

Distance Measures and Correlation

Finding patterns in data

Kris Gunsalus

11/12/2020

Contents

Common analysis tasks and questions	1
Gene expression analysis	1
Distance measures	2
Metrics versus measures	2
Euclidean distance	2
Manhattan distance	3
Minkowski distance	4
Euclidean or Manhattan?	5
Similarity measures	5
Covariance	5
Linear correlation	6
Coefficient of determination	8
Rank-based methods	9
Examples	10
Whitlock & Schluter Resources: Correlation	10
Stream invertebrates (W&S Figure 16.4-1)	10
Exercise	14

Common analysis tasks and questions

We are often interested in identifying **large-scale patterns** in our data as a first step to develop more specific questions.

Gene expression analysis

Common **applications** to measure gene activity include:

- Time course:
 - During development
 - After exposure to some environmental stimulus (chemical, light/dark, etc.)
- Comparison between different cell types or tissues
 - Different developmental lineages
 - Normal vs. diseased tissue

We usually want to ask different types of **questions** of the same dataset:

- Identify up- and down-regulated genes.
 - *Example:* Genes that change in expression in diseased vs. normal cells, or in response to an experimental treatment.
- Find groups of genes with similar expression profiles.
 - *Example:* How many genes respond to circadian rhythms? (~6% in *Arabidopsis*)

- *Example:* Group yeast genes based on their response to a heat shock, to look for regulatory elements shared by genes with similar response patterns
- Find groups of experiments / conditions with similar expression profiles.
 - *Example:* Group cancer patients based on their gene expression profiles, to refine diagnosis and/or treatment choices.
- Find genes that explain observed differences among conditions (feature selection).
 - *Example:* Identify sentinel genes for different cell types.

To begin analyzing our data, we need some way to **group** genes and conditions based on **similarities** in gene expression profiles. The first step is to find a way to **measure distances** between genes and groups.

Distance measures

Distance measures between two variables are commonly used for such applications as constructing phylogenetic trees (e.g. based on sequence divergence, a proxy for evolutionary distance) and some clustering algorithms (e.g. k-means clustering). Distance can also be used to evaluate overall differences in gene expression, for example between different samples.

There are many ways to measure distance. The most common is **Euclidean** distance, which measures the Cartesian distance between data points. On the flip side, instead of measuring how far apart things are, we can measure how similar they are. **Correlation** (which we will discuss below) provides a measure of the **association** between variables with respect to their total variation. Things that are very far apart (large distance) have a very low correlation, and things that are very similar have a high correlation and a very small distance between them.

Metrics versus measures

We are all familiar with Euclidean distance, i.e. the shortest path between two points. It's probably a good time to note that Euclidean distance qualifies as a bona fide distance *metric*. A **metric** satisfies four conditions:

1. **Non-negativity**
 - It is always equal to or greater than zero: $d(p, q) \geq 0$.
2. **Identity**
 - If $d(p, q) = 0$, then p and q are identical: $d(p, q) = 0 \Leftrightarrow p = q$.
3. **Symmetry**
 - It is symmetric, i.e. the distance from p to q always equals the distance from q to p : $d(p, q) = d(q, p)$.
4. **Triangle inequality**
 - For three points, $d(x, z) \leq d(x, y) + d(y, z)$.

Not all distance measures satisfy these criteria. Such measures may be called *semi-metric* or *quasi-metric*, or simply *measures*. You should keep this in mind for future reference.

Euclidean distance

The most common and familiar distance metric; it gives the **root mean squared error (RMSE)**.

Euclidean distance in two dimensions Euclidean distance is simply the geometric distance between the two points in Cartesian space. Recall the basic Pythagorean formula:

$$a^2 + b^2 = c^2$$

where a and b are sides of a right triangle and c is the hypotenuse. Specifically, a^2 is the difference between the two points in the first dimension squared, and b^2 is the difference in the second dimension squared. To get the Euclidean distance c we take the $\sqrt{c^2}$, or simply $\sqrt{a^2 + b^2}$.

To calculate the distance between any two points, let's call our data points \mathbf{p} and \mathbf{q} , and let's say we have two measurements for each of them, in dimensions D_1 and D_2 (for example, weight and height, which gives body mass index). Then, $\mathbf{p} = (p_1, p_2)$ and $\mathbf{q} = (q_1, q_2)$, and we can write the distance between them, d_{pq} , as:

$$d_{pq} = d_{qp} = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2}$$

Figure 1 shows an illustration of Euclidian distance in two dimensions.

