

Love at First Stat — Rediscover the Fundamentals

Love Data Week 2026

Statistical Consulting and Survey Research Unit
University of Waterloo

2026-2-12

Agenda

In this workshop, we will cover:

1. Types of variables
 - ▶ continuous, categorical, ...
2. Outliers and extremem values
3. Hypothesis testing: the standard procedure
4. T-test
5. ANOVA tests
6. Correlation test
7. Linear regression
8. Model diagnostics

1. Types of variables

It is common to see different variables in a data set. There are many types of variables, but we can generally classify the variables as:

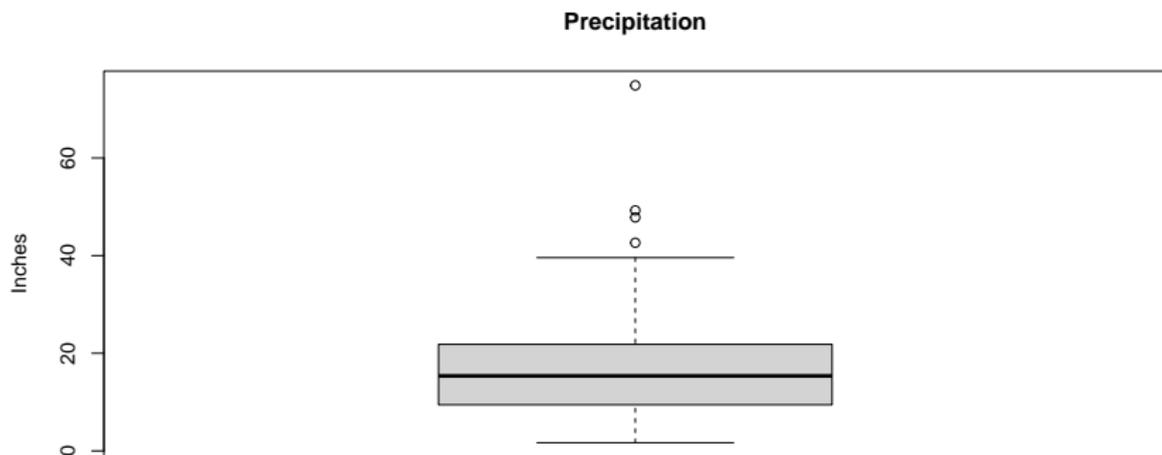
- ▶ **Discrete/categorical:** only take values over a finite set of values (or levels), e.g.,
 - ▶ A university student's major.
 - ▶ A person's blood type.
 - ▶ The entrees on a menu.
 - ▶ A person's eye colour.
- ▶ **Binary variable:** categorical variables with only 2 levels.
- ▶ **Continuous:** take any value over a continuous range, e.g.
 - ▶ age in years,
 - ▶ number of work hours,
 - ▶ midterm scores, etc.

2. Outliers and extreme values

- ▶ The boxplot is often used to visualize the distribution of a numeric variable, and potential outliers.
- ▶ The outliers are presented as dots or points beyond the box and its whiskers.
- ▶ The rule used to identify the outliers is called the $1.5 \times IQR$ rule, where $IQR = Q_3 - Q_1$.

2.1 Use `boxplot()` to identify outliers

```
boxplot(rain_df$PRECIP, main="Precipitation",  
        ylab= "Inches")
```



