

ETC5521: Exploratory Data Analysis

Going beyond two variables, exploring
high dimensions

Lecturer: *Di Cook*

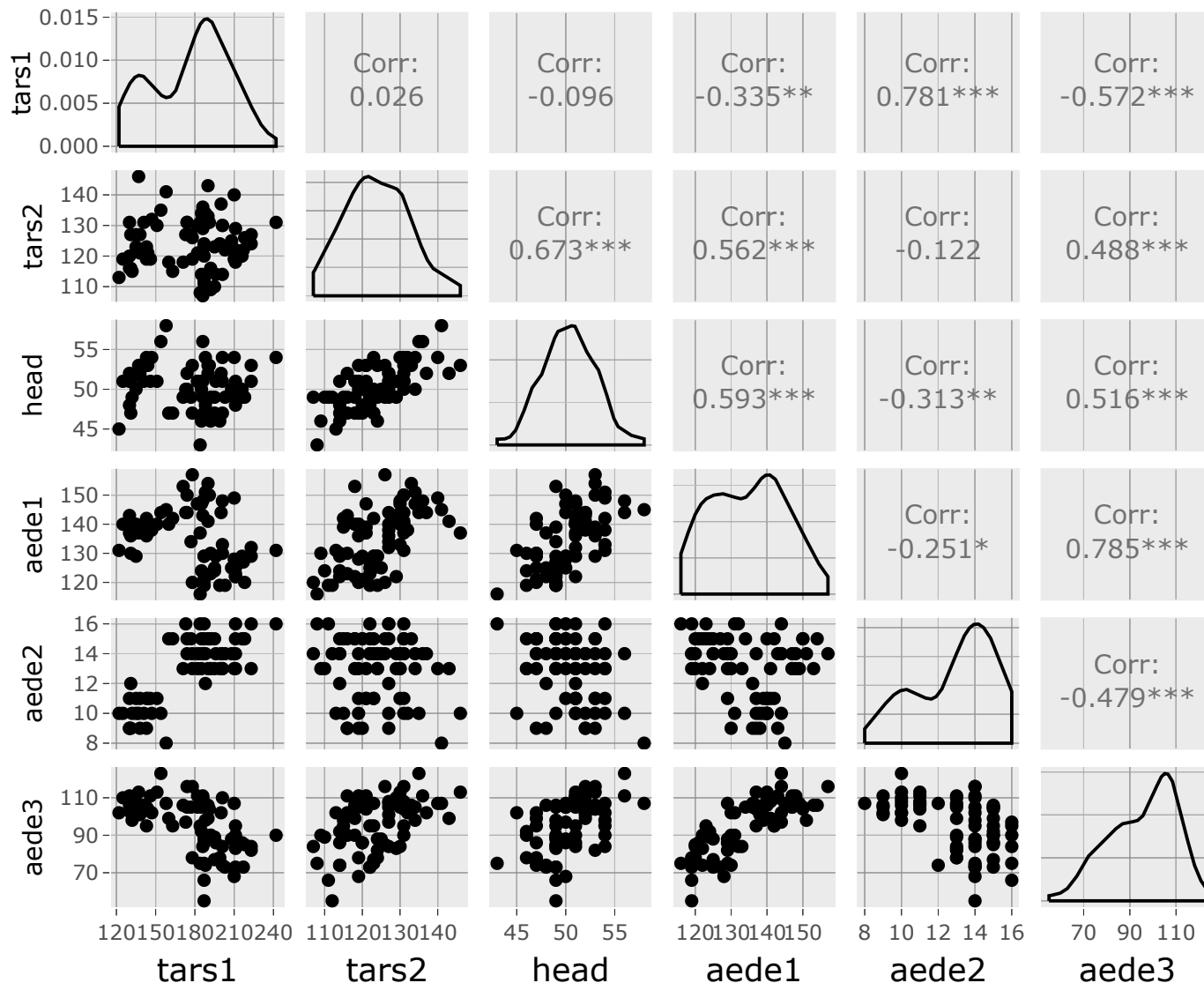
Department of Econometrics and Business Statistics

✉ ETC5521.Clayton-x@monash.edu

Week 6 - Session 2



Linking between plots



If you have interactive plots you can investigate whether

- an outlier in one or two variables is one of the extremes in other variables → then it is a multivariate outlier

- clusters of observations in one plot are concentrated in a cluster in other variables → the data is multimodal in multivariate space, and there are likely sub-populations.



