distribution there are a set number of trials, let's say 8 trials with a 20% probability of Success:

1	2	3	4	5	6	7	8
F	S	F	F	S	S	F	F

In this case there are 3 successes, so let's set $k = 3$ and instead however the region was divided into smaller spaces ($n \rightarrow \infty$):

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
X	S	X	X	S	S	F	F	X	S	X	X	S	S	F	F

If we kept dividing this up we would have

The poisson is the limit as we increase the number of trials but try to keep k constant??

5. Confidence Intervals

(a) Binomial Just use bootstrapping, but also you can just use an approximate standardisation:

$$\begin{aligned}\sigma &= \sqrt{p(1-p)} \\ \sigma_{\bar{x}} &= \frac{\sigma}{\sqrt{n}} \\ &= \frac{\sqrt{p(1-p)}}{\sqrt{n}} \\ \Rightarrow p - z_{0.025} \times \sigma_{\bar{x}} &< p < p + z_{0.975} \times \sigma_{\bar{x}} \\ \Rightarrow p - 1.96\sigma_{\bar{x}} &< p < p + 1.96\sigma_{\bar{x}}\end{aligned}$$

So an approximate 95% confidence interval could be

This is the normal data, but remember that this is estimating binomial by standard normal and so will only be good for large values of n because binomial is discrete by nature.

See [the Correlation Notes](#) and [Khan Academy](#) and this section section []]

This means that for any sample drawn from the population, the true population value would be found within this interval for 0.95 of those samples

(b) Poisson This can also be done with Poisson by bootstrapping or using the same trick of λ as the variance:

$$\begin{aligned}\sigma &= \sqrt{\lambda} \\ \sigma_{\bar{x}} &= \frac{\sigma}{\sqrt{n}} \\ &= \frac{\sqrt{\lambda}}{\sqrt{n}} \\ \Rightarrow \lambda - z_{0.025} \times \sigma_{\bar{x}} &< \lambda < \lambda + z_{0.975} \times \sigma_{\bar{x}} \\ \lambda - 1.96\sigma_{\bar{x}} &< \lambda < \lambda + 1.96\sigma_{\bar{x}} \\ p - 1.96\sqrt{\frac{\lambda}{n}} &< p < p + 1.96\sqrt{\frac{\lambda}{n}}\end{aligned}$$

6. Summary

- Binomial for independent trials
 - mean: np
 - variance: $np(1-p)$
 - * standard error roughly is $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$
- Poisson for number of events in a given period
 - mean λ
 - variance λ
 - Standard error roughly is $\sqrt{\frac{\hat{\lambda}}{n}}$
- Choropleth maps are useful for visualising changes over area.

Tutorial

- 05_RMD File
- 05_MD File
- 05_PDF File
- 05_HTML File

DONE (6) Paired t-test (Observation or Experiment)

wk6

- 06 Lecture Notes
- 06 Practical
 - RMD File

TODO (7) Corellation (Do Taller People Earn More)

wk7

- 07 Lecture Notes
- 07 Practical
 - RMD File

DONE Lecture

In the past we did categorical and categorical-continuous.

Now we're doing purely continuous.

1. **DONE** How to Derive the Correlation Coefficient Refer to this paper The Correlation coefficient can be interpreted in one of two ways:

- The covariance scaled relative to the x and y variance
 - $\rho = \frac{S_{x,y}}{s_x \cdot s_y}$
- The rate of change of the line of best fit of the standardised data
 - This is equivalent to the rate of change of the line of best fit divided by $(\frac{s_y}{s_x})$:
 - * $\rho = b \cdot \frac{s_x}{s_y} \iff \hat{y}_i = bx_i + c$

This can be seen by performing linear regression in **R**:

```

1  head(cars)
2  cars_std <- as.data.frame(scale(cars))
3  y <- cars$dist
4  x <- cars$speed
5
6  ## Correlation Coefficient
7  cor(x = cars$speed, y = cars$dist)
8
9  ## Covariance
10 cov(x = cars$speed, y = cars$dist)/sd(cars$speed)/sd(cars$dist)
11
12 ## Standardised Rate of Change
13 lm(dist ~ speed, data = cars_std)$coefficients[2]
14
15 ### Using Standardised Rate of change
16 lm(dist ~ speed, data = cars)$coefficients[2] / (sd(y)/sd(x))

```

```

      speed dist
1         4    2
2         4   10
3         7    4
4         7   22
5         8   16
6         9   10