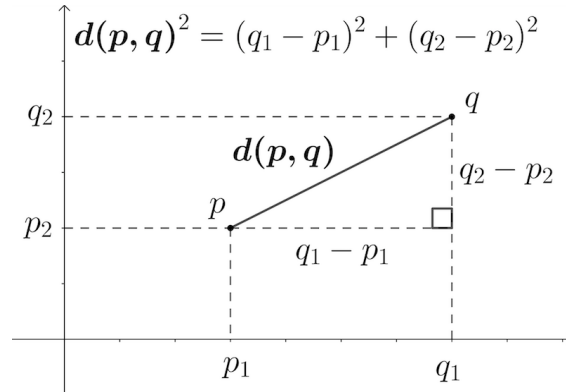


Figure 1: Euclidian distance between two points

Euclidean distance in multiple dimensions To make this statement more general, we can write $\mathbf{p} = (p_1, p_2)$ and $\mathbf{q} = (q_1, q_2)$, and express the Euclidean distance as follows:

$$d_{pq} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

where $n = 2$, so $i = (1, 2)$.

The Euclidean distance is easily generalized to multiple dimensions using the above formula. Now i takes on values 1 to n , where n is the number of dimensions: $i = (1 \dots n)$. For three dimensions, $n = 3$ (Figure 2).

Random variables are more commonly denoted using letters toward the end of the alphabet (letters in the middle tend to be used for things like the number of rows or columns in a matrix). So, instead of \mathbf{p} and \mathbf{q} , let us denote our two data points as $\mathbf{x} = (x_1, \dots, x_n)$ and $\mathbf{y} = (y_1, \dots, y_n)$. Now we can write:

$$d(x, y) = \sqrt{\sum_{i=1}^n (y_i - x_i)^2}$$

Most of the time we will omit writing out the limits of the summation explicitly, and instead just write \sum .

It will eventually become natural for you to think of these as *vectors* in n dimensions, or *n-dimensional vectors*. Vector and matrix algebra are the workhorses of *linear algebra*, which is foundational for *big data* analytics. This has become so popular a field of study that it now has its own name: *data science*.

Manhattan distance

- a.k.a. *taxicab metric*, *grid distance*, and *rectilinear distance*
- Gives the *mean absolute error* (MAE)

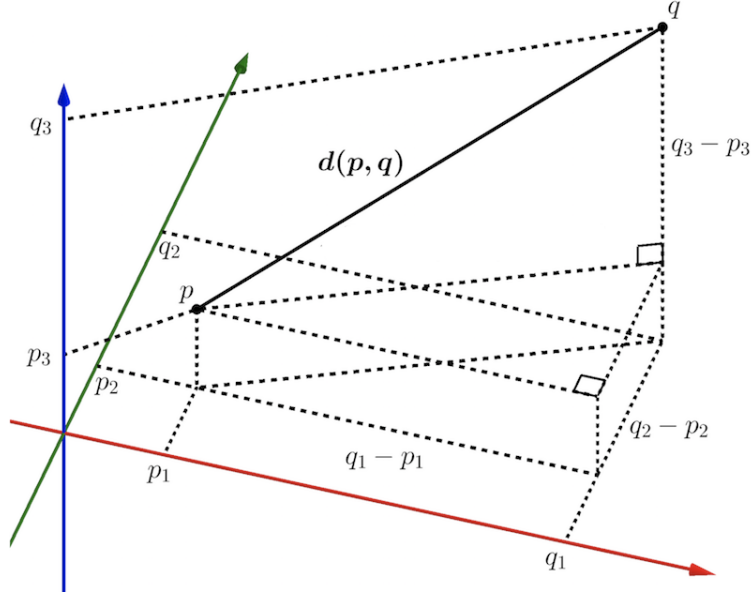


Figure 2: Euclidian distance in 3D

The Manhattan distance gets its name because it uses a “city block” approach to measuring distance. It is calculated using absolute distances in orthogonal dimensions (rather than squared distances):

$$d(x, y) = \sum_{i=1}^n |x_i - y_i|$$



Figure 3: Manhattan distance

Minkowski distance

The above two distance measures belong to the **Minkowski family** of distance measures. The Minkowski distance is defined, for vectors **a** and **b** of length p in vector space R^p , as:

$$d_{Minkowski}(\mathbf{a}, \mathbf{b}) = \left[\sum_{i=1}^p |a_i - b_i|^r \right]^{1/r}$$

By substitution, we can see that Euclidean distance corresponds to the special case of $r = 2$ and the Manhattan distance to $r = 1$.