Your turn, **cut and paste the code** into your R console,
and **select** regions of the resulting plot to examine where
these points lie in other plots.

```
library(tidyverse)
library(tourr)
library(plotly)
highlight_key(flea) %>%
  GGally::ggpairs(aes(color = species), columns = 1:6) %>%
  ggplotly() %>%
  highlight("plotly_selected")
```

05:00

History

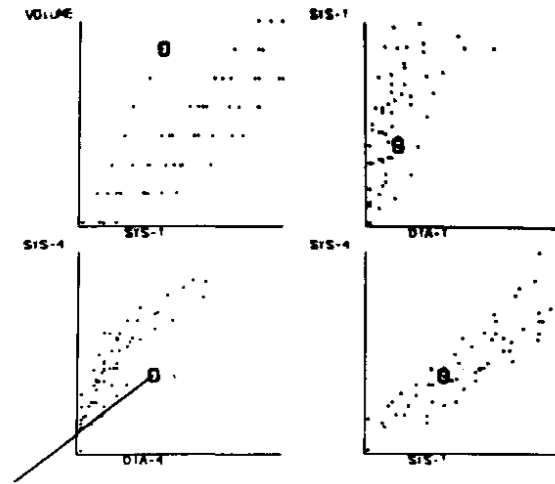
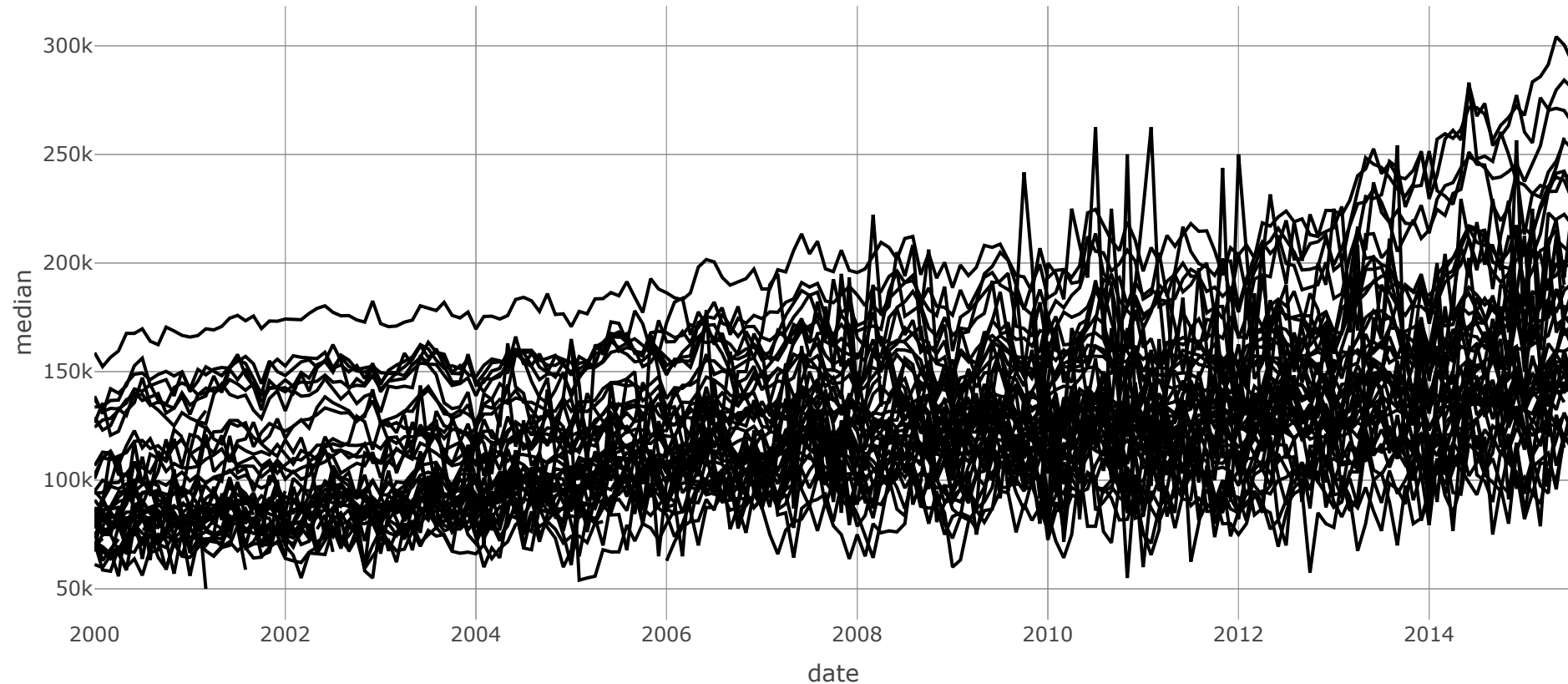


