

Uncertainty: (how to) deal with it

Peter Ralph

1 October 2020 – Advanced Biological Statistics

1

Course overview

2 . 1



image: Frank Klausz, woodandshop.com

Steps in data analysis

1. Care, or at least think, about the data.
2. Look at the data.
3. Query the data.
4. Check the results.
5. Communicate.

Steps in data analysis

1. Care, or at least think, about the data.
2. Look at the data.
3. Query the data.
4. Check the results.
5. Communicate.

Often “statistics” focuses on *querying*. Doing that effectively requires all the other steps, too.

Prerequisites

We'll be assuming that you have some familiarity with

- programming, and
- statistics

2 / 4

Prerequisites

We'll be assuming that you have some familiarity with

- programming, and
- statistics

For instance, you should be able to figure out what this means:

```
x = c(2, 4, 3, 6)
y = c(5, 12, 4, 10, 2)
t.test(x, y)
```

```
##
## Welch Two Sample t-test
##
## data: x and y
## t = -1.3761, df = 5.4988, p-value = 0.2222
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -8.031728 2.331728
## sample estimates:
## mean of x mean of y
## 3.75 6.60
```

2 / 4

Overview and mechanics

See [the course website](#).

2 / 5

Break

3 . 1

Please take 10 minutes to

1. answer the “Welcome Survey” on Canvas,
2. [get the course repository from github](#),
3. [install Rstudio](#) and/or
4. move around.

3 . 2

Questions?

3 . 3

Some core statistical concepts

4 - 1

Statistics or parameters?

A statistic is

a numerical description of a dataset.

4 - 2

Statistics or parameters?

A statistic is

a numerical description of a dataset.

A parameter is

a numerical attribute of a model of reality.

4 - 2

Statistics or parameters?

A statistic is

a numerical description of a dataset.

A parameter is

a numerical attribute of a model of reality.

Often, *statistics* are used to estimate *parameters*.

4 - 2

The two heads of classical statistics

estimating parameters, with uncertainty (*confidence intervals*)

evaluating (in-)consistency with a particular situation (*p-values*)

4 - 3

The two heads of classical statistics

estimating parameters, with uncertainty (*confidence intervals*)

evaluating (in-)consistency with a particular situation (*p-values*)

1. What do these data tell us about the world?
2. How strongly do we believe it?

4 - 3

The two heads of classical statistics

estimating parameters, with uncertainty (*confidence intervals*)

evaluating (in-)consistency with a particular situation (*p-values*)

1. What do these data tell us about the world?

2. How strongly do we believe it?

This week: digging in, with simple examples.

4 . 3

Lurking, behind everything:

is *uncertainty*

4 . 4

Lurking, behind everything:

is *uncertainty*

thanks to *randomness*.

4 . 4

Lurking, behind everything:

is *uncertainty*

thanks to *randomness*.

How do we understand randomness, concretely and quantitatively?

4 . 4

Lurking, behind everything:

is *uncertainty*

thanks to *randomness*.

How do we understand randomness, concretely and quantitatively?

With *models*.

4 . 4

A quick look at some data

5 . 1

Some data

AirBnB hosts in Portland, OR: [website](#) and [download link](#).

```
airbnb <- read.csv("../Datasets/portland-airbnb-listings.csv")
nrow(airbnb)

## [1] 5634

names(airbnb)

## [1] "id" "listing_url"
## [5] "name" "summary"
## [9] "experiences_offered" "neighborhood_overview"
## [13] "access" "interaction"
## [17] "medium_url" "picture_url"
## [21] "host_url" "host_name"
## [25] "host_about" "host_response_time"
## [29] "host_is_superhost" "host_thumbnail_url"
## [33] "host_listings_count" "host_total_listings_count"
## [37] "host_identity_verified" "street"
## [41] "neighbourhood_group_cleansed" "city"
## [45] "market" "smart_location"
## [49] "latitude" "longitude"
## [53] "room_type" "accommodates"
## [57] "beds" "bed_type"
## [61] "price" "weekly_price"
## [65] "cleaning_fee" "guests_included"
## [69] "maximum_nights" "minimum_minimum_nights"
## [73] "maximum_maximum_nights" "minimum_nights_avg_ntm"
```

```
## [77] "has_availability" "availability_30"
## [81] "availability_365" "calendar_last_scraped"
```

Questions: how much does an AirBnB typically cost in Portland? Do “instant bookable” ones cost more?

