# Statistics

Introduction to R for Public Health Researchers

# Statistics

Now we are going to cover how to perform a variety of basic statistical tests in R.

· Correlation

· T-tests/Rank-sum tests

· Linear Regression

· Logistic Regression

· Proportion tests

· Chi-squared

· Fisher's Exact Test

Note: We will be glossing over the statistical theory and "formulas" for these tests. There are plenty of resources online for learning more about these tests, as well as dedicated Biostatistics series at the School of Public Health

# Correlation

`cor()` performs correlation in R

```
cor(x, y = NULL, use = "everything",
    method = c("pearson", "kendall", "spearman"))
```

Like other functions, if there are NAs, you get NA as the result. But if you specify use only the complete observations, then it will give you correlation on the non-missing data.

```
library(readr)
circ = read_csv("http://johnmuschelli.com/intro_to_r/data/Charm_City_Circulato
cor(circ$orangeAverage, circ$purpleAverage, use="complete.obs")
```

```
[1] 0.9195356
```

# Correlation with `corrr`

The `corrr` package allows you to do correlations easily:

```
library(corrr); library(dplyr)
circ %>%
  select(orangeAverage, purpleAverage) %>%
  correlate()
```

```
Correlation method: 'pearson'
Missing treated using: 'pairwise.complete.obs'

# A tibble: 2 x 3
  rowname        orangeAverage purpleAverage
  <chr>                  <dbl>         <dbl>
1 orangeAverage             NA         0.920
2 purpleAverage          0.920            NA
```
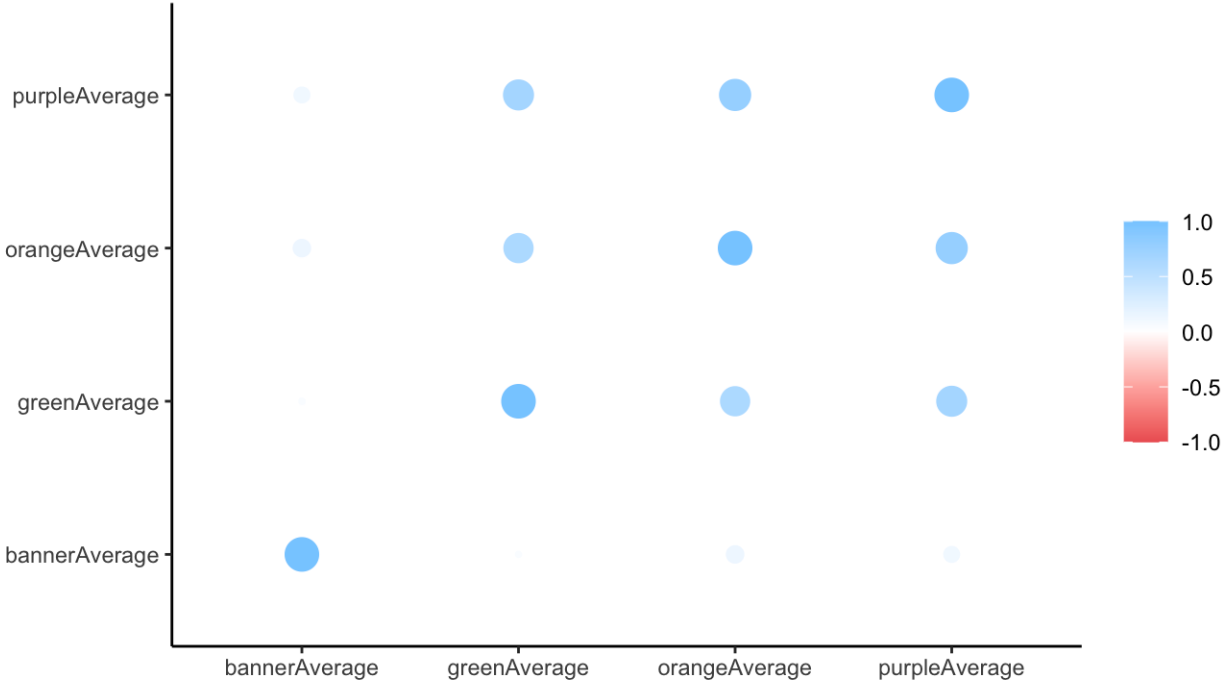
# Correlation

