

A Very Incomplete Survey of Basic Statistical Tests in R

EDP 619 Week 5

Dr. Abhik Roy

Getting Ready

To follow along, please make sure to do the following

1. Run Basic Statistical Tests install.R
2. Open up a blank .R script
3. Load up the following packages

```
library(tidyverse)
library(car)
library(foreign)
library(lme4)
library(MASS)
library(CCA)
library(psych)
```

Purpose

This is not a traditional presentation per say, rather it is a collection of common statistical tests used for survey analyses that you may need to use at some point. Review the content simply to see what is included and then come back as needed.

The examples are of course presented with R in mind, but obviously not restricted to the platform.

The content is intended to cast a wide net so some may look familiar while others will not

You can think of this as a toolkit

Note

Here are some things to consider as you go through this

- Items in here will apply depending on your focus and strengths
- Never memorize formulas - that's what the Internet is for. Rather it is important to know what test applies to which scenario
- Since this is not a statistics course, only test names and examples are provided
- Some of the syntax may be structured in an odd way. This is only done so so that they be fit within the slide. Remember in general spacing doesn't mean much in R
- These are presented in alphabetical order according to the standard name of the test

Decisions Decisions Decisions

When deciding which test is appropriate to use, it is important to consider the type of variables that you have. Please load in the following data sets (and look at them by using `View()` or `head()`)

```
some_ed_data <-  
  read_csv("some_ed_data.csv")
```

```
some_exercise_data <-  
  read_csv("some_exercise_data.csv")
```

```
some_survey_data <-  
  read_csv("some_survey_data.csv")
```

The Datas (1/3)

```
some_ed_data
```

```
## # A tibble: 200 × 11
##       id female race   ses schtyp  prog  read write  math science socst
##   <dbl>  <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
## 1    70     0     4     1     1     1    57    52    41    47    57
## 2   121     1     4     2     1     3    68    59    53    63    61
## 3    86     0     4     3     1     1    44    33    54    58    31
## 4   141     0     4     3     1     3    63    44    47    53    56
## 5   172     0     4     2     1     2    47    52    57    53    61
## 6   113     0     4     2     1     2    44    52    51    63    61
## 7    50     0     3     2     1     1    50    59    42    53    61
## 8    11     0     1     2     1     2    34    46    45    39    36
## 9    84     0     4     2     1     1    63    57    54    58    51
## 10   48     0     3     2     1     2    57    55    52    50    51
## # ... with 190 more rows
```

The Datas (2/3)

```
some_exercise_data
```

```
## # A tibble: 90 × 6
##       id   diet exertype pulse  time highpulse
##   <dbl> <dbl>    <dbl> <dbl> <dbl>     <dbl>
## 1     1     1        1    85     1        0
## 2     1     1        1    85     2        0
## 3     1     1        1    88     3        0
## 4     2     1        1    90     1        0
## 5     2     1        1    92     2        0
## 6     2     1        1    93     3        0
## 7     3     1        1    97     1        0
## 8     3     1        1    97     2        0
## 9     3     1        1    94     3        0
## 10    4     1        1    80     1        0
## # ... with 80 more rows
```

The Datas (3/3)

```
some_survey_data
```

```
## # A tibble: 20 × 4
##   respondent gender item_1 item_2
##       <dbl>    <dbl>    <dbl>    <dbl>
## 1           1        1        1        5
## 2           2        2        2        4
## 3           3        3        1        2
## 4           4        4        2        4
## 5           5        5        2        2
## 6           6        6        1        3
## 7           7        7        1        4
## 8           8        8        2        5
## 9           9        9        2        3
## 10          10       10        2        3
## 11          11       11        1        2
## 12          12       12        2        3
## 13          13       13        2        1
## 14          14       14        1        4
## 15          15       15        2        2
## 16          16       16        1        2
## 17          17       17        1        3
## 18          18       18        1        5
## 19          19       19        1        3
## 20          20       20        2        5
```

An Incomplete Table of Approaches (1/4)

Number of Dependent Variables	Number and Type of Independent Variables	Type of Dependent Variables	Test(s)
1	0 IVs (1 population)	interval & normal	one-sample t-test
1	0 IVs (1 population)	ordinal or interval	one-sample median
1	0 IVs (1 population)	categorical (2 categories)	binomial test
1	0 IVs (1 population)	categorical	Chi-square goodness-of-fit
1	1 IV with 2 levels (independent groups)	interval & normal	2 independent sample t-test
1	1 IV with 2 levels (independent groups)	ordinal or interval	Wilcoxon-Mann Whitney test
1	1 IV with 2 levels (independent groups)	categorical	Chi-square test
1	1 IV with 2 levels (independent groups)	categorical	Fisher's exact test