Figure 2: The user has designated a point on one graph by means of the light-pen. It and corresponding points on the other graphs then are identified by a circle, and the values of all variables for that case are as follows:

First linked brushing across a set of scatterplots done by Carol Newton (1978) "Graphics from Alpha to Omega in Data Analysis" in Graphical Representation of Multivariate Data.

Parallel coordinate plots



You can think about parallel coordinate plots like time series plots. Each variable is treated as a parallel axis. Observations are drawn as a line across the axes.

```
# load the `txhousing` dataset
data(txhousing, package = "ggplot2")

# declare `city` as the SQL 'query by' column
tx <- highlight_key(txhousing, ~city)

# initiate a plotly object
base <- plot_ly(tx, color = I("black")) %>%
  group_by(city)

# create a time series of median house price
base %>%
  group_by(city) %>%
  add_lines(x = ~date, y = ~median)
```

What can you learn?

- 🗨 lines that are parallel indicate positive linear association, possibly across many variables
- 🗨 lines that cross indicate negative linear association
- 🗨 lines that go up and down differently from any other lines indicate multivariate outliers
- 🗨 lines that go up and down together, but differently from other lines indicate multivariate clustering

*You need to have an **interactive** parallel coordinate plot for them to be effective for exploring data*



Your turn, **cut and paste the code** into your R console,
and **click** in the resulting plot to examine the line for an
observation.