Euclidean or Manhattan?

Most of the time you will want to use Euclidean distance. It turns out that Manhattan distance is preferable when *dimensions are not comparable*, particularly for high-dimensional data. It is also less sensitive to *outliers* because distances are measured in absolute magnitude rather than squared magnitude. In fact, reducing the exponent to $r < 1$ can even be an improvement over Manhattan distance when working in high dimensions. However if data are distributed approximately normally, so that extreme values are uncommon, then Euclidean distance usually works well.

A short digression (you can safely ignore this section, but read on if you are interested):

If you continue in data science you will probably hear more about *L-norms* (based on Lebesgue spaces). In linear algebra, a *norm* is a function that assigns a positive length or size to any non-zero vector.

- For **L2 statistics** (based on Euclidean distance), the *mean* is the value that minimizes the *mean squared error* across the dataset.
- For **L1 statistics** (based on Manhattan distance), the *median* is the value that minimizes the sum of *absolute deviations* across the dataset.
- A family of **regularization** techniques for multivariate regression problems is often applied to large datasets to simplify computational solutions by applying penalties to a solution's parameter values (don't worry if this sounds like a lot of gobbledy-gook for now; you may find it useful later).
 - *L1 regularization* (like the LASSO method) encourage solutions with a lot of zero values, whereas
 - *L2 regularization* (used in ridge regression) encourage solutions where many values are very small.
 - *Elastic net regularization* uses a combination of L1 and L2 norms.

Similarity measures

It is often useful to measure the similarity, or **association**, between two variables instead of how far apart they are.

Covariance

In order to look at how much two variables change in tandem, we can calculate the **covariance**. Just as the **variance** of a single variable represents its **dispersion**:

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1} = \frac{\sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{n - 1}$$

so the **covariance** measures how much two values vary **together**:

$$cov_{x,y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

Now, instead of taking the square of a single variable to measure its deviation from the mean, we now multiply the difference between two variables and their respective means.

We divide by $n - 1$ because for n data points, we have $n - 1$ **independent** observations, or degrees of freedom (and not, for example, $(n - 2)$ or $(n - 1)^2$). *It may help you to think of the degrees of freedom as belonging to the dimensions of the system, not as belonging to each variable within it.*

Linear correlation¹

For some applications, in order to make comparisons between datasets, it's helpful to get them all on the same scale. We can do this by standardizing, or normalizing, the data.

The linear correlation coefficient measures the **strength** and the **direction** of a linear relationship between two variables. This is sometimes referred to as the *Pearson product moment correlation coefficient*, or more simply the **Pearson correlation coefficient**, in honor of its developer Karl Pearson.

Recall that for a series of observations in **one dimension** that are normally distributed, the data can be *normalized* by subtracting the mean from each data point and dividing by the standard deviation of the population. The *Z - score* is thus defined as:

$$Z = \frac{\sum x_i - \mu_x}{\sigma_x}, \text{ where } Z \sim \mathcal{N}(0, 1)$$

In two dimensions, the **correlation** of two variables in a population is defined as their **covariance** divided by the product of their **standard deviations**:

$$\rho = \frac{Cov(X, Y)}{\sigma_x \sigma_y}$$

Thus, the correlation coefficient is just the covariance, normalized using the variance of the individual variables. For n observations, r becomes:

$$\begin{aligned} r &= \frac{Cov(X, Y)}{s_X s_Y} \\ &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(n - 1)s_X s_Y} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1 \sqrt{Var(X)Var(Y)}} \\ &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \end{aligned}$$

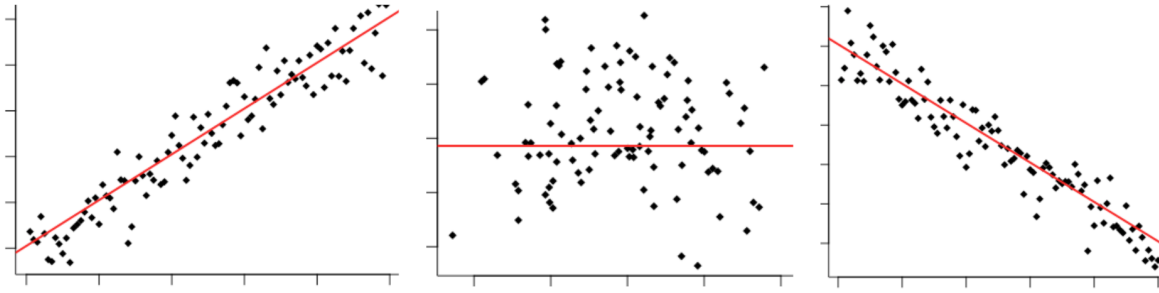
This last expression is the form of the correlation coefficient given in your book, which is the sum of products divided by the individual sums of squares. Notice that if $y = x$, then we can write:

$$r = \frac{\sum (x_i - \bar{x})(x_i - \bar{x})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (x_i - \bar{x})^2}} = 1$$

Similarly, if $y = -x$, then $r = -1$. So r ranges from -1 to $+1$.