An Incomplete Table of Approaches (2/4)

Number of Dependent Variables	Number and Type of Independent Variables	Type of Dependent Variables	Test(s)
1	1 IV with 2 or more levels (independent groups)	interval & normal	one-way ANOVA
1	1 IV with 2 or more levels (independent groups)	ordinal or interval	Kruskal Wallis
1	1 IV with 2 or more levels (independent groups)	categorical	Chi-square test
1	1 IV with 2 levels (dependent/matched groups)	interval & normal	paired t-test
1	1 IV with 2 levels (dependent/matched groups)	ordinal or interval	Wilcoxon signed ranks test
1	1 IV with 2 levels (dependent/matched groups)	categorical	McNemar
1	1 IV with 2 or more levels (dependent/matched groups)	interval & normal	one-way repeated measures ANOVA
1	1 IV with 2 or more levels (dependent/matched groups)	ordinal or interval	Friedman test
1	1 IV with 2 or more levels (dependent/matched groups)	categorical (2 categories)	repeated measures logistic regression

An Incomplete Table of Approaches (3/4)

Number of Dependent Variables	Number and Type of Independent Variables	Type of Dependent Variables	Test(s)
1	2 or more IVs (independent groups)	interval & normal	factorial ANOVA
1	2 or more IVs (independent groups)	ordinal or interval	ordered logistic regression
1	2 or more IVs (independent groups)	categorical (2 categories)	factorial logistic regression
1	1 interval IV	interval & normal	correlation
1	1 interval IV	interval & normal	simple linear regression
1	1 interval IV	ordinal or interval	non-parametric correlation
1	1 interval IV	categorical	simple logistic regression

An Incomplete Table of Approaches (4/4)

Number of Dependent Variables	Number and Type of Independent Variables	Type of Dependent Variables	Test(s)
1	1 or more interval IVs and/or 1 or more categorical IVs	interval & normal	multiple regression
1	1 or more interval IVs and/or 1 or more categorical IVs	interval & normal	analysis of covariance
1	1 or more interval IVs and/or 1 or more categorical IVs	categorical	multiple logistic regression
1	1 or more interval IVs and/or 1 or more categorical IVs	categorical	discriminant analysis
2+	1 IV with 2 or more levels (independent groups)	interval & normal	one-way MANOVA
2+	2+	interval & normal	multivariate multiple linear regression
2+	0	interval & normal	factor analysis
2 sets of 2+	0	interval & normal	canonical correlation

Tests

ANOVA

```
summary(aov(some_ed_data$write ~ some_ed_data$prog + some_ed_data$read))

##          Df Sum Sq Mean Sq F value    Pr(>F)
## some_ed_data$prog     1    586     586   10.2 0.00164 **
## some_ed_data$read     1   5965    5965  103.7 < 2e-16 ***
## Residuals        197 11327      57
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Binomial Test

```
prop.test(sum(some_ed_data$female), length(some_ed_data$female), p = 0.5)

##
##      1-sample proportions test with continuity correction
##
## data: sum(some_ed_data$female) out of length(some_ed_data$female), null probability 0.5
## X-squared = 1.445, df = 1, p-value = 0.2293
## alternative hypothesis: true p is not equal to 0.5
## 95 percent confidence interval:
##  0.4733037 0.6149394
## sample estimates:
##      p
## 0.545
```

Canonical Correlation

```
cc(cbind(some_ed_data$read, some_ed_data$write),  
  cbind(some_ed_data$math, some_ed_data$science))  
## $cor  
## [1] 0.7728409 0.0234784  
##  
## $names  
## $names$Xnames  
## NULL  
##  
## $names$Ynames  
## NULL  
##  
## $names$ind.names  
## NULL  
##  
##  
## $xcoef  
## [,1] [,2]  
## [1,] -0.06326131 -0.1037908  
## [2,] -0.04924918  0.1219084  
##  
## $ycoef  
## [,1] [,2]  
## [1,] -0.06698268  0.1201425  
## [2,] -0.04824063 -0.1208860  
##  
## $scores
```

Chi-square Test

```
chisq.test(table(some_ed_data$female, some_ed_data$schtyp))

## Pearson's Chi-squared test with Yates' continuity correction
## data: table(some_ed_data$female, some_ed_data$schtyp)
## X-squared = 0.00054009, df = 1, p-value = 0.9815
```

Chi-square Goodness of Fit