`correlate` is usually better with more than 2 columns:

```
avgs = circ %>% select(ends_with("Average"))
cobj = avgs %>% correlate(use = "complete.obs", diagonal = 1)
cobj %>% fashion(decimals = 3)
```

```
        rowname orangeAverage purpleAverage greenAverage bannerAverage
1 orangeAverage         1.000          .908         .840          .545
2 purpleAverage          .908         1.000         .867          .521
3  greenAverage          .840          .867        1.000          .453
4 bannerAverage          .545          .521         .453         1.000
```

# Correlation

```
cobj %>% rplot()
```

# Correlation

You can also use `cor.test()` to test for whether correlation is significant (ie non-zero). Note that linear regression may be better, especially if you want to regress out other confounders.

```
ct = cor.test(circ$orangeAverage, circ$purpleAverage,
              use = "complete.obs")
ct
```

```
        Pearson's product-moment correlation

data:  circ$orangeAverage and circ$purpleAverage
t = 73.656, df = 991, p-value < 2.2e-16
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 0.9093438 0.9286245
sample estimates:
      cor
0.9195356
```

# Correlation

For many of these testing result objects, you can extract specific slots/results as numbers, as the `ct` object is just a list.

```
# str(ct)
names(ct)
```

```
[1] "statistic"   "parameter"   "p.value"      "estimate"     "null.value"
[6] "alternative" "method"      "data.name"    "conf.int"
```

```
ct$statistic
```

```
       t
73.65553
```

```
ct$p.value
```

```
[1] 0
```

# Broom package

The `broom` package has a `tidy` function that puts most objects into `data.frames` so that they are easily manipulated:

```
library(broom)
tidy_ct = tidy(ct)
tidy_ct
```
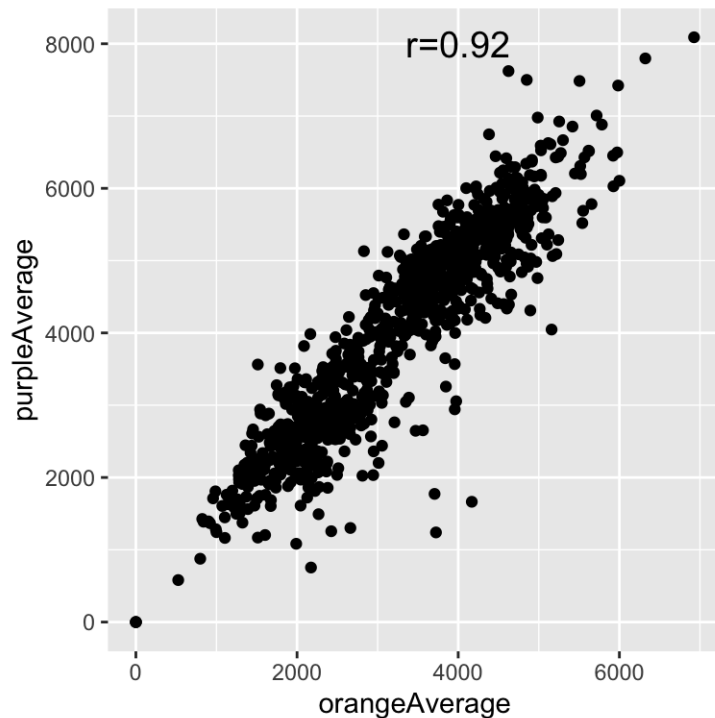
```
# A tibble: 1 x 8
  estimate statistic p.value parameter conf.low conf.high method      alternati
     <dbl>     <dbl>   <dbl>     <int>    <dbl>     <dbl> <chr>       <chr>
1    0.920      73.7       0       991    0.909     0.929 Pearson's... two.sided
```

# Correlation

Note that you can add the correlation to a plot, via the `annotate`

```
library(ggplot2)
txt = paste0("r=", signif(ct$estimate,3))
q = qplot(data = circ, x = orangeAverage, y = purpleAverage)
q + annotate("text", x = 4000, y = 8000, label = txt, size = 5)
```

# T-tests

The T-test is performed using the `t.test()` function, which essentially tests for the difference in means of a variable between two groups.