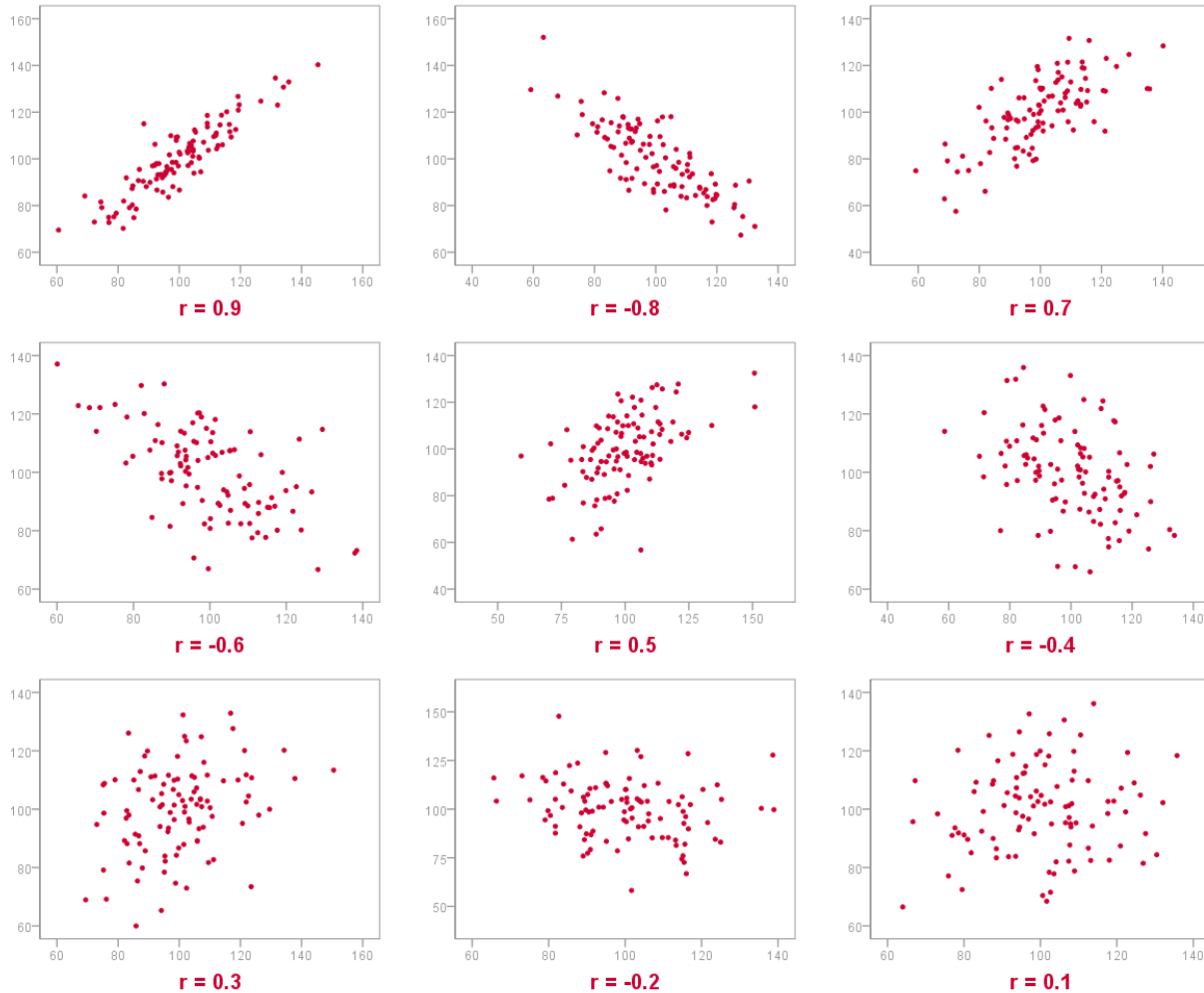
Properties of r The value of r is such that $-1 \leq r \leq +1$. The $+$ and $-$ signs signify positive and negative linear correlations, respectively. The figure below illustrate data with a strong positive correlation (left), uncorrelated data (middle), and data with a strong negative correlation (right).