```
chisq.test(table(some_ed_data$race), p = c(10, 10, 10, 70)/100)

## Chi-squared test for given probabilities
## data: table(some_ed_data$race)
## X-squared = 5.0286, df = 3, p-value = 0.1697
```

Correlation

```
cor(some_ed_data$read, some_ed_data$write)

## [1] 0.5967765

cor.test(some_ed_data$read, some_ed_data$write)

##
##      Pearson's product-moment correlation
##
## data: some_ed_data$read and some_ed_data$write
## t = 10.465, df = 198, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  0.4993831 0.6792753
## sample estimates:
##      cor
## 0.5967765
```

Discriminant Analysis

```
lda(factor(some_ed_data$prog) ~ some_ed_data$read + some_ed_data$write + some_ed_data$math, data = some_ed_data)

## Call:
## lda(factor(some_ed_data$prog) ~ some_ed_data$read + some_ed_data$write +
##      some_ed_data$math, data = some_ed_data)
##
## Prior probabilities of groups:
##       1      2      3
## 0.225 0.525 0.250
##
## Group means:
##   some_ed_data$read some_ed_data$write some_ed_data$math
## 1      49.75556      51.33333     50.02222
## 2      56.16190      56.25714     56.73333
## 3      46.20000      46.76000     46.42000
##
## Coefficients of linear discriminants:
##                               LD1        LD2
## some_ed_data$read  0.02919876  0.04385321
## some_ed_data$write 0.03832289 -0.13698224
## some_ed_data$math  0.07034625  0.07931008
##
## Proportion of trace:
##      LD1      LD2
## 0.9874 0.0126
```

Factor Analysis

```
fa(r = cor(model.matrix(~read + write + math + science + socst - 1, data = some_ed_data)), rotate = "none", fm = "pa", 2)

## maximum iteration exceeded

## Factor Analysis using method = pa
## Call: fa(r = cor(model.matrix(~read + write + math + science + socst -
##     1, data = some_ed_data)), nfactors = 2, rotate = "none",
##     fm = "pa")
## Standardized loadings (pattern matrix) based upon correlation matrix
##          PA1    PA2    h2   com
## read      0.81  0.06  0.66  0.34  1.0
## write     0.76  0.00  0.58  0.42  1.0
## math      0.80  0.17  0.67  0.33  1.1
## science   0.75  0.26  0.62  0.38  1.2
## socst     0.79 -0.48  0.85  0.15  1.6
##
##          PA1    PA2
## SS loadings   3.06  0.33
## Proportion Var 0.61  0.07
## Cumulative Var 0.61  0.68
## Proportion Explained 0.90  0.10
## Cumulative Proportion 0.90  1.00
##
## Mean item complexity = 1.2
## Test of the hypothesis that 2 factors are sufficient.
```

Factorial ANOVA (Analysis of Variance)

```
anova(lm(write ~ female * ses, data = some_ed_data))

## Analysis of Variance Table
##
## Response: write
##           Df  Sum Sq Mean Sq F value    Pr(>F)
## female      1 1176.2 1176.21 14.7212 0.0001680 ***
## ses         1 1042.3 1042.32 13.0454 0.0003862 ***
## female:ses  1     0.0     0.04  0.0005 0.9827570
## Residuals 196 15660.3    79.90
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Factorial Logistic Regression

```
summary(glm(female ~ prog * schtyp, data = some_ed_data, family = binomial))

## 
## Call:
## glm(formula = female ~ prog * schtyp, family = binomial, data = some_ed_data)
## 
## Deviance Residuals:
##    Min      1Q  Median      3Q     Max 
## -1.698  -1.247   1.069   1.109   1.572 
## 
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)    
## (Intercept) -2.2765    1.8857  -1.207   0.227    
## prog         1.2303    0.9398   1.309   0.191    
## schtyp       2.2405    1.7017   1.317   0.188    
## prog:schtyp -1.1313    0.8622  -1.312   0.189    
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
## Null deviance: 275.64 on 199 degrees of freedom
## Residual deviance: 273.65 on 196 degrees of freedom
## AIC: 281.65
## 
## Number of Fisher Scoring iterations: 4
```

Friedman Test

```
friedman.test(cbind(some_ed_data$read, some_ed_data$write, some_ed_data$math))

## Friedman rank sum test
##
## data: cbind(some_ed_data$read, some_ed_data$write, some_ed_data$math)
## Friedman chi-squared = 0.64491, df = 2, p-value = 0.7244
```