In this syntax, x and y are the column of data for each group.

```
tt = t.test(circ$orangeAverage, circ$purpleAverage)
tt
```

```
	Welch Two Sample t-test

data:  circ$orangeAverage and circ$purpleAverage
t = -17.076, df = 1984, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -1096.7602  -870.7867
sample estimates:
mean of x mean of y
 3033.161  4016.935
```

# T-tests

Using `t.test` treats the data as independent. Realistically, this data should be treated as a paired t-test. The `paired = TRUE` argument to do a paired test

```
t.test(circ$orangeAverage, circ$purpleAverage, paired = TRUE)
```

```
    Paired t-test

data:  circ$orangeAverage and circ$purpleAverage
t = -42.075, df = 992, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -799.783 -728.505
sample estimates:
mean of the differences
              -764.144
```

# T-tests

`t.test` saves a lot of information: the difference in means `estimate`, confidence interval for the difference `conf.int`, the p-value `p.value`, etc.

```
names(tt)
```

```
[1] "statistic"   "parameter"   "p.value"     "conf.int"   "estimate"
[6] "null.value"  "stderr"      "alternative" "method"     "data.name"
```

# T-tests

```
tidy(tt)
```

```
# A tibble: 1 x 10
  estimate estimate1 estimate2 statistic  p.value parameter conf.low conf.high
     <dbl>     <dbl>     <dbl>     <dbl>    <dbl>     <dbl>    <dbl>     <dbl>
1    -984.     3033.     4017.     -17.1 4.20e-61     1984.   -1097.     -871.
# … with 2 more variables: method <chr>, alternative <chr>
```

# T-tests

You can also use the 'formula' notation. In this syntax, it is $y \sim x$, where $x$ is a factor with 2 levels or a binary variable and $y$ is a vector of the same length.

```
library(tidyr)
long = circ %>%
  select(date, orangeAverage, purpleAverage) %>%
  gather(key = line, value = avg, -date)
tt = t.test(avg ~ line, data = long)
tidy(tt)
```

```
# A tibble: 1 x 10
  estimate estimate1 estimate2 statistic  p.value parameter conf.low conf.high
     <dbl>     <dbl>     <dbl>     <dbl>    <dbl>     <dbl>    <dbl>     <dbl>
1    -984.     3033.     4017.     -17.1 4.20e-61     1984.   -1097.     -871.
# … with 2 more variables: method <chr>, alternative <chr>
```

# Wilcoxon Rank-Sum Tests

Nonparametric analog to t-test (testing medians):

```
tidy(wilcox.test(avg ~ line, data = long))
```

```
# A tibble: 1 x 4
  statistic  p.value method                                  alternati
      <dbl>    <dbl> <chr>                                   <chr>
1   336714. 4.56e-58 Wilcoxon rank sum test with continuity correct… two.sided
```

# Lab Part 1

[Website](Website)

# Analysis of Variance

The `aov` function exists for ANOVA, but we'd recommend `lm` (Linear models):

```
long3 = circ %>%
  select(date, orangeAverage, purpleAverage, bannerAverage) %>%
  gather(key = line, value = avg, -date)
anova_res = aov(avg ~ line, data = long3)
anova_res
```

```
Call:
   aov(formula = avg ~ line, data = long3)

Terms:
                        line  Residuals
Sum of Squares   2214242251 3724379811
Deg. of Freedom           2       2396

Residual standard error: 1246.762
Estimated effects may be unbalanced
1039 observations deleted due to missingness
```

# Linear Regression

Now we will briefly cover linear regression. I will use a little notation here so some of the commands are easier to put in the proper context.