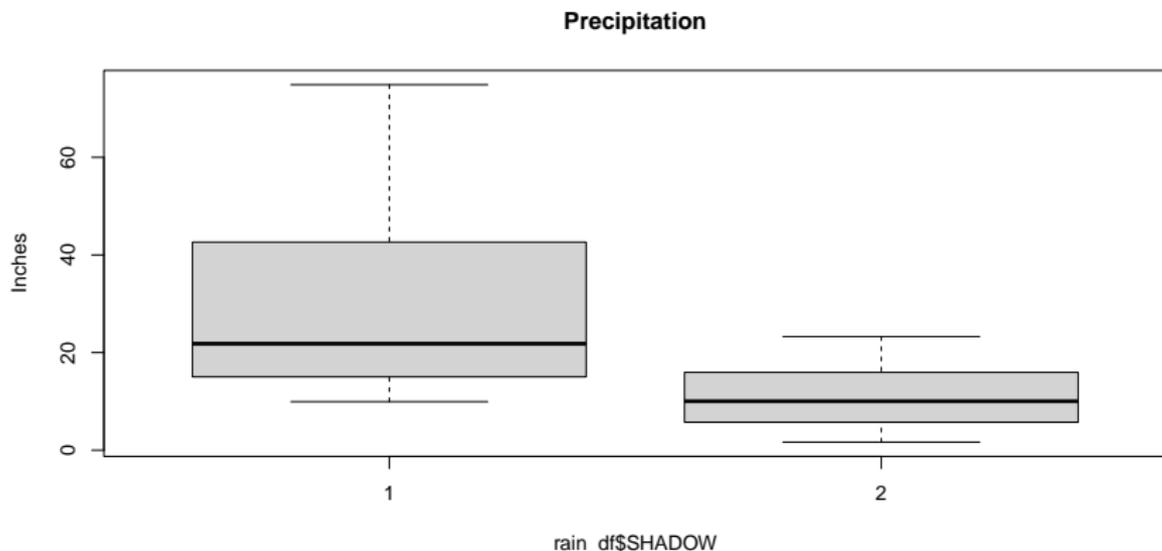
Why 1.5?

- ▶ Tukey who introduced the $1.5 \times IQR$ rule claimed that, “2 was too big, and 1 was too small.”
- ▶ This rule suggests whether a value is potentially an outlier.
- ▶ However, experts opinion is crucial when deciding whether an observation is an outlier, or an extreme value.
- ▶ An extreme value may contain interesting information and must not be dismissed without careful thoughts.

2.2 Side-by-side boxplots

In the side-by-side boxplots, notice that there is no potential outlier. What happened to the circles in the first boxplot?

```
boxplot(rain_df$PRECIP~rain_df$SHADOW, main="Precipitation",  
        ylab= "Inches")
```



3. Hypothesis Test

Oftentimes, we are interested to investigate the relationships between multiple variables. For example,

- ▶ Does the distance from Pacific Ocean affect precipitation?
- ▶ Is the midterm average this term higher than the average last term?
- ▶ Is my average grocery purchase every week around \$50?

The questions lead us to *hypothesis testing*.

3.1 Steps for hypothesis testing

1. Formulate the null and alternate hypothesis.
2. Choose and evaluate the appropriate test statistic (with R).
3. Assess the strength of the evidence against the null hypothesis.
4. Interpret the results.

Note that Step (2) is a tedious (and sometimes iterative) process. It cannot be checked off using “a few clicks”.

Step 1: Null and alternate hypotheses

In hypotheses testing, we begin by translating a question of interest into the appropriate null and alternate hypotheses:

- ▶ Null hypothesis: Status quo statement that is commonly denoted as H_0 . (An assertion that you want to prove wrong.)
- ▶ Alternate hypothesis: The answer we are looking for, commonly denoted as H_1 , H_A or H_a

Example 1: Is my average grocery purchase every week around \$50?

- ▶ H_0 : My average grocery purchase is \$50.
- ▶ H_A : My average grocery purchase is not \$50.

Step 2: Statistical tests

Step 2 involves choosing the appropriate test statistics. In this workshop, we will briefly discuss several common statistical tests:

- ▶ t-Test
- ▶ Analysis of Variance (ANOVA)
- ▶ Pearson Correlation Test
- ▶ Linear Regression

Step 3: Strength of evidence

- ▶ In most statistical tests, a p-value will be produced.
- ▶ The p-value is the probability of finding results equal or more extreme than the observed results (data), given that the null hypothesis (H_0) is true.
- ▶ The smaller the p-value, the more evidence we have against the null hypothesis.
- ▶ The default significance levels are 0.01, 0.05 and 0.10.
 - ▶ When the p-value is less than the significance level (of your choice), we say that we have evidence against the null hypothesis in favor of the alternate hypothesis.
 - ▶ When the p-value is greater than the default value, we say that we do not have sufficient evidence against the null hypothesis. Sometimes, we say “we do not reject the null hypothesis”.
 - ▶ However, we almost always avoid saying “we accept the null hypothesis”.