```

```
[1] 0.8068949
```

```
[1] 0.8068949
```

```

      speed
0.8068949

```

```

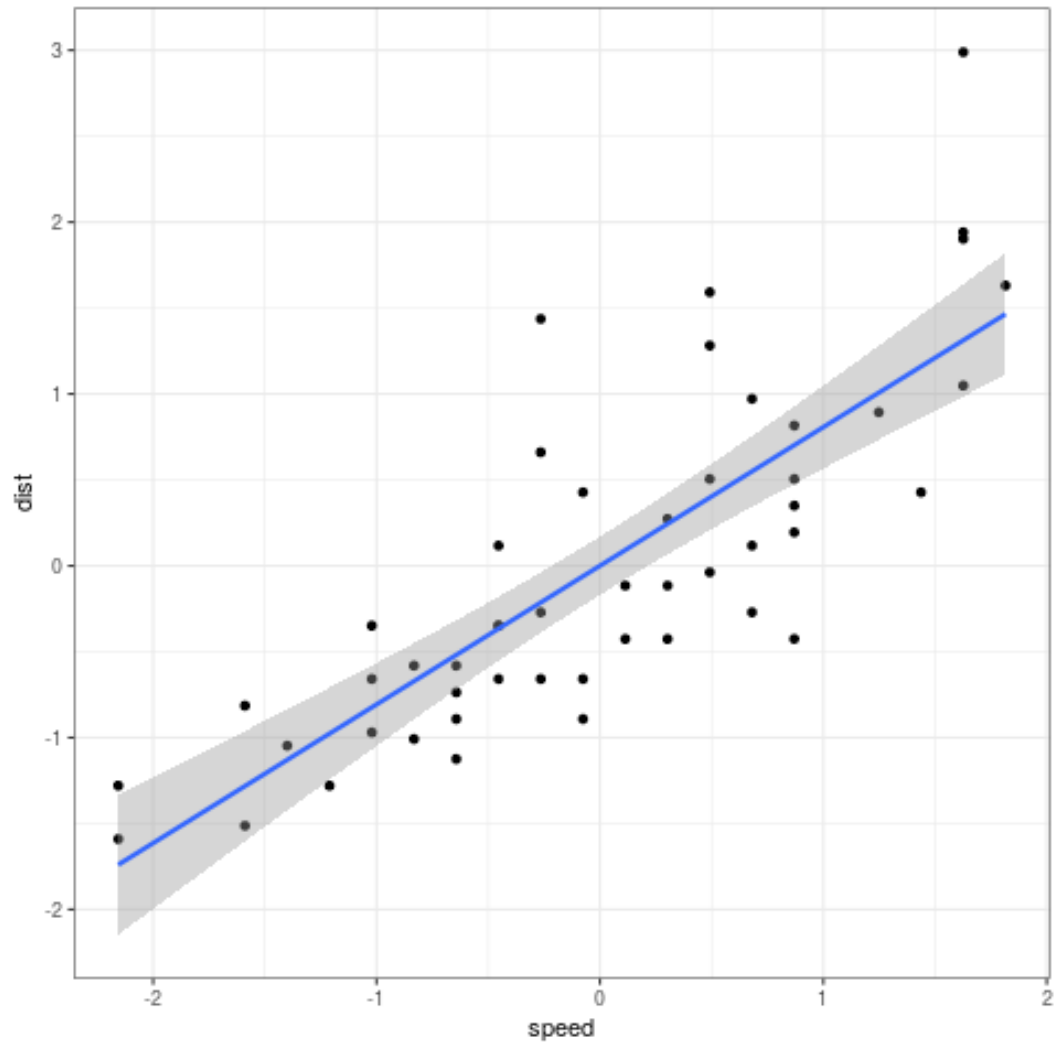
      speed
0.8068949

```

```

1  #+BEGIN_SRC R :cache yes :exports both :results output graphics
   ↪ :file ./test.png
2  library("ggplot2")
3  cars_std <- as.data.frame(scale(cars))
4
5  ggplot(cars_std, aes(x = speed, y = dist)) +
6  geom_point() +
7  geom_smooth(method = "lm") +
8  theme_bw()

```



This shows that despite noise data can still have a correlation coefficient of 1 if the noise is evenly distributed in a way that the correlation coefficient can have a rate of change of 1.

2. **TODO** Prove the Correlation Coefficient and email Laurence
3. Bootstrapping The big assumption with bootstrapping is that the population can be seen as equiv to an infinite repetition of the sample size.
 So assume that a population that is an infinite repetition of the sample, then take a sample of that infinite population and you have a bootstrap.
 So we could either create a population from the sample with a size of ∞ , which might be difficult, or, we could instead just resample the observation for each replication.
4. Confidence Intervals


```

1  load("Notes/DataSci/ThinkingAboutData/TAD.rdata ")
2  r = cor(crabsmolt$postsz, crabsmolt$presz)
3  a <- crabsmolt
4  N <- nrow(crabsmolt)
5  pos <- sample(N, size = N, replace = TRUE)
6  aboot <- a[pos,]
7
8  cor(aboot$postsz, aboot$presz)
9
10
11 # replicate(10^4, {})

```

- Attached RMD

5. **TODO** Questions [In this part](#) would it simply be equivalent to take the mean of all observations?

TODO (8) No Really do they Earn More?

wk8

TODO Lecture

TODO Tutorial

- [RMD File](#)

wk9 Break

wk9

TODO (9) Do redheads have a lower pain threshold?

wk10

TODO (10) What is Normal?

wk11

TODO (11) Normality as opposed to deviant etc.

wk12

TODO (12) When it all goes Wrong

wk13

TODO (13) Exam Prep

wk14

Symlinks to Material

Lecture

- Lecture 04
- Lecture 05
- Lecture 06
- Lecture 07
- Lecture 08
- Lecture 09
- Lecture 10
- Lecture 11
- Lecture 12
- Lecture 13
- Lecture 14

Tutorial

- Worksheet 01
- Worksheet 02
- Worksheet 03
- Worksheet 04
- Worksheet 05
- Worksheet 06
- Worksheet 07
- Worksheet 08
- Worksheet 09
- Worksheet 10
- Worksheet 11
- Worksheet 12
- Worksheet 13
- Worksheet 14

Central Limit Theorem

The central Limit theorem provides us the sampling distribution of \bar{X} even when we don't know what the original population of X looks like:

1. If the population is normal, the sample mean of that population will be

normally distributed, $\bar{X} \sim \mathcal{N}\left(\mu, \left(\frac{\sigma}{\sqrt{n}}\right)\right)$

1. As sample size n increases, the distribution of sample means converges to the population mean μ

- i.e. the *standard error of the mean*

$\sigma_{\bar{x}} = \left(\frac{\sigma}{\sqrt{n}}\right)$ will become smaller

1. If the sample size (of sample means) is large enough ($n \geq 30$) the sample means will be normally distributed even if the original population is non-normal