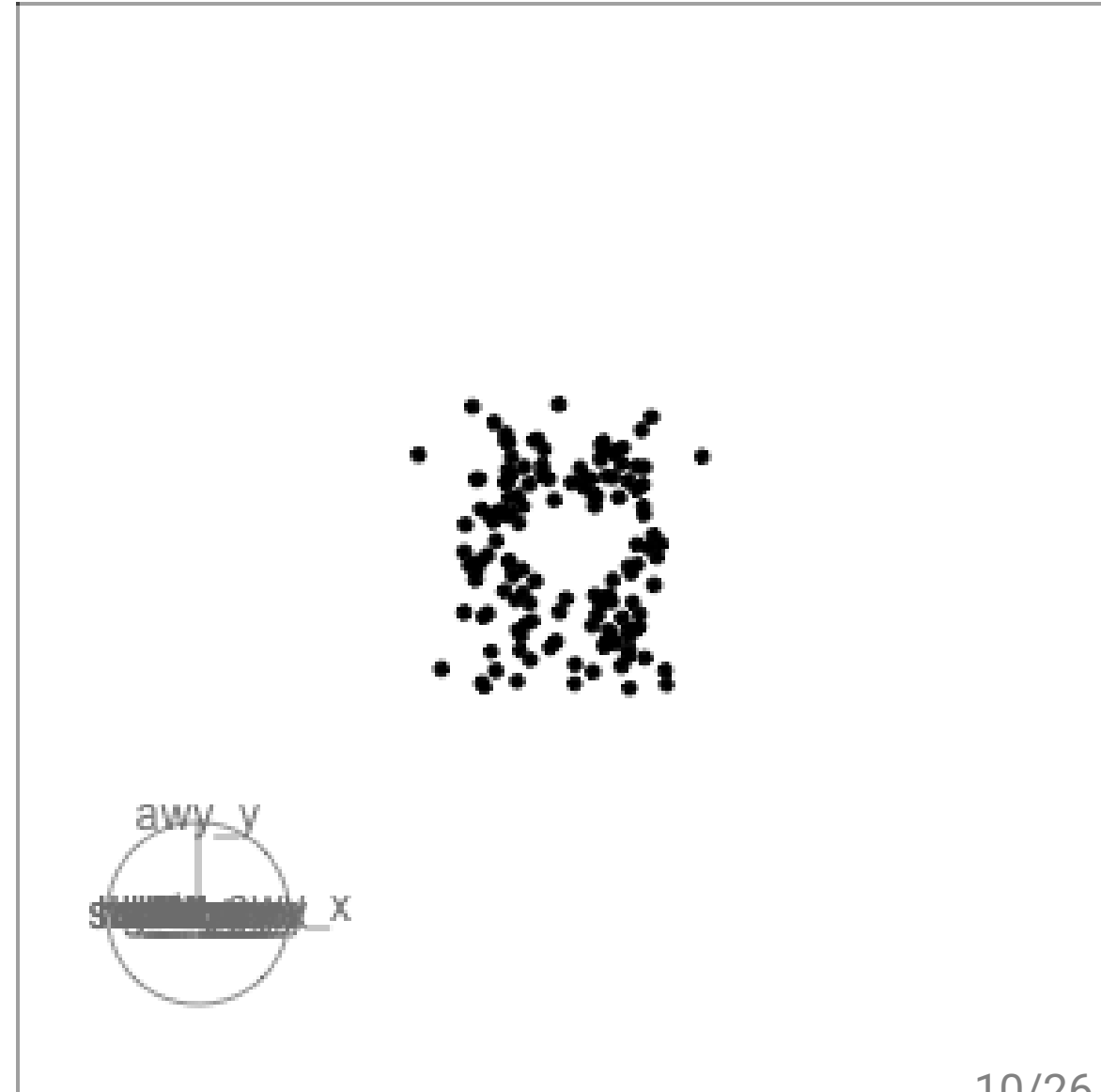
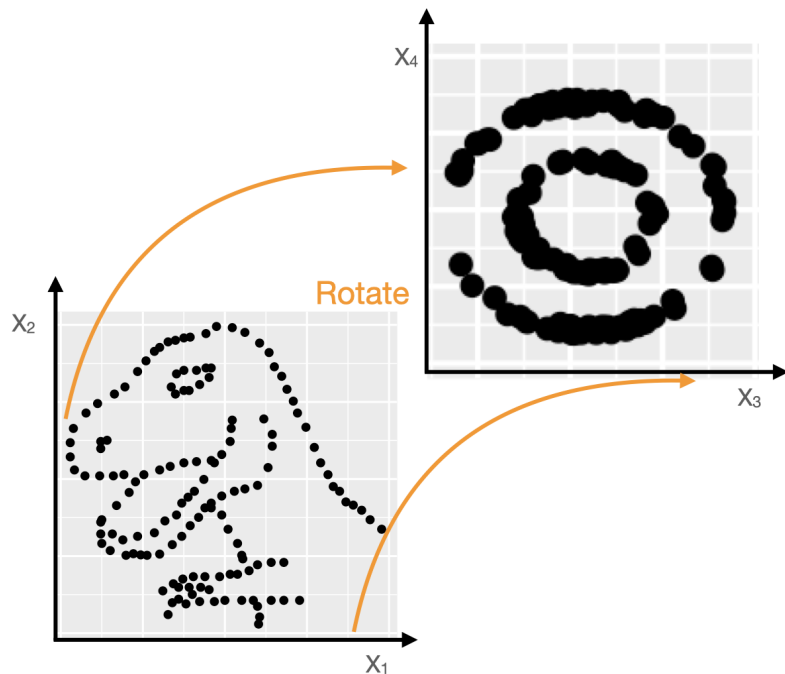
```
ui <- fluidPage(  
  plotlyOutput("parcoords"),  
  verbatimTextOutput("data")  
)  
  
server <- function(input, output, session) {  
  
  flea_numeric <- dplyr::select_if(flea, is.numeric)  
  
  output$parcoords <- renderPlotly({  
    dims <- Map(function(x, y) {  
      list(values = x, range = range(x), label = y)  
    }, flea_numeric, names(flea_numeric), USE.NAMES = FALSE)  
    plot_ly(type = 'parcoords', dimensions = dims, source = "pcoords") %>%  
      layout(margin = list(r = 30)) %>%  
      event_register("plotly_restyle")  
  })  
}
```