Step 4: Drawing conclusion

The final step in hypothesis testing is to draw conclusion in the words of the problem.

Example: Is my average grocery purchase every week around \$50?

- ▶ H_0 : My average grocery purchase is \$50.
- ▶ H_A : My average grocery purchase is not \$50.

If the p-value is less than 0.05, we say that there is evidence against H_0 , in favour of H_A , i.e. the data suggests that my true average grocery purchase is not \$50.

If the p-value is greater than 0.05, we say that there is not enough evidence against H_0 , i.e. the data suggests that my true average grocery purchase is \$50.

3.2 Statistical significance vs practical significance

- ▶ Statistical inference techniques test for statistical significance.
- ▶ Statistical significance means that the effect observed in a sample is very unlikely to occur if the null hypothesis is true.
- ▶ Whether this observed effect has practical importance is an entirely different question. The experts in the field of interest determine whether these results have any practical importance.

3.3 Danger of over reliance on p-values

The ASA's Statement on p-values:

- ▶ P-values can indicate how incompatible the data are with a specified statistical model.
- ▶ P-values do not measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone.
- ▶ Scientific conclusions and business or policy decisions should not be based only on whether a p-value passes a specific threshold.
- ▶ Proper inference requires full reporting and transparency
- ▶ A p-value, or statistical significance, does not measure the size of an effect or the importance of a result.
- ▶ By itself, a p-value does not provide a good measure of evidence regarding a model or hypothesis.

4. T-Tests

- ▶ One of the most commonly used statistical test is the t-test.
- ▶ The **one-sample t-test** is used to compare the mean of a variable to a hypothetical value. In most cases, the hypothetical value comes from theory.
- ▶ The **two-samples t-test** is used to compare the means of two variables. It is often used to determine whether a treatment has an effect on the population of interest, or whether two groups are different from one another.
- ▶ The **paired t-test** is commonly used to investigate the difference of a variable pre- and post-treatment. Oftentimes, every subject of the study produces a pair of observations, or two similar subjects will be paired up.

4.1 One-sample t-test

- ▶ Compare the sample mean with a **hypothetical value**, often obtained from literature.
- ▶ Example: Is the average annual precipitation in California around 25 inches?
- ▶ Assumptions
 - ▶ Each sample was drawn from an **identically independently distributed** (IID) Normal distribution.
 - ▶ If the sample size is large enough, normality is not necessary because Central Limit Theorem applies.

4.1 Example

H_0 : The average annual precipitation in California is 25 inches.

H_A : The average annual precipitation in California is not 25 inches.

- ▶ We can perform this two-sided one-sample t-test in R:

```
t.test(rain_df$PRECIP, mu=25)
```

One Sample t-test

```
data: rain_df$PRECIP
```

```
t = -1.7112, df = 29, p-value = 0.09773
```

```
alternative hypothesis: true mean is not equal to 25
```

```
95 percent confidence interval:
```

```
13.60088 26.01379
```

```
sample estimates:
```

```
mean of x
```

```
19.80733
```

4.1 Example: interpretation

- ▶ The p-value is greater than 0.05 and hence, we claim no evidence against H_0 , i.e. the average annual precipitation in California is 25 inches.
- ▶ The average annual precipitation is 19.81 inches. How is this 25 inches?

4.2 Two-samples t-test

- ▶ Compare the means of 2 samples or groups.
- ▶ Example: Do precipitation differs due to shadow?
- ▶ Another way to think about this is that we are interested to investigate the relationship between a continuous and a categorical variable.
- ▶ Assumptions
 - ▶ Each sample was drawn from an **identically independently distributed** (IID) Normal distribution.
 - ▶ The population variances of the 2 samples are similar.
 - ▶ This is not as important when using R because the default in R is to assume unequal population variances. The results of tests for equal and unequal population variances will be the same if the population variances are the same.

4.2 Example: two-sided t-test

Is the average annual precipitation affected by SHADOW?