¹Based in part on <https://mathbits.com/MathBits/TISection/Statistics2/correlation.htm>



- **Positive correlation:** If x and y have a strong positive linear correlation, r is close to $+1$. An r value of exactly $+1$ indicates a perfect positive fit. Positive values indicate a relationship between x and y variables such that as x increases, y also increases.
- **Negative correlation:** If x and y have a strong negative linear correlation, r is close to -1 . An r value of exactly -1 indicates a perfect negative fit. Negative values indicate a relationship between x and y such that as values for x increase, values for y decrease.
- **No correlation:** If there is no linear correlation or a weak linear correlation, r is close to 0 . A value near zero means that there no relationship between the two variables.
- Note that r is a dimensionless quantity; that is, it does not depend on the units employed.
- A perfect correlation of ± 1 occurs only when the data points all lie exactly on a straight line. If $r = +1$, the slope of this line is positive. If $r = -1$, the slope of this line is negative.
- A correlation greater than 0.8 is generally described as strong, whereas a correlation less than 0.5 is generally described as weak. These values can vary based upon the “type” of data being examined. A study utilizing scientific data may require a stronger correlation than a study using social science data.

The following diagram illustrates data showing a range of correlation coefficients:



Interactive visualizations

- Whitlock & Schluter have an **interactive exercise** online that can help you get a feel for correlation coefficients for point clouds with different shapes.
- Another one may be found [here](#).

Pearson correlation distance The Pearson correlation can be turned into a distance measure simply by subtracting it from 1:

$$d = 1 - r$$

Variables with high correlation have a short distance between them, and vice versa. So if two genes have a perfect positive correlation, i.e. $r = 1$, then the distance between them is $d = 1 - 1 = 0$. For a perfect negative correlation, the distance will be $d = 1 - (-1) = 2$.

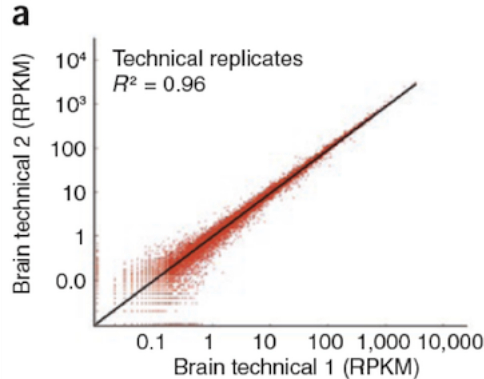
Coefficient of determination²

The **coefficient of determination** is denoted as R^2 . For the special case of two variables, R^2 is equivalent to r^2 – the square of the linear correlation coefficient (we will come back to this shortly when we discuss linear regression). As such, it varies between 0 and 1.

²From <https://mathbits.com/MathBits/TISection/Statistics2/correlation.htm>

The coefficient of determination is useful because it gives the **proportion of the variance** (fluctuation) of one variable that is predictable from the other variable. It is a measure that allows us to determine how certain one can be in making predictions using a particular model.

Notice that the following graph, showing the RPKM for two replicate samples, gives their R^2 value. This is a measure of the reproducibility of the data.



Properties of r^2

- r^2 is the ratio of the explained variation to the total variation.
- r^2 is such that $0 < r^2 < 1$, and denotes the strength of the linear association between \mathbf{x} and \mathbf{y} .
- r^2 represents the percent of the data that is explained by the line of best fit. For example, if $r = 0.922$, then $r^2 = 0.850$, which means that 85% of the total variation in \mathbf{y} can be explained by the linear relationship between \mathbf{x} and \mathbf{y} (as described by the regression equation, which we will discuss later). The other 15% of the total variation in \mathbf{y} remains unexplained.
- r^2 is a measure of how well the regression line represents the data. If the regression line passes exactly through every point on the scatter plot, it would be able to explain all of the variation. The further the line is away from the points, the less it is able to explain.