Second, look at the data

```
summary(airbnb$price)
```

```
##      Length      Class      Mode  
##      5634 character character
```

5 . 4

```
summary(airbnb$instant_bookable)
```

```
##      Length      Class      Mode  
##      5634 character character
```

5 . 5

Whoops

```
airbnb$price <- as.numeric(gsub("$", "", airbnb$price, fixed=TRUE))
```

```
## Warning: NAs introduced by coercion
```

```
airbnb$instant_bookable <- (airbnb$instant_bookable == "t")
```

5 . 6

```
summary(airbnb$price)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.   NA's  
##      0.0   69.0   95.0   119.5   136.0   999.0    32
```

```
summary(airbnb$instant_bookable)
```

```
##      Mode FALSE  TRUE  
## logical  2960  2674
```

5 / 7

```
table(airbnb$bed_type) # fm
```

```
##  
##      Airbed      Couch      Futon Pull-out Sofa      Real Bed  
##         9         7         42         21         5555
```

5 / 8

How much is a typical night?

```
mean(airbnb$price, na.rm=TRUE)
```

```
## [1] 119.5396
```

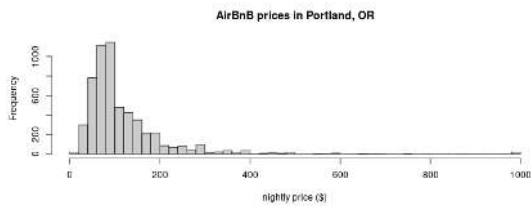
5 / 9

```
hist(airbnb$price, breaks=40, xlab='nightly price ($)', col=grey(.8),
     xlim=range(airbnb$price, finite=TRUE), main='AirBnB prices in Portland, OR')
```



5 / 10

```
hist(airbnb$price, breaks=40, xlab='nightly price ($)', col=grey(.8),
     xlim=range(airbnb$price, finite=TRUE), main='AirBnB prices in Portland, OR')
```

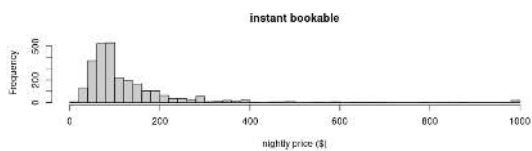
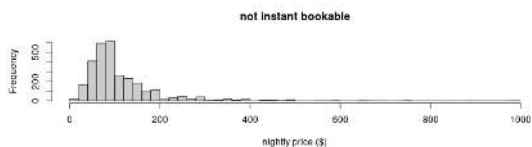


Conclusion?

5 / 10

Do "instant bookable" charge more?

```
layout(1:2)
instant <- airbnb$price[airbnb$instant_bookable]
not_instant <- airbnb$price[!airbnb$instant_bookable]
hist(not_instant, breaks=40, xlab='nightly price ($)', col=grey(.8),
     xlim=range(airbnb$price, finite=TRUE), main='not instant bookable')
hist(instant, breaks=40, xlab='nightly price ($)', col=grey(.8), main='instant
bookable')
```



5 / 11

```
(tt <- t.test(instant, not_instant))
```

```
##  
## Welch Two Sample t-test  
##  
## data: instant and not_instant  
## t = 3.6482, df = 5039.8, p-value = 0.0002667  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
## 4.475555 14.872518  
## sample estimates:  
## mean of x mean of y  
## 124.6409 114.9668
```

5 / 12

Conclusion

*Instant bookable hosts cost more than others
($P=0.00027$, t-test with $df=5039.7695486$).*

5 / 13

Conclusion

*Instant bookable hosts cost more than others
($P=0.00027$, t-test with $df=5039.7695486$).*

Critique this conclusion, and write your own.

Scribe: person with the smallest sample.int(1000, 1).

5 / 12

Don't forget Steps 1 and 5!

1. Care, or at least think, about the data.
2. Communicate.

5 / 14

Don't forget Steps 1 and 5!

1. Care, or at least think, about the data.
2. Communicate.

How *big* is the difference? How sure are we?

5 / 14

Don't forget Steps 1 and 5!

1. Care, or at least think, about the data.
2. Communicate.

How *big* is the difference? How sure are we?

Statistical significance does not imply real-world significance.

5 / 14

Revised conclusion (in class)

Instant bookable hosts cost on average \$10 more than not instant bookable, with a 95% confidence interval of \$4.50 to \$15. The distribution of prices in the two groups were very similar: for instance, the first and third quantiles of instant bookable hosts are \$70 and \$145, and those of not instant bookable hosts are \$68 and \$130, respectively. The average instant bookable cost was about \$125, with a 95% confidence interval of +/- about \$4; non-instant bookable hosts cost on average \$115 per night, with a 95% CI of about +/- \$3. Note that the difference of \$10 is smallish compared to the price of a room, but the difference is highly significant ($p=.0003$, t-test with 5039 degrees of freedom) because of the large sample sizes.

5 . 15

So: what did we just do?

5 . 16

Hypothesis testing and p -values

6 . 1

A p -value is

6 - 2

A p -value is

the probability of seeing a result at least as surprising as what was observed in the data, if the null hypothesis is true.

6 - 2

A p -value is

the probability of seeing a result at least as surprising as what was observed in the data, if the null hypothesis is true.

Usually, this means

- *a result* - numerical value of a statistic
- *surprising* - big
- *null hypothesis* - the model we use to calculate the p -value

which can all be defined to suit the situation.

6 - 2

What does a small p -value mean?

If the null hypothesis were true, then you'd be really unlikely to see something like what you actually did.

6 . 3

What does a small p -value mean?

If the null hypothesis were true, then you'd be really unlikely to see something like what you actually did.

So, either the “null hypothesis” is not a good description of reality or something surprising happened.

6 . 3

What does a small p -value mean?