$$y_i = \alpha + \beta x_i + \varepsilon_i$$

where:

- $y_i$ is the outcome for person i
- $\alpha$ is the intercept
- $\beta$ is the slope
- $x_i$ is the predictor for person i
- $\varepsilon_i$ is the residual variation for person i

# Linear Regression

The `R` version of the regression model is:

```
y ~ x
```

where:

- y is your outcome
- x is/are your predictor(s)

# Linear Regression

For a linear regression, when the predictor is binary this is the same as a t-test:

```
fit = lm(avg ~ line, data = long)
fit
```

```
Call:
lm(formula = avg ~ line, data = long)

Coefficients:
      (Intercept)   linepurpleAverage
           3033.2               983.8
```

'(Intercept)' is $\alpha$

'linepurpleAverage' is $\beta$

# Linear Regression

```
fit = lm(avg ~ line, data = long3)
fit
```

```
Call:
lm(formula = avg ~ line, data = long3)

Coefficients:
     (Intercept)   lineorangeAverage   linepurpleAverage
           827.3              2205.9              3189.7
```

# Linear Regression

The `summary` command gets all the additional information (p-values, t-statistics, r-square) that you usually want from a regression.

```
sfit = summary(fit)
print(sfit)
```

```
Call:
lm(formula = avg ~ line, data = long3)

Residuals:
    Min       1Q   Median       3Q      Max
-4016.9  -1008.4      5.6    973.1   4072.6

Coefficients:
                   Estimate Std. Error t value Pr(>|t|)
(Intercept)          827.27      75.88   10.90   <2e-16 ***
lineorangeAverage   2205.89      84.41   26.13   <2e-16 ***
linepurpleAverage   3189.67      85.57   37.27   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1247 on 2396 degrees of freedom
  (1039 observations deleted due to missingness)
Multiple R-squared:  0.3729,    Adjusted R-squared:  0.3723
F-statistic: 712.2 on 2 and 2396 DF,  p-value: < 2.2e-16
```

# Linear Regression

We can `tidy` linear models as well and it gives us all of this in a tibble:

```
tidy(fit)
```

```
# A tibble: 3 x 5
  term                estimate std.error statistic    p.value
  <chr>                  <dbl>     <dbl>     <dbl>      <dbl>
1 (Intercept)             827.      75.9      10.9 4.76e- 27
2 lineorangeAverage      2206.      84.4      26.1 1.16e-132
3 linepurpleAverage      3190.      85.6      37.3 2.98e-240
```

# Linear Regression

The `confint` argument allows for confidence intervals

```
tidy(fit, conf.int = TRUE)
```

```
# A tibble: 3 x 7
  term                estimate std.error statistic   p.value conf.low conf.high
  <chr>                  <dbl>     <dbl>     <dbl>     <dbl>    <dbl>     <dbl>
1 (Intercept)             827.      75.9      10.9 4.76e- 27     678.      976.
2 lineorangeAverage      2206.      84.4      26.1 1.16e-132    2040.     2371.
3 linepurpleAverage      3190.      85.6      37.3 2.98e-240    3022.     3357.
```