05:00

Touring between pairs of variables

Dinosaur data as multivariate data.

Rotate one pair of variables into the other - special type of interpolation/animation called a *little tour*



brushing



Brushing on a scatterplot matrix, Rick Becker (Becker and Cleveland, 1987 Brushing Scatterplots, Technometrics)

Tours of multivariate data

The simplest explanation of is that we look at combinations of variables, as well as the individual variables. The grand tour does this in an elegant way so that there's a chance of seeing all possible combinations of the variables, and it glues these together so that there is a smooth change from one to another. With an important note: that always when scatterplots (or higher dimensional combinations are shown) that the axes are always orthogonal.

Using combinations of variables, it can be possible to see separations between groups, outliers, linear and non-linear dependencies that were not visible in the single variable plots.

■ *Yes, 3D plots aren't good enough. You need tours.*

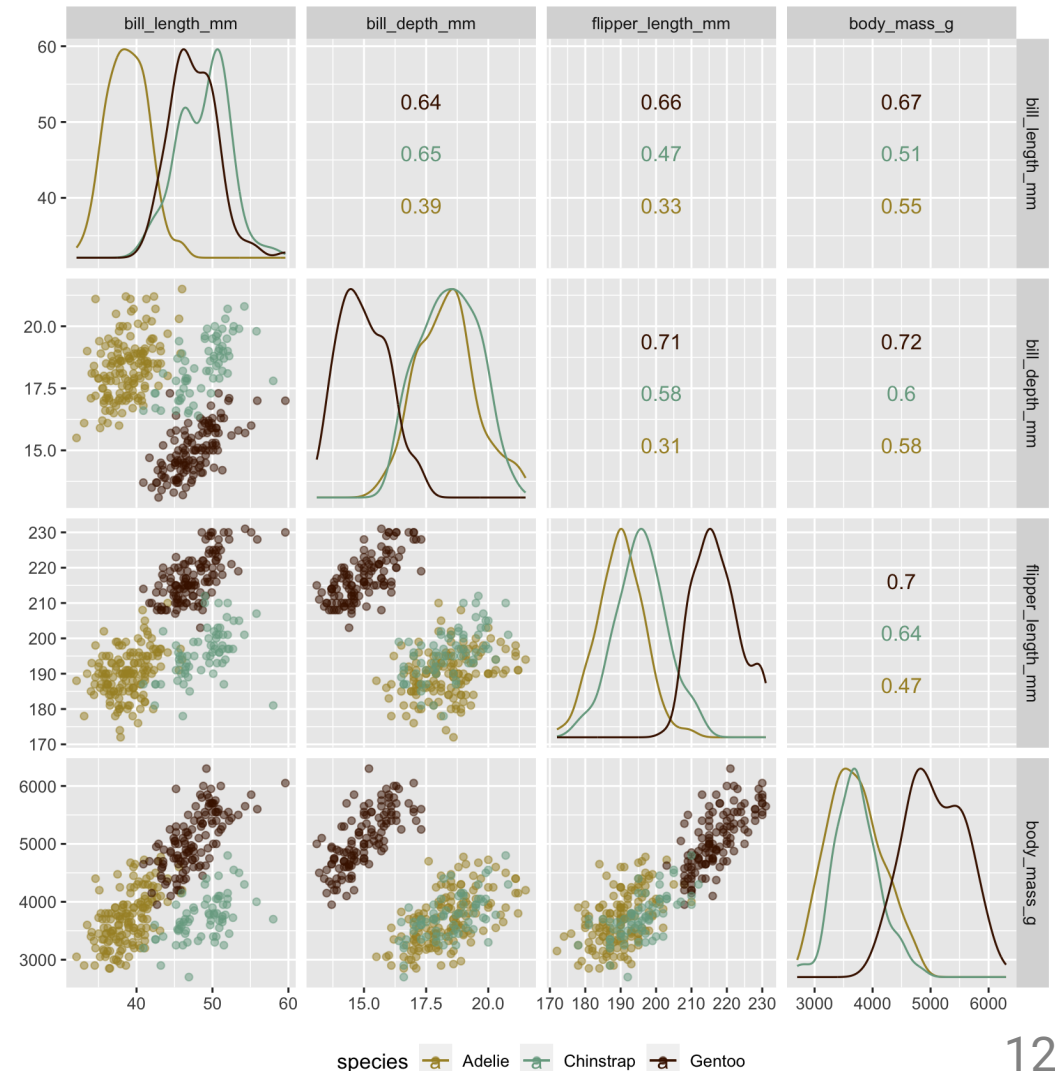
Case study 4 Penguins



Size measurements for adult foraging penguins near Palmer Station, Antarctica

- species: a factor denoting penguin species (Adélie, Chinstrap and Gentoo)
- bill_length_mm: a number denoting bill length (millimeters)
- bill_depth_mm: a number denoting bill depth (millimeters)
- flipper_length_mm: an integer denoting flipper length (millimeters)
- body_mass_g: an integer denoting body mass (grams)