- ▶ H_0 : The precipitation is not affected by SHADOW.
- ▶ H_A : The precipitation is affected by SHADOW.

In another words,

- ▶ H_0 : The precipitation on the Leeward side is the same as the precipitation on the Westward side.
- ▶ H_A : The precipitation on the Leeward side is not the same as the precipitation on the Westward side.

4.2 Example: two-sided t-test

```
t.test(rain_df$PRECIP~rain_df$SHADOW)
```

Welch Two Sample t-test

data: rain_df\$PRECIP by rain_df\$SHADOW

t = 3.5309, df = 14.01, p-value = 0.003321

alternative hypothesis: true difference in means between group 1

95 percent confidence interval:

7.743558 31.702957

sample estimates:

mean in group 1 mean in group 2

30.98385

11.26059

4.2 Example: one-sided t-test

Is there higher precipitation on the Westward side?

- ▶ H_0 : The precipitation on the Leeward side is the same as the precipitation on the Westward side.
- ▶ H_A : The precipitation on the Leeward side is less than the precipitation on the Westward side.

or

- ▶ H_A : The precipitation on the Westward side is more than the precipitation on the Leeward side.

4.2 Example: one-sided t-test

- ▶ In R, the default test is the two-sided test.
- ▶ To perform a one-sided test, we need to make changes to the `alternative` in the code.
- ▶ If the alternate hypothesis is $\mu_{Group1} < \mu_{Group2}$, then we have `alternative = "less"`. Otherwise, set `alternative = "greater"`.
- ▶ By default, `Group1` is the category in which the category's name comes first in alphabetical order.

4.2 Example: one-sided t-test

```
t.test(rain_df$PRECIP~rain_df$SHADOW, alternative = "less")
```

Welch Two Sample t-test

data: rain_df\$PRECIP by rain_df\$SHADOW

t = 3.5309, df = 14.01, p-value = 0.9983

alternative hypothesis: true difference in means between group 1

95 percent confidence interval:

-Inf 29.56122

sample estimates:

mean in group 1 mean in group 2

30.98385

11.26059

5. Analysis of variance (ANOVA)

- ▶ The name analysis of variance can be misleading. It is actually a test on means across 2+ samples (groups).
- ▶ Consider the oats data set from the MASS library, with 4 variables: B (Blocks), V (Varieties), N (Nitrogen treatment) and Y (Yield of crop)

```
library(MASS)
str(oats)
```

```
'data.frame': 72 obs. of 4 variables:
 $ B: Factor w/ 6 levels "I","II","III",...: 1 1 1 1 1 1 1 1 1 1
 $ V: Factor w/ 3 levels "Golden.rain",...: 3 3 3 3 1 1 1 1 1 1
 $ N: Factor w/ 4 levels "0.0cwt","0.2cwt",...: 1 2 3 4 1 2
 $ Y: int 111 130 157 174 117 114 161 141 105 140 ...
```

5.1 Explanatory and response variables (recap)

The most common goal in research is to understand relationship between variables. These variables are typically categorized as:

- ▶ Response variable (or dependent variable): An outcome of the study or of interest.
- ▶ Explanatory variable (or independent variable): A measure in the study used to explain, predict or influence the response variable.

In this workshop, we will only consider response variables that are continuous.

5.2 One-way ANOVA

- ▶ Compare the response variable across samples grouped by one explanatory variable.
- ▶ Example: Suppose we are interested to find out whether the average yield (response variable) of each variety (within the independent variable, V) are the same.
- ▶ There are 3 varieties of oats:

```
levels(oats$V)
```

```
[1] "Golden.rain" "Marvellous"  "Victory"
```

5.2 One-way ANOVA

It is tempting to compare the average yield of the varieties in pairs using t-test:

- ▶ Golden.rain vs Marvellous
- ▶ Golden.rain vs Victory
- ▶ Marvellous vs Victory

However, such paired comparisons have limitations:

- ▶ the process can be tedious when there are many pairs
- ▶ the risk of a type I error increases when making multiple statistical tests. Type I error means rejecting the null hypothesis when it's actually true, i.e., false negative.