# Using Cars Data

```
http_data_dir = "http://johnmuschelli.com/intro_to_r/data/"
cars = read_csv(
  paste0(http_data_dir, "kaggleCarAuction.csv"),
  col_types = cols(VehBCost = col_double()))
head(cars)
```

```
# A tibble: 6 x 34
  RefId IsBadBuy PurchDate Auction VehYear VehicleAge Make  Model Trim  SubMod
  <dbl>    <dbl> <chr>     <chr>     <dbl>      <dbl> <chr> <chr> <chr> <chr>
1     1        0 12/7/2009 ADESA      2006          3 MAZDA MAZD… i     4D SED
2     2        0 12/7/2009 ADESA      2004          5 DODGE 1500… ST    QUAD (
3     3        0 12/7/2009 ADESA      2005          4 DODGE STRA… SXT   4D SED
4     4        0 12/7/2009 ADESA      2004          5 DODGE NEON  SXT   4D SED
5     5        0 12/7/2009 ADESA      2005          4 FORD  FOCUS ZX3   2D COU
6     6        0 12/7/2009 ADESA      2004          5 MITS… GALA… ES    4D SED
# … with 24 more variables: Color <chr>, Transmission <chr>, WheelTypeID <chr>
#   WheelType <chr>, VehOdo <dbl>, Nationality <chr>, Size <chr>,
#   TopThreeAmericanName <chr>, MMRAcquisitionAuctionAveragePrice <chr>,
#   MMRAcquisitionAuctionCleanPrice <chr>,
#   MMRAcquisitionRetailAveragePrice <chr>,
#   MMRAcquisitonRetailCleanPrice <chr>, MMRCurrentAuctionAveragePrice <chr>,
#   MMRCurrentAuctionCleanPrice <chr>, MMRCurrentRetailAveragePrice <chr>,
#   MMRCurrentRetailCleanPrice <chr>, PRIMEUNIT <chr>, AUCGUART <chr>,
#   BYRNO <dbl>, VNZIP1 <dbl>, VNST <chr>, VehBCost <dbl>, IsOnlineSale <dbl>,
#   WarrantyCost <dbl>
```

# Linear Regression

We'll look at vehicle odometer value by vehicle age:

```
fit = lm(VehOdo~VehicleAge, data = cars)
print(fit)

Call:
lm(formula = VehOdo ~ VehicleAge, data = cars)

Coefficients:
(Intercept)    VehicleAge
      60127          2723
```

# Linear Regression

Note that you can have more than 1 predictor in regression models.The interpretation for each slope is change in the predictor corresponding to a one-unit change in the outcome, holding all other predictors constant.

```
fit2 = lm(VehOdo ~ IsBadBuy + VehicleAge, data = cars)
tidy(fit2)
```

```
# A tibble: 3 x 5
  term          estimate std.error statistic  p.value
  <chr>            <dbl>     <dbl>     <dbl>    <dbl>
1 (Intercept)    60142.      135.      446.   0.
2 IsBadBuy        1329.      158.        8.42 3.83e-17
3 VehicleAge      2680.       30.3      88.5  0.
```

# Linear Regression: Interactions

The * does interactions:

```
fit3 = lm(VehOdo ~ IsBadBuy * VehicleAge, data = cars)
tidy(fit3)

# A tibble: 4 x 5
  term                   estimate std.error statistic p.value
  <chr>                     <dbl>     <dbl>     <dbl>   <dbl>
1 (Intercept)              60140.      143.    420.        0
2 IsBadBuy                  1347.      456.      2.95  0.00315
3 VehicleAge                2681.       32.5    82.4       0
4 IsBadBuy:VehicleAge       -3.79       88.7    -0.0427 0.966
```

# Linear Regression: Interactions

You can take out main effects with minus

```
fit4 = lm(VehOdo ~ IsBadBuy * VehicleAge -IsBadBuy , data = cars)
tidy(fit4)

# A tibble: 3 x 5
  term                  estimate std.error statistic  p.value
  <chr>                    <dbl>     <dbl>     <dbl>    <dbl>
1 (Intercept)             60272.      136.      443.  0.
2 VehicleAge               2653.      31.1      85.2  0.
3 IsBadBuy:VehicleAge       242.      30.7       7.89 3.19e-15
```

# Linear Regression

Factors get special treatment in regression models - lowest level of the factor is the comparison group, and all other factors are relative to its values.

```
fit3 = lm(VehOdo ~ factor(TopThreeAmericanName), data = cars)
tidy(fit3)
```

```
# A tibble: 5 x 5
  term                                estimate std.error statistic    p.value
  <chr>                                  <dbl>     <dbl>      <dbl>      <dbl>
1 (Intercept)                           68248.      93.0    734.     0.
2 factor(TopThreeAmericanName)FORD       8523.     158.      53.8    0.
3 factor(TopThreeAmericanName)GM         4952.     129.      38.4    2.74e-319
4 factor(TopThreeAmericanName)NULL      -2005.    6362.      -0.315  7.53e-  1
5 factor(TopThreeAmericanName)OTHER       585.     160.       3.66   2.55e-  4
```