Rank-based methods

In order to use a linear correlation coefficient, data should be randomly sampled from a population and should follow a **bivariate normal distribution** with the following features:

- Individual variables are normally distributed
- Variables show a linear relationship

When these do not hold, it may be possible to make data more normal using a transformation such as taking the log, square, or arcsin. Alternatively, non-parametric (rank-based) methods may be applied:

Spearman's rho (r_s) Whitlock & Schluter explains Spearman's rho and how it is used for significance testing. Very briefly, features of Spearman's rho are:

- Assumes a **monotonic** relationship between random samples (variables increase or decrease together, or one does not change)
- Correlation test uses a t -test, and so incorporates variation (standard error) of r_s

Kendall's tau (τ) An alternative rank-based measure that you will come across is **Kendall's τ** . The **Wikipedia** entry for Kendall's tau has a good explanation of how it is computed and how it is used for significance testing, so I won't go into the details here, but just give a very brief overview:

- Measures difference between **concordant** (the ranks are the same) vs. *discordant* (different ranks) pairs of data points, normalized by the binomial coefficient (the **total number of pairwise combinations**)
- Values usually smaller than Spearman's rho
- Advantages: Insensitive to error, and gives more accurate p-values for small sample sizes

Spearman's rho and Kendall's tau often give similar interpretations of significance. **Both measures vary between -1 and 1**, and so intuitively are similar to the linear correlation coefficient.

Examples

Whitlock & Schluter Resources: Correlation

- R code for the Chapter 16 **examples**
- **Lab 11**: Correlation and regression
- **Interactive demo** for correlation

Stream invertebrates (W&S Figure 16.4-1)

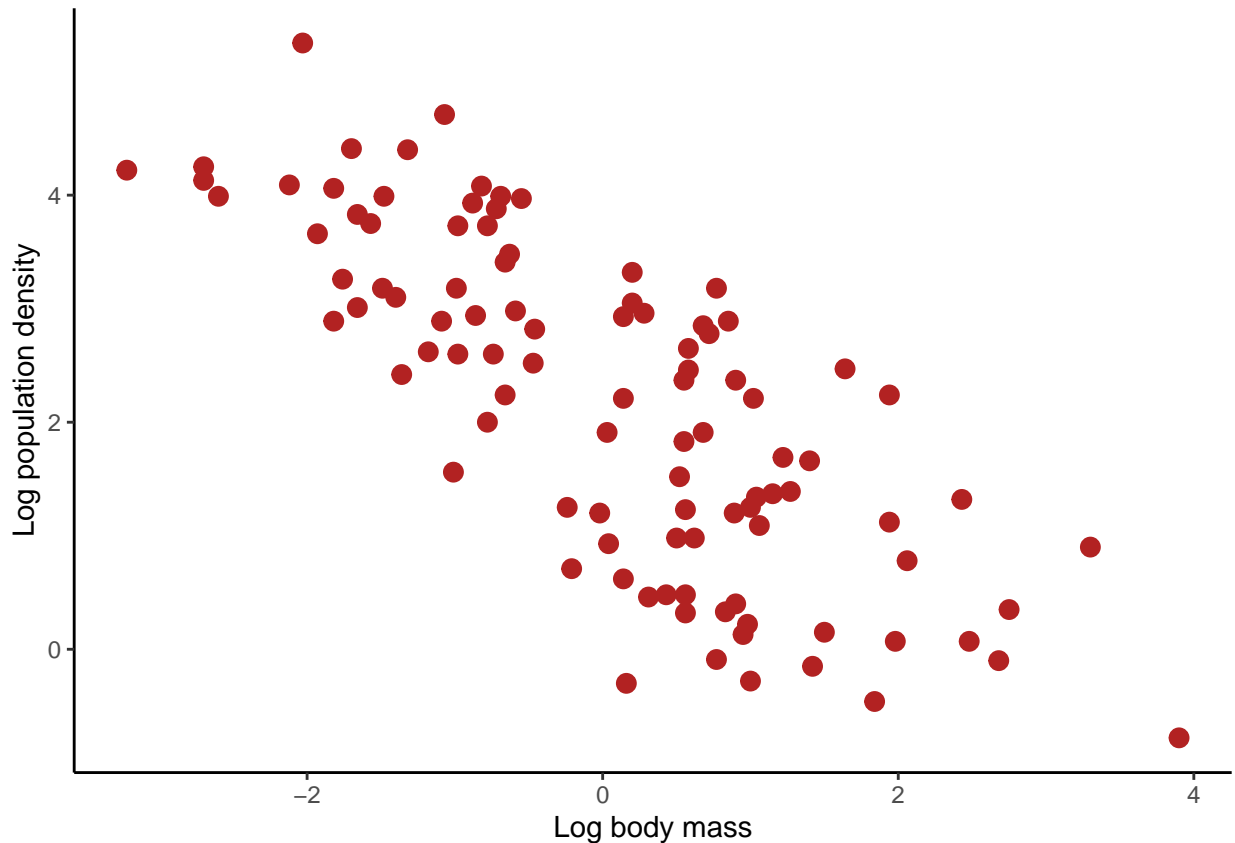
R contains several functions for covariance and correlation. We will use the example from **Figure 16.4-1** to illustrate these. This example shows a negative linear correlation between log-transformed data on the population density of different invertebrates in streams and their body mass (Schmid et al., 2000).

```
streamInvert <- read.csv(url("https://whitlockschluter3e.zoology.ubc.ca/Data/chapter16/chap16f4_1StreamInvertebrates.csv"))
head(streamInvert)
```