5.2 One-way ANOVA

In one-way ANOVA, we test the null hypothesis that k populations all have the same mean

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_k$$

against the alternative hypothesis that the population means are not all equal.

- ▶ Assumptions: same as those of the pooled-variance two-samples t-test:
 - ▶ The samples are independent simple random samples from the populations.
 - ▶ The populations are normally distributed.
 - ▶ The population variances are equal. ANOVA works poorly if the variances are extremely different.

5.2 One-way ANOVA

```
oats_aov <- aov(Y~V,data=oats)
summary(oats_aov)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
V	2	1786	893.2	1.228	0.299
Residuals	69	50200	727.5		

- ▶ Assuming the model fit is good, the corresponding p-value is 0.299, indicating no sufficient evidence that the three varieties of oats have different yields.

5.3 Two-way ANOVA

- ▶ Compare the response variable across samples grouped by two explanatory variables, and understand the combined impact.
- ▶ Suppose we are interested to find out whether milk content (`hasMilk`) and temperature (`temp`) can affect the amount of calories.
- ▶ Instead of performing two individual t-tests, we would perform a two-way ANOVA.
- ▶ In general, when the response variables of the one-way ANOVA are the same, we try to use one model.
- ▶ As mentioned earlier, performing multiple statistical tests is not only tedious, but also increases the risk of a Type I error.

5.3 Two-way ANOVA

```
drinks_aov <- aov(calories~factor(hasMilk) + temp,  
                 data=drinks_df)  
summary(drinks_aov)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
factor(hasMilk)	1	191753	191753	24.893	1.9e-05	***
temp	1	807	807	0.105	0.748	
Residuals	33	254197	7703			

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

- ▶ The variable `hasMilk` has corresponding p-value smaller than 0.05. This indicates strong evidence that the groups within that variable have different averages of calories.
- ▶ The variable `temp` has corresponding p-value larger than 0.05. This indicates that holding everything else constant, there is no evidence the Hot and Cold drinks have different calories on average.

5.3 Two-way ANOVA

- ▶ A common practice is to remove the insignificant variable from the model. However, we advise relying on significance entirely.
- ▶ Every variable contributes differently to a model.
- ▶ Some variable contribute by explaining the variability of the response variable. These variables will appear with small p-values.
- ▶ Some variable contribute by holding the structure of a model. They may not be significant, but they help the model meets its assumptions.
- ▶ It is important to check the model fit when variable(s) is/are added or removed from the model.

5.4 Post-hoc Tests

If there is strong evidence that not all the population means are equal for one variable, the next question is which categories are different. The ANOVA does not tell us which population means are different.

To explore where the difference lies, we perform the post-hoc tests. The post-hoc tests control for family-wise error rate. Here are a few common post-hoc tests:

- ▶ Fisher's Least Significant Difference (LSD), `LSD.test()`,
- ▶ Bonferroni correction, `pairwise.t.test(x, g, p.adjust.method="bonferroni")`,
- ▶ Tukey's Honest Significant Different, `TukeyHSD()`, and
- ▶ Scheffe's, `ScheffeTest()`.

6. Correlation

- ▶ Pearson correlation, r , is used to measure the linear relationship between two continuous variables.
- ▶ It is also unitless.
- ▶ Correlation is between -1 and 1 .
 - ▶ If $r \approx 0$, the linear relationship between two variables is weak.
 - ▶ If $r \approx 1$, there is a strong positive linear relationship between two variables.
 - ▶ If $r \approx -1$, there is a strong negative linear relationship between two variables.

6. Correlation

Suppose we are interested to evaluate the correlation between ALTITUDE and PRECIP in R:

```
cor(rain_df$PRECIP, rain_df$ALTITUDE)
```

```
[1] 0.3020067
```

- ▶ It does not matter which variable is first because the correlation between X and Y is the same as the correlation between Y and X.

6.1 Pearson's correlation test

The goal of this hypothesis test is to test the null hypothesis that the true correlation is equal to zero.

$$H_0 : r = 0$$

If $p\text{-value} < 0.05$, we say that there is evidence against the null hypothesis in favour of the alternate hypothesis. In another word, we have evidence that the true correlation cannot be zero and hence there exists linear relationship between the two variables.