Kruskal Wallis Test

```
kruskal.test(some_ed_data$write, some_ed_data$prog)

## Kruskal-Wallis rank sum test
##
## data: some_ed_data$write and some_ed_data$prog
## Kruskal-Wallis chi-squared = 34.045, df = 2, p-value = 4.047e-08
```

McNemar Test

```
# Some made up data in matrix form
made_up_matrixdata <- matrix(c(150, 22, 21, 12), 2, 2)

mcnemar.test(made_up_matrixdata)

##      McNemar's Chi-squared test with continuity correction
##
## data:  made_up_matrixdata
## McNemar's chi-squared = 0, df = 1, p-value = 1
```

Multiple Regression

```
lm(some_ed_data$write ~  
    some_ed_data$female + some_ed_data$read + some_ed_data$math + some_ed_data$science + some_ed_data$socst)  
  
##  
## Call:  
## lm(formula = some_ed_data$write ~ some_ed_data$female + some_ed_data$read +  
##       some_ed_data$math + some_ed_data$science + some_ed_data$socst)  
##  
## Coefficients:  
##             (Intercept)  some_ed_data$female      some_ed_data$read      some_ed_data$math  
##                 6.1388                  5.4925                  0.1254                  0.2381  
## some_ed_data$science      some_ed_data$socst  
##                 0.2419                  0.2293
```

Multivariate Multiple Regression

```
mmrlm <- lm(cbind(write, read) ~ female + math + science + socst, data = some_ed_data)

summary(Anova(mmrlm))

## 
## Type II MANOVA Tests:
## 
## Sum of squares and products for error:
##      write     read
## write 7258.783 1091.057
## read  1091.057 8699.762
## 
## -----
## 
## Term: female
## 
## Sum of squares and products for the hypothesis:
##      write     read
## write 1413.5284 -133.48461
## read  -133.4846  12.60544
## 
## Multivariate Tests: female
##                   Df test stat approx F num Df den Df    Pr(>F)
## Pillai          1 0.1698853 19.85132      2     194 1.4335e-08 ***
## Wilks           1 0.8301147 19.85132      2     194 1.4335e-08 ***
```

Non-parametric Correlation

```
cor.test(some_ed_data$read, some_ed_data$write, method = "spearman")

## Warning in cor.test.default(some_ed_data$read, some_ed_data$write, method = "spearman"):
## Cannot compute exact p-value with ties

##
##      Spearman's rank correlation rho
##
## data: some_ed_data$read and some_ed_data$write
## S = 510993, p-value < 2.2e-16
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
##      rho
## 0.6167455
```

One Sample *t*-test

```
t.test(some_ed_data$read, mu = 50)

##
##      One Sample t-test
##
## data: some_ed_data$read
## t = 3.0759, df = 199, p-value = 0.002394
## alternative hypothesis: true mean is not equal to 50
## 95 percent confidence interval:
##  50.80035 53.65965
## sample estimates:
## mean of x
##      52.23
```

One-way Analysis of Variance (ANOVA)

```
summary(aov(some_ed_data$read ~ some_ed_data$prog))

##           Df Sum Sq Mean Sq F value Pr(>F)
## some_ed_data$prog     1    381   381.1   3.674 0.0567 .
## Residuals       198  20538   103.7
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

One-way Multivariate Analysis of Variance (MANOVA)

```
summary(  
  manova(  
    cbind(some_ed_data$read, some_ed_data$write, some_ed_data$math) ~ some_ed_data$prog  
  )  
)  
  
##           Df Pillai approx F num Df den Df Pr(>F)  
## some_ed_data$prog  1 0.035319   2.392      3     196 0.06984 .  
## Residuals       198  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

One-way Repeated Measures Analysis of Variance (ANOVA)

```
model <- lm(gender ~ item_1 + item_2, data = some_survey_data)
analysis <- Anova(model, idata = factor_surveydata, idesign = ~s)
print(analysis)

## Anova Table (Type II tests)
##
## Response: gender
##           Sum Sq Df F value Pr(>F)
## item_1     0.0601  1  0.2396 0.6307
## item_2     0.7268  1  2.8974 0.1069
## Residuals 4.2642 17
```