# Logistic Regression and GLMs

Generalized Linear Models (GLMs) allow for fitting regressions for non-continuous/normal outcomes. The `glm` has similar syntax to the `lm` command. Logistic regression is one example. See `?family` for

```
glmfit = glm(IsBadBuy ~ VehOdo + VehicleAge, data=cars, family = binomial())
tidy(glmfit)
```

```
# A tibble: 3 x 5
  term            estimate    std.error statistic  p.value
  <chr>               <dbl>        <dbl>     <dbl>    <dbl>
1 (Intercept) -3.78          0.0638          -59.2 0.
2 VehOdo       0.00000834 0.000000853          9.78 1.33e-22
3 VehicleAge   0.268         0.00677           39.6 0.
```

# Tidying GLMs

```
tidy(glmfit, conf.int = TRUE)
```

```
# A tibble: 3 x 7
  term           estimate    std.error statistic  p.value     conf.low  conf.hig
  <chr>             <dbl>        <dbl>     <dbl>    <dbl>         <dbl>      <dbl
1 (Intercept) -3.78        0.0638        -59.2  0.          -3.90       -3.65
2 VehOdo       0.00000834  0.000000853     9.78 1.33e-22    0.00000667  0.000010
3 VehicleAge   0.268       0.00677        39.6  0.           0.255      0.281
```

# Tidying GLMs

```
tidy(glmfit, conf.int = TRUE, exponentiate = TRUE)
```

```
# A tibble: 3 x 7
  term         estimate   std.error statistic  p.value conf.low conf.high
  <chr>           <dbl>       <dbl>     <dbl>    <dbl>    <dbl>     <dbl>
1 (Intercept)    0.0229 0.0638          -59.2  0.         0.0202    0.0259
2 VehOdo         1.00   0.000000853       9.78 1.33e-22   1.00      1.00
3 VehicleAge     1.31   0.00677          39.6  0.         1.29      1.32
```

# Logistic Regression

Note the coefficients are on the original scale, we must exponentiate them for odds ratios:

```
exp(coef(glmfit))
```

```
(Intercept)         VehOdo   VehicleAge
 0.02286316   1.00000834   1.30748911
```

# Chi-squared tests

`chisq.test()` performs chi-squared contingency table tests and goodness-of-fit tests.

```
chisq.test(x, y = NULL, correct = TRUE,
           p = rep(1/length(x), length(x)), rescale.p = FALSE,
           simulate.p.value = FALSE, B = 2000)
```

```
tab = table(cars$IsBadBuy, cars$IsOnlineSale)
tab
```

```
       0     1
0  62375  1632
1   8763   213
```

# Chi-squared tests

You can also pass in a table object (such as `tab` here)

```
cq = chisq.test(tab)
cq
```

```
	Pearson's Chi-squared test with Yates' continuity correction

data:  tab
X-squared = 0.92735, df = 1, p-value = 0.3356
```

```
names(cq)
```

```
[1] "statistic" "parameter" "p.value"    "method"     "data.name" "observed"
[7] "expected"  "residuals" "stdres"
```

```
cq$p.value
```

```
[1] 0.3355516
```

# Chi-squared tests

Note that does the same test as `prop.test`, for a 2x2 table (`prop.test` not relevant for greater than 2x2).

```
chisq.test(tab)
```

```
    Pearson's Chi-squared test with Yates' continuity correction

data:  tab
X-squared = 0.92735, df = 1, p-value = 0.3356
```

```
prop.test(tab)
```

```
    2-sample test for equality of proportions with continuity correction

data:  tab
X-squared = 0.92735, df = 1, p-value = 0.3356
alternative hypothesis: two.sided
95 percent confidence interval:
 -0.005208049  0.001673519
sample estimates:
   prop 1    prop 2
0.9745028 0.9762701
```