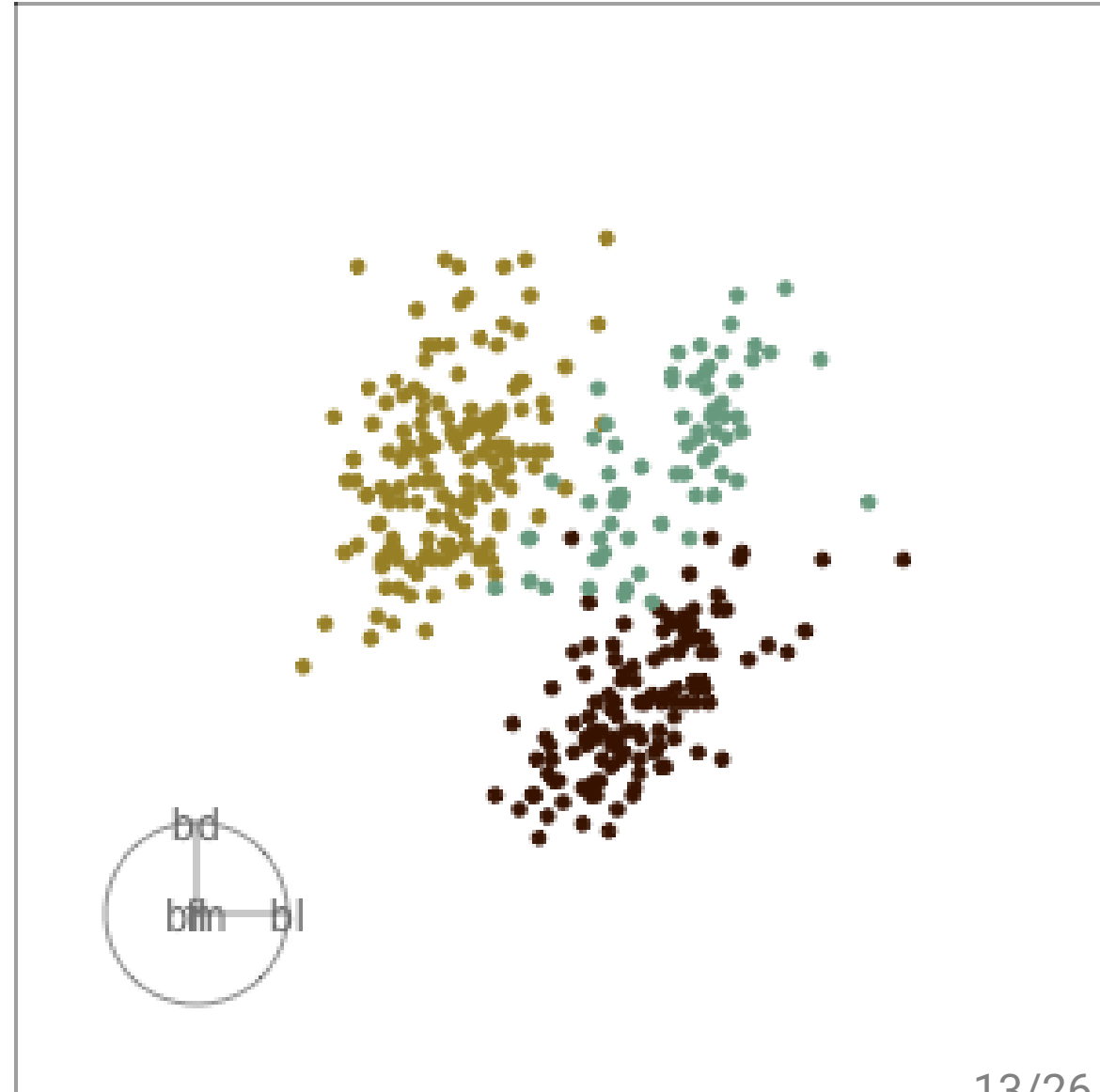
Gorman KB, Williams TD, Fraser WR (2014) Ecological Sexual Dimorphism and Environmental Variability within a Community of Antarctic Penguins (Genus *Pygoscelis*). PLoS ONE 9(3): e90081. doi:10.1371/journal.pone.0090081



Our first grand tour

Here's the code

```
clrs <- ochre_pal(  
  palette="nolan_ned")(3)  
col <- clrs[  
  as.numeric(  
    penguins$species)]  
animate_xy(penguins[,3:6],  
  col=col,  
  axes="off",  
  fps=15)
```



Case study 4 Penguins

What do you learn about this data?

💬 clusters ✓

💬 outliers ✓

💬 linear dependence ✓

💬 elliptical clusters with slightly different shapes ✓

💬 separated elliptical clusters with slightly different shapes ✓

What is a tour?

A grand tour is by definition a movie of low-dimensional projections constructed in such a way that it comes arbitrarily close to showing all possible low-dimensional projections; in other words, a grand tour is a space-filling curve in the manifold of low-dimensional projections of high-dimensional data spaces.

$\mathbf{x}_i \in \mathcal{R}^p$, i^{th} data vector

F is a $p \times d$ orthonormal matrix, $F'F = I_d$, where d is the projection dimension.

The projection of \mathbf{x}_i onto F is $\mathbf{y}_i = F'\mathbf{x}_i$.