Load dataset

```
##   log10Mass log10Density
## 1      -3.22         4.22
## 2      -2.70         4.25
## 3      -2.70         4.13
## 4      -2.60         3.99
## 5      -2.03         5.34
## 6      -2.12         4.09
```

```
ggplot(streamInvert, aes(log10Mass, log10Density)) +
  geom_point(size = 3, col = "firebrick") +
  labs(x = "Log body mass", y = "Log population density") +
  theme_classic()
```



Scatterplot

Correlation The correlation coefficient is calculated with the `cor()` function. The default method is to compute the Pearson correlation coefficient, r . The rank-based methods Spearman's rho (r_s) and Kendall's tau (τ) may be specified as an alternative. All of these vary between -1 and 1. Notice that the magnitude of Kendall's tau is indeed smaller than Spearman's rho.

```
cor(streamInvert$log10Density, streamInvert$log10Mass, method = "pearson")
```

```
## [1] -0.7651667
```

```
cor(streamInvert$log10Density, streamInvert$log10Mass, method = "spearman")
```

```
## [1] -0.7675654
```

```
cor(streamInvert$log10Density, streamInvert$log10Mass, method = "kendall")
```

```
## [1] -0.5643002
```

Correlation test Hypothesis testing for correlation allows quantification of statistical significance. The `cor.test()` function can be applied to different types of correlation measures, but it only gives an approximate confidence interval for Pearson's correlation.

```
# test for significance
```

```
cor.test(~ log10Density + log10Mass, data = streamInvert)
```

```
##
```

```
## Pearson's product-moment correlation
```

```
##
```

```
## data: log10Density and log10Mass
```

```
## t = -11.765, df = 98, p-value < 2.2e-16
```

```
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.8359480 -0.6693464
## sample estimates:
##      cor
## -0.7651667

cor.test(~ log10Density + log10Mass, data = streamInvert, method="spearman")

## Warning in cor.test.default(x = c(4.22, 4.25, 4.13, 3.99, 5.34, 4.09, 4.41, :
## Cannot compute exact p-value with ties

##
## Spearman's rank correlation rho
##
## data: log10Density and log10Mass
## S = 294565, p-value < 2.2e-16
## alternative hypothesis: true rho is not equal to 0
## sample estimates:
##      rho
## -0.7675654

cor.test(~ log10Density + log10Mass, data = streamInvert, method="kendall")

##
## Kendall's rank correlation tau
##
## data: log10Density and log10Mass
## z = -8.2867, p-value < 2.2e-16
## alternative hypothesis: true tau is not equal to 0
## sample estimates:
##      tau
## -0.5643002

# just the correlation coefficient
cor.test(~ log10Density + log10Mass, data = streamInvert)$estimate

##      cor
## -0.7651667
```

Covariance Correlation allows the comparison of relative changes for variables with different magnitudes, since it is standardized relative to the total variation of the individual variables.

Covariance also measures the strength of association between two variables, though it is often less useful, particularly when data are not on the same scale or their ranges are of very different magnitudes. Notice that covariance computed by the different methods gives values of very different magnitudes.

```
cov(streamInvert$log10Density, streamInvert$log10Mass, method = "pearson")

## [1] -1.546049

cov(streamInvert$log10Density, streamInvert$log10Mass, method = "spearman")

## [1] -645.947

cov(streamInvert$log10Density, streamInvert$log10Mass, method = "kendall")

## [1] -5564
```

Correlograms The `ggpairs` function from the `GGally` package also provides a nice snapshot of correlation between multiple columns in a data matrix and is customizable to show different representations of the data. Let's look at a simple example for illustration.