# Fisher's Exact test

`fisher.test()` performs contingency table test using the hypogeometric distribution (used for small sample sizes).

```
fisher.test(x, y = NULL, workspace = 200000, hybrid = FALSE,
           control = list(), or = 1, alternative = "two.sided",
           conf.int = TRUE, conf.level = 0.95,
           simulate.p.value = FALSE, B = 2000)
```
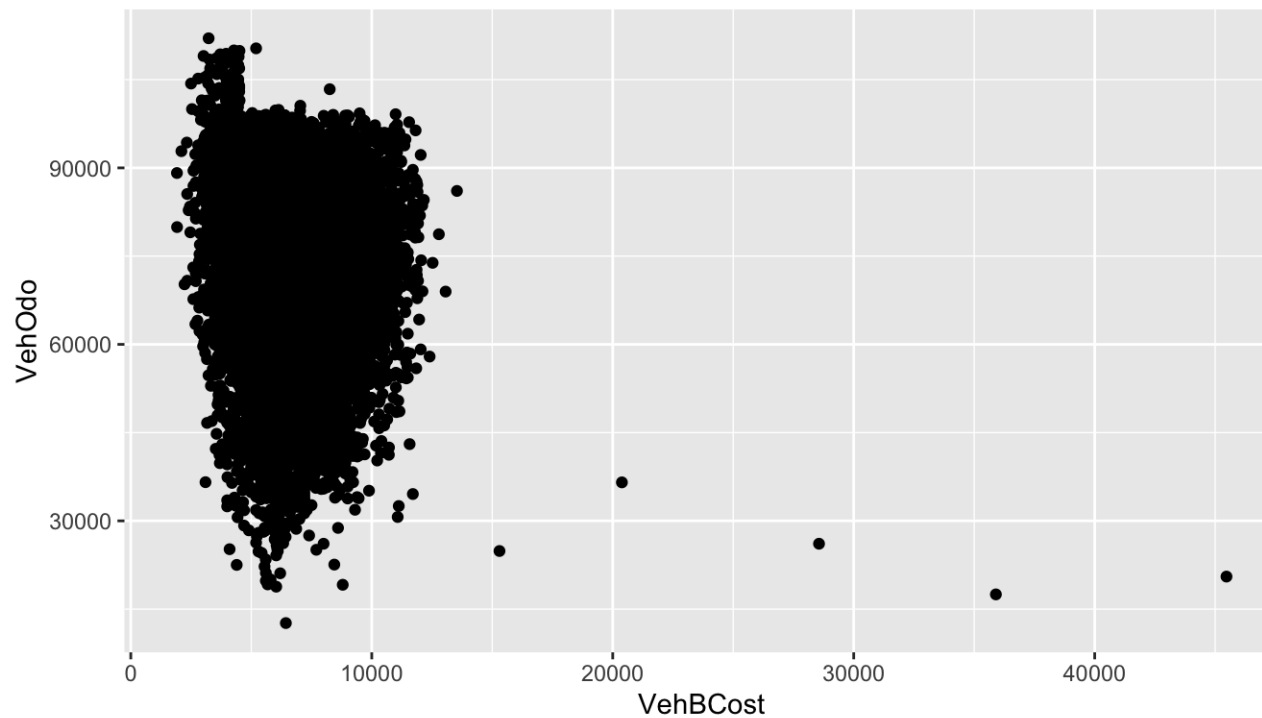
```
fisher.test(tab)
```

```
    Fisher's Exact Test for Count Data

data:  tab
p-value = 0.3324
alternative hypothesis: true odds ratio is not equal to 1
95 percent confidence interval:
 0.8001727 1.0742114
sample estimates:
odds ratio
 0.9289923
```

# Sampling

Also, if you want to only plot a subset of the data (for speed/time or overplotting)

```
samp.cars = sample_n(cars, size = 10000)
samp.cars = sample_frac(cars, size = 0.2)
ggplot(aes(x = VehBCost, y = VehOdo),
       data = samp.cars) + geom_point()
```

# Lab Part 2

[Website](Website)