6.1 Pearson's correlation test

```
cor.test(rain_df$PRECIP, rain_df$ALTITUDE)
```

Pearson's product-moment correlation

```
data: rain_df$PRECIP and rain_df$ALTITUDE
t = 1.6763, df = 28, p-value = 0.1048
alternative hypothesis: true correlation is not equal to 0
95 percent confidence interval:
 -0.0653756  0.5972887
sample estimates:
      cor
0.3020067
```

- ▶ The output showed that the p-value is around 0.10, which is greater than 0.05. This implies that we do not have sufficient evidence against the null hypothesis. In another word, precipitation and altitude are not correlated.

7. Linear Regression

The linear regression is also known as linear model. It is widely used in data analysis because:

- ▶ the model assumptions are often found satisfactory among many data sets; and
- ▶ the interpretation of each parameter in the model is easy and clear.

Assumptions:

- ▶ given the predictors, the expectation of the response is a linear function.
- ▶ the errors are normally distributed.
- ▶ the errors are independent of one another.
- ▶ the errors have mean zero and equal variance.

7.2 Simple Linear Model

- ▶ A simple linear model investigates possible linear relationship between two random variables.
- ▶ The response variable is a continuous variable.
- ▶ The explanatory variable can be of any type.

Example: Suppose we want to know whether the altitude of the station affect annual precipitation?

- ▶ The dependent variable here is annual precipitation, whereas the independent variable is the independent variable.
- ▶ To fit this linear regression model in R,

```
model <- lm(PRECIP~ALTITUDE, data=rain_df)
```

7.2 Simple Linear Model: R output

```
summary(model)
```

Call:

```
lm(formula = PRECIP ~ ALTITUDE, data = rain_df)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-20.620	-8.479	-2.729	4.555	58.271

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	16.514799	3.539141	4.666	6.9e-05 ***
ALTITUDE	0.002394	0.001428	1.676	0.105

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 16.13 on 28 degrees of freedom

Multiple R-squared: 0.09121, Adjusted R-squared: 0.05875

F-statistic: 2.81 on 1 and 28 DF. p-value: 0.1048

7.2 Simple Linear Model: R output

This result is the same as that of a correlation test:

```
cor.test(rain_df$PRECIP,rain_df$ALTITUDE)
```

Pearson's product-moment correlation

data: rain_df\$PRECIP and rain_df\$ALTITUDE

t = 1.6763, df = 28, p-value = 0.1048

alternative hypothesis: true correlation is not equal to 0

95 percent confidence interval:

-0.0653756 0.5972887

sample estimates:

cor

0.3020067

7.2 Simple Linear Model: Model diagnostics

When the assumptions of the linear regression model are satisfied, the model is powerful in terms of inference and interpretation. How do we know whether the assumptions are satisfied?

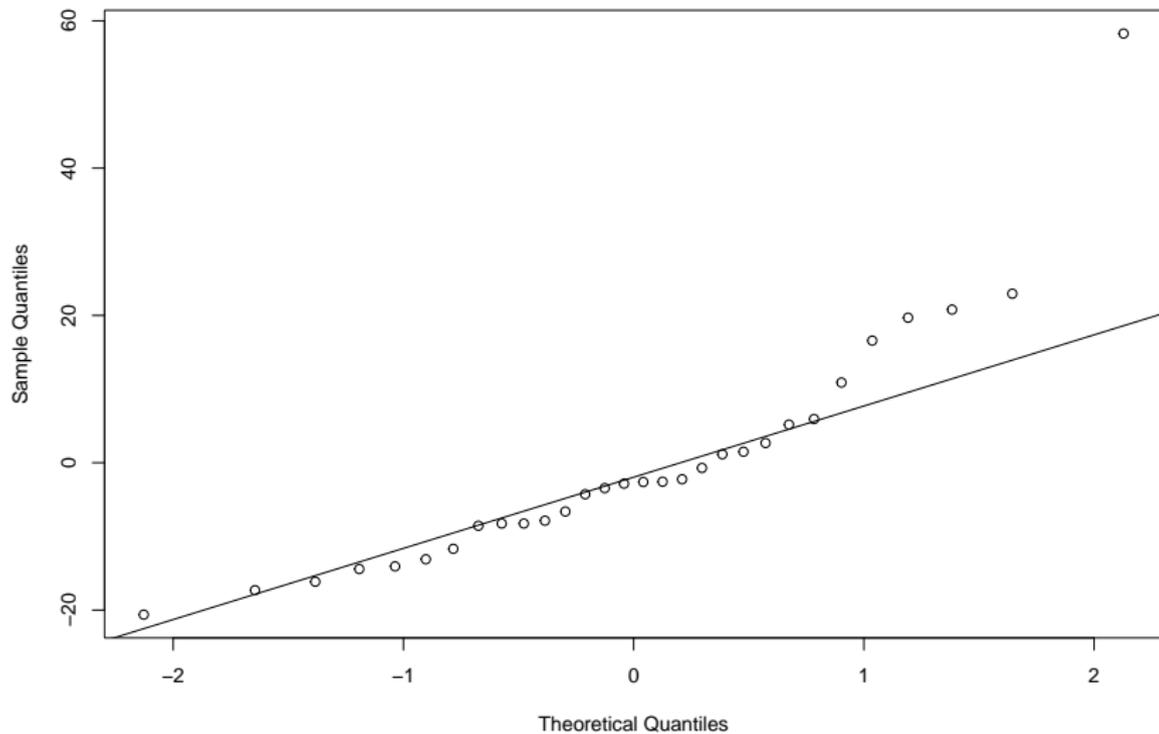
- ▶ Model diagnostics.

There are a variety of model diagnostics test for different model assumptions. Due to time constraint, we will only cover two simple tools:

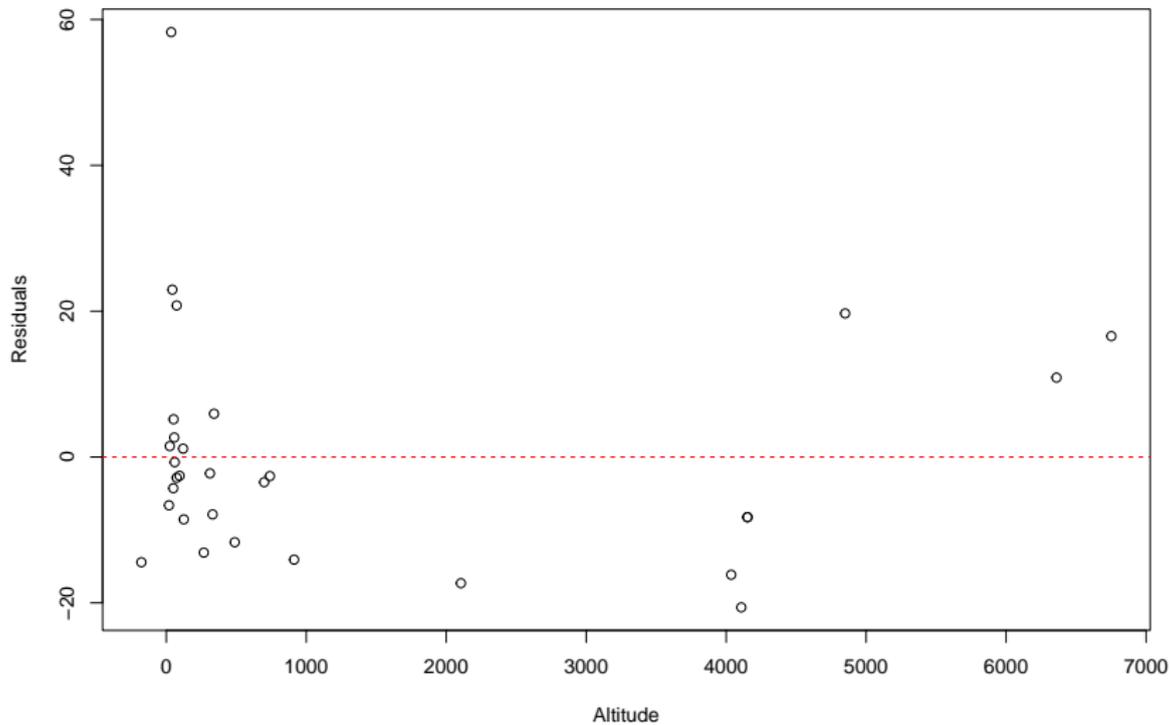
- ▶ **quantile-quantile plot** (QQ-plot). When the assumptions of residuals normality is met, we expect the points to lie on a straight line.
- ▶ **residuals against the explanatory variables**. When the assumption of independent error is met, we expect the points to scatter randomly around the the horizontal line $y = 0$.

7.2 Simple Linear Model: Model diagnostics – QQ plot

Normal Q-Q Plot



7.2 Simple Linear Model: Model diagnostics – Residual against explanatory



7.2 Simple Linear Model: Output and interpretation