Ordered Logistic Regression

```
# Create ordered variable write3 as
# a factor with levels 1, 2, and 3
some_ed_data$write3 <-
  cut(some_ed_data$write, c(0, 48, 57, 70),
    right = TRUE,
    labels = c(1,2,3))

table(some_ed_data$write3)

# fit ordered logit model & store results 'some_write_data'
some_write_data <-
  polr(write3 ~
    female + read + socst, data = some_ed_data,
    Hess=TRUE)

summary(some_write_data)
```

```
## 
##   1   2   3
## 61 61 78

## Call:
## polr(formula = write3 ~ female + read + socst, data = some_ed_d...
##           Hess = TRUE)
## 
## Coefficients:
##              Value Std. Error t value
## female  1.28543  0.32445  3.962
## read    0.11772  0.02136  5.512
## socst   0.08019  0.01944  4.124
## 
## Intercepts:
##              Value Std. Error t value
## 1|2  9.7037  1.1968    8.1080
## 2|3 11.8001  1.3041    9.0486
```

Principal Components Analysis (PCA)

```
princomp(formula = ~read + write + math + science + socst, data = some_ed_data)

## Call:
## princomp(formula = ~read + write + math + science + socst, data = some_ed_data)
##
## Standard deviations:
##      Comp.1     Comp.2     Comp.3     Comp.4     Comp.5
## 18.252929  7.677044  6.213371  5.774331  5.429881
##
## 5 variables and 200 observations.
```

Repeated Measures Logistic Regression

```
glmer(highpulse ~ diet + (1 | id), data = some_exercise_data, family = binomial)

## Generalized linear mixed model fit by maximum likelihood (Laplace Approximation) [
## glmerMod]
## Family: binomial ( logit )
## Formula: highpulse ~ diet + (1 | id)
## Data: some_exercise_data
##      AIC    BIC  logLik deviance df.resid
## 105.4679 112.9674 -49.7340   99.4679      87
## Random effects:
## Groups Name        Std.Dev.
## id     (Intercept) 1.821
## Number of obs: 90, groups: id, 30
## Fixed Effects:
## (Intercept)      diet
##             -3.148     1.145
```

Simple Linear Regression

```
lm(some_ed_data$write ~ some_ed_data$read)

##
## Call:
## lm(formula = some_ed_data$write ~ some_ed_data$read)
##
## Coefficients:
## (Intercept)  some_ed_data$read
##           23.9594          0.5517
```

Simple Logistic Regression

```
glm(some_ed_data$female ~ some_ed_data$read, family = binomial)

## 
## Call: glm(formula = some_ed_data$female ~ some_ed_data$read, family = binomial)
## 
## Coefficients:
## (Intercept) some_ed_data$read
##          0.72609        -0.01044
## 
## Degrees of Freedom: 199 Total (i.e. Null);  198 Residual
## Null Deviance:      275.6
## Residual Deviance: 275.1    AIC: 279.1
```

Two Independent Samples *t*-test

```
t.test(some_ed_data$read ~ some_ed_data$female)

##
##      Welch Two Sample t-test
##
## data: some_ed_data$read by some_ed_data$female
## t = 0.74506, df = 188.46, p-value = 0.4572
## alternative hypothesis: true difference in means between group 0 and group 1 is not equal to 0
## 95 percent confidence interval:
## -1.796263 3.976725
## sample estimates:
## mean in group 0 mean in group 1
##           52.82418          51.73394
```

Wilcoxon-Mann-Whitney Test

```
wilcox.test(some_ed_data$read ~ some_ed_data$female)

## Wilcoxon rank sum test with continuity correction
## data: some_ed_data$read by some_ed_data$female
## W = 5300, p-value = 0.4029
## alternative hypothesis: true location shift is not equal to 0
```

Wilcoxon Signed Rank Sum Test

```
wilcox.test(some_ed_data$write, some_ed_data$read, paired = TRUE)

## Wilcoxon signed rank test with continuity correction
## data: some_ed_data$write and some_ed_data$read
## V = 9261, p-value = 0.3666
## alternative hypothesis: true location shift is not equal to 0
```

Other Approaches

There are so many other approaches that are for specific cases or use statistical approaches, but aren't themselves statistics. With that said, the approaches given in this overview cover the gambit of what you will likely need

Some Extra Things

Reporting

After running a statistical test successfully, it can be difficult to know how to report the results. The `report` package automatically produces reports of models and dataframes according to best practices guidelines. [Click here](#) for more information.

Visualizations

Interested in making incredible visuals? Check out #tidytuesday on Twitter. You do not need an account for access.

Something useless

If you are a fan of the show Rick & Morty, consider downloading the most pointless package `mortyr` to do pointless statistics on pointless data. More about the package [here](#).

Source

A majority of the information included in this survey of approached was scraped from the web using `R` via the UCLA Institute for Digital Research & Education site using the `xml2` package. They also fully support SAS, SPSS (for some reason), Stata, and Mplus.

Thats it!

If you have any questions, please reach out



This work is licensed under a
Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License