If the null hypothesis were true, then you'd be really unlikely to see something like what you actually did.

So, either the “null hypothesis” is not a good description of reality or something surprising happened.

How useful this is depends on the null hypothesis.

6 . 3

For instance

```
##  
## Welch Two Sample t-test  
##  
## data: instant and not_instant  
## t = 3.6482, df = 5039.8, p-value = 0.0002667  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
## 4.475555 14.872518  
## sample estimates:  
## mean of x mean of y  
## 124.6409 114.9668
```

6 . 4

Also for instance

```
t.test(airbnb$price)
```

```
##  
## One Sample t-test  
##  
## data: airbnb$price  
## t = 91.32, df = 5601, p-value < 2.2e-16  
## alternative hypothesis: true mean is not equal to 0  
## 95 percent confidence interval:  
## 116.9734 122.1058  
## sample estimates:  
## mean of x  
## 119.5396
```

6 . 5

Also for instance

```
t.test(airbnb$price)
```

```
##  
## One Sample t-test  
##  
## data: airbnb$price  
## t = 91.32, df = 5601, p-value < 2.2e-16  
## alternative hypothesis: true mean is not equal to 0  
## 95 percent confidence interval:  
## 116.9734 122.1058  
## sample estimates:  
## mean of x  
## 119.5396
```

Is that p -value useful?

6 . 5

Exercise:

My hypothesis: People tend to have longer index fingers on the hand they write with because writing stretches the ligaments.

(class survey) How many people have a longer index finger on the hand they write with?

6 / 6

Exercise:

My hypothesis: People tend to have longer index fingers on the hand they write with because writing stretches the ligaments.

(class survey) How many people have a longer index finger on the hand they write with?

(class survey) Everyone flip a coin:

```
ifelse(runif(1) < 0.5, "H", "T")
```

and put the result in [this google doc](#).

6 / 6

Exercise:

My hypothesis: People tend to have longer index fingers on the hand they write with because writing stretches the ligaments.

(class survey) How many people have a longer index finger on the hand they write with?

(class survey) Everyone flip a coin:

```
ifelse(runif(1) < 0.5, "H", "T")
```

and put the result in [this google doc](#).

We want to estimate the parameter

$$\theta = \mathbb{P}(\text{random person has writing finger longer}),$$

and now we have a *fake dataset* with $\theta = 1/2$.

6 / 6

Let's get some more data:

```
n <- 37 # class size
sum(ifelse(runif(1) < 1/2, "H", "T") == "H")
```

and put the result in [the same google doc](#).

6 . 7

Let's get some more data:

```
n <- 37 # class size
sum(ifelse(runif(1) < 1/2, "H", "T") == "H")
```

and put the result in [the same google doc](#).

Now we can estimate the p -value for the hypothesis that $\theta = 1/2$.

6 . 7

A faster method:

```
replicate(1000, sum(rbinom(n, 1, 1/2) > 0))
```

6 . 8

A faster method:

```
replicate(1000, sum(rbinom(n, 1, 1/2) > 0))
```

or, equivalently,

```
rbinom(1000, n, 1/2)
```

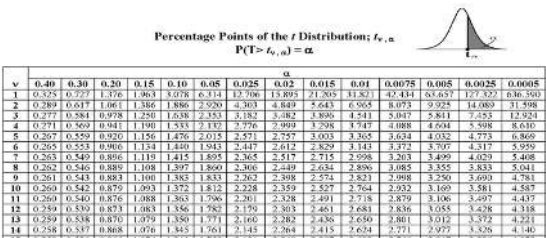
(in class)

```
mean(replicate(10000, rbinom(1, 36, 1/2) >= 20))  
  
## [1] 0.3057
```

Here, we’ve estimated that the difference in numbers of people with a longer finger on each hand is not statistically significant ($p \approx 0.3$, by simulation).

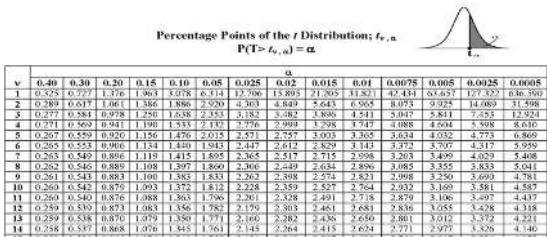
So, where do p -values come from?

Either math:



So, where do *p*-values come from?

Either math:



Or, computers. (maybe math, maybe simulation, maybe both)

So, where did *this* *p*-value come from?

```
(tt <- t.test(instant, not_instant))

##
## Welch Two Sample t-test
##
## data: instant and not_instant
## t = 3.6482, df = 5039.8, p-value = 0.0002667
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  4.475555 14.872518
## sample estimates:
## mean of x mean of y
## 124.6409 114.9668
```

So, where did *this* *p*-value come from?

```
(tt <- t.test(instant, not_instant))

##
## Welch Two Sample t-test
##
## data: instant and not_instant
## t = 3.6482, df = 5039.8, p-value = 0.0002667
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  4.475555 14.872518
## sample estimates:
## mean of x mean of y
## 124.6409 114.9668
```

The *t* distribution! (see separate slides)