```
summary(model)
```

Call:

```
lm(formula = PRECIP ~ ALTITUDE, data = rain_df)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-20.620	-8.479	-2.729	4.555	58.271

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	16.514799	3.539141	4.666	6.9e-05 ***
ALTITUDE	0.002394	0.001428	1.676	0.105

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 16.13 on 28 degrees of freedom

Multiple R-squared: 0.09121, Adjusted R-squared: 0.05875

F-statistic: 2.81 on 1 and 28 DF, p-value: 0.1048

- ▶ There is no sufficient evidence that ALTITUDE affects precipitation, $\beta = 0.00$, $t(28) = 1.676$, $p = 0.105$.

8. Model fitting procedure

The MIDI steps of data analysis

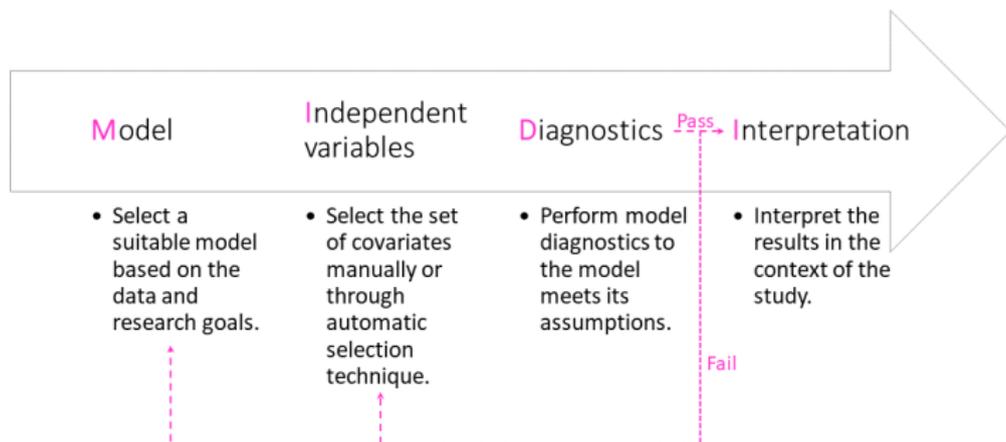


Figure 1: Recommended steps to data analysis

9. Our advice

- ▶ Visualize your data with suitable graphs.
- ▶ Check the model assumptions.
- ▶ If the normality assumption is not reasonable, there are other options available.
- ▶ Avoid extrapolating.
- ▶ Correlation does not imply causation.
- ▶ Every model has its strengths and limitations. When in doubt, get help. The SCSRU offers free 1-1 consultation to all UWaterloo researchers.

"All models are wrong, but some are useful." — George Box

10. Next steps

Now that you have reviewed how to perform some common statistics with R, you can learn more about the different models, or explore other models such as

- ▶ exploratory data analysis,
- ▶ linear regression with R,
- ▶ generalized linear models,
- ▶ count data analysis, etc.

The SCSRU organize similar workshops to this on a regular basis to improve quality of research and data literacy among the UWaterloo community. We also provide 1-1 free consultation to all researchers on campus. More information are available on our website.

Thank you!

The Statistical Consulting and Survey Research Unit (SCSRU) is the unit through which the Department of Statistics and Actuarial Science provides statistical advice to those working on research problems.