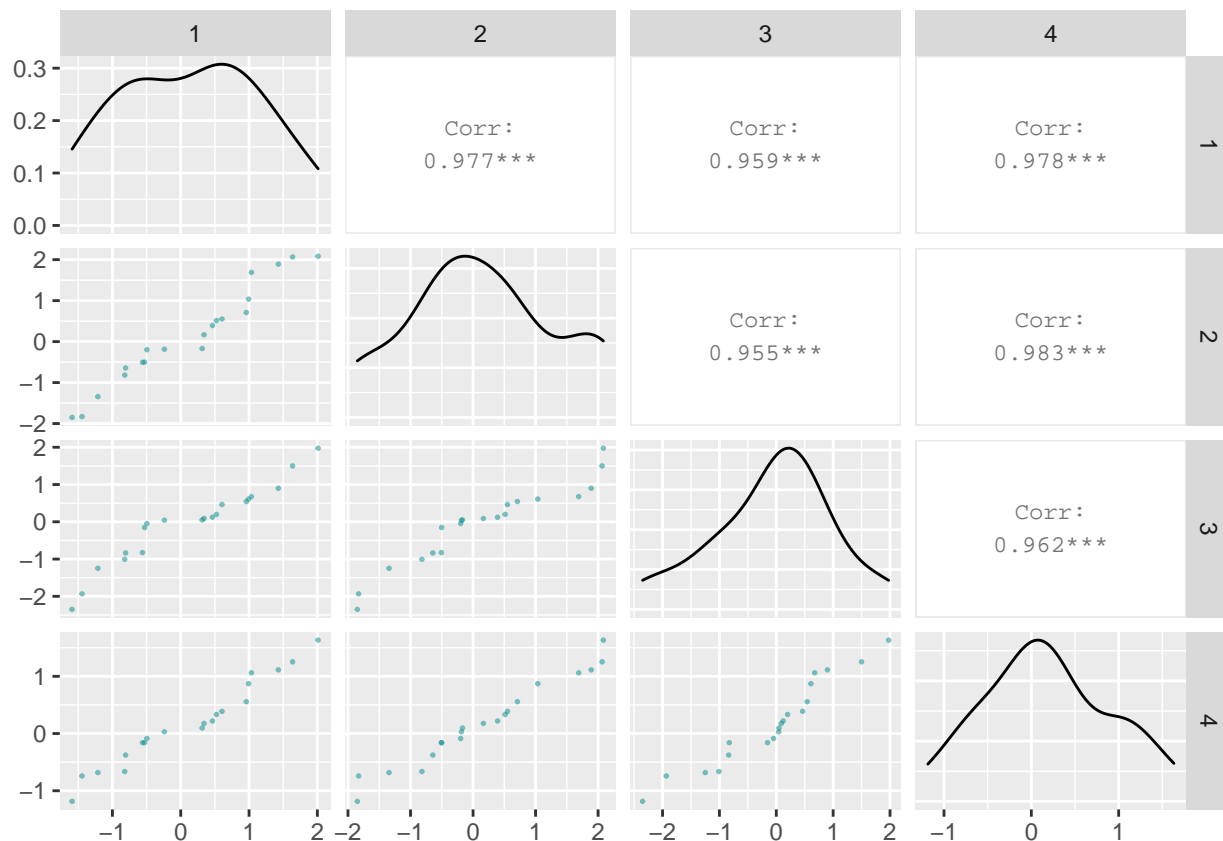
```
#install.packages("GGally")
library(GGally)

## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2

# create a simulated dataset of experimental replicates
exp.reps = cbind(sort(runif(20, -1, 1)),
                  sort(runif(20, -1, 1)),
                  sort(runif(20, -1, 1)),
                  sort(runif(20, -1, 1)) )

exp.reps = cbind(sort(rnorm(20, 0, 1)),
                  sort(rnorm(20, 0, 1)),
                  sort(rnorm(20, 0, 1)),
                  sort(rnorm(20, 0, 1)) )

# plot the correlogram
ggpairs(exp.reps, # log-transformed counts + 1
        upper = list(continuous = wrap("cor",size = 3)), # upper triangle: correlations
        lower = list(continuous = wrap("points",
                                       alpha=0.5,
                                       size = 0.3,
                                       color = "cyan4")) # lower triangle: scatterplot of data points
```



Heat maps You can use the `geom_tile()` geometry from `ggplot` to make a heat map of the overall distance (or similarity) between datasets, for example to compare gene expression profiles between different treatments. We will see some examples in this week's homework and in recitation next week; we will make more heatmaps when we get to clustering.

Exercise

The accompanying exercise illustrates the mechanics of computing distance and correlation, with some simple examples in R.