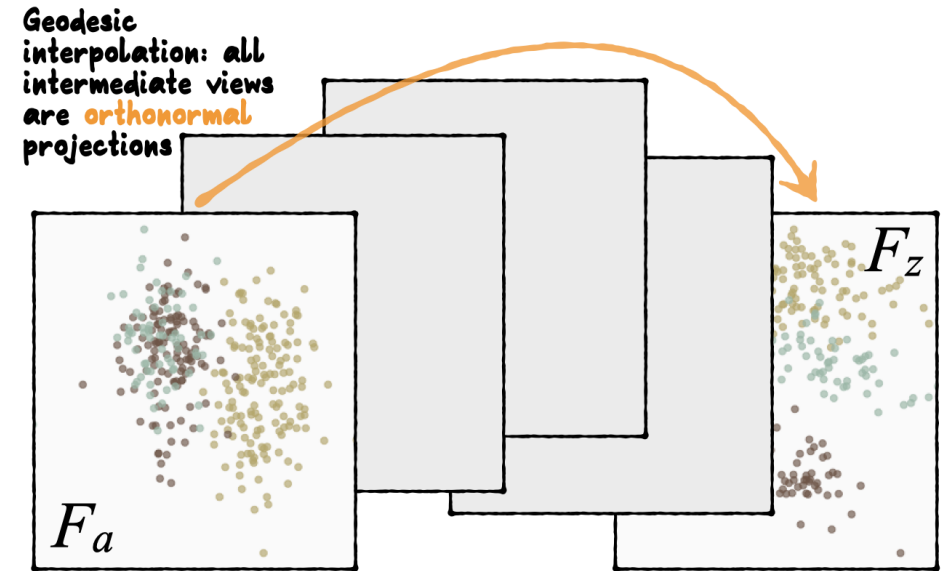
Tour is indexed by time, $F(t)$, where $t \in [a, z]$. Starting and target frame denoted as $F_a = F(a)$, $F_z = F(t)$.

The animation of the projected data is given by a path $\mathbf{y}_i(t) = F'(t)\mathbf{x}_i$.

Geodesic interpolation between planes

Tour is indexed by time, $F(t)$, where $t \in [a, z]$. Starting and target frame denoted as $F_a = F(a)$, $F_z = F(t)$.

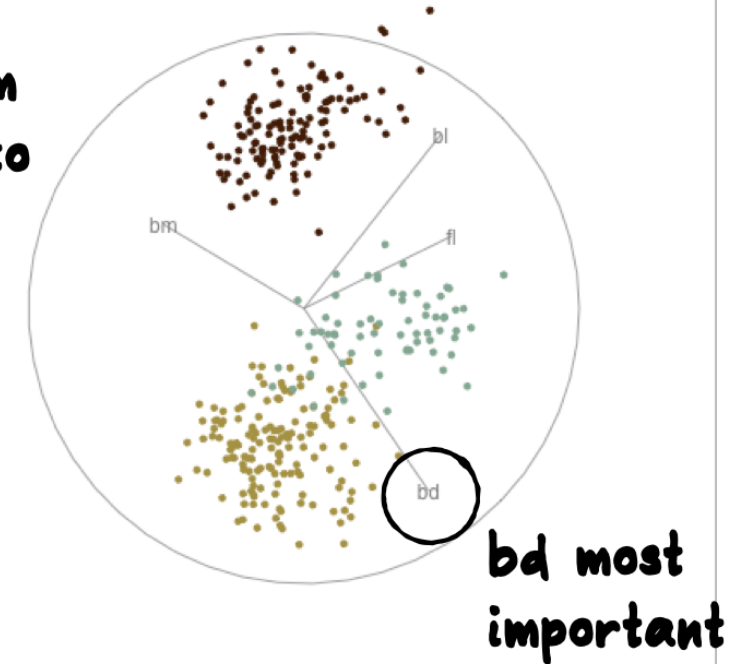
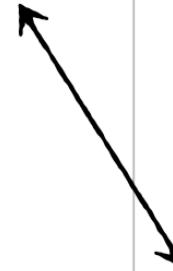
The animation of the projected data is given by a path $\mathbf{y}_i(t) = F'(t)\mathbf{x}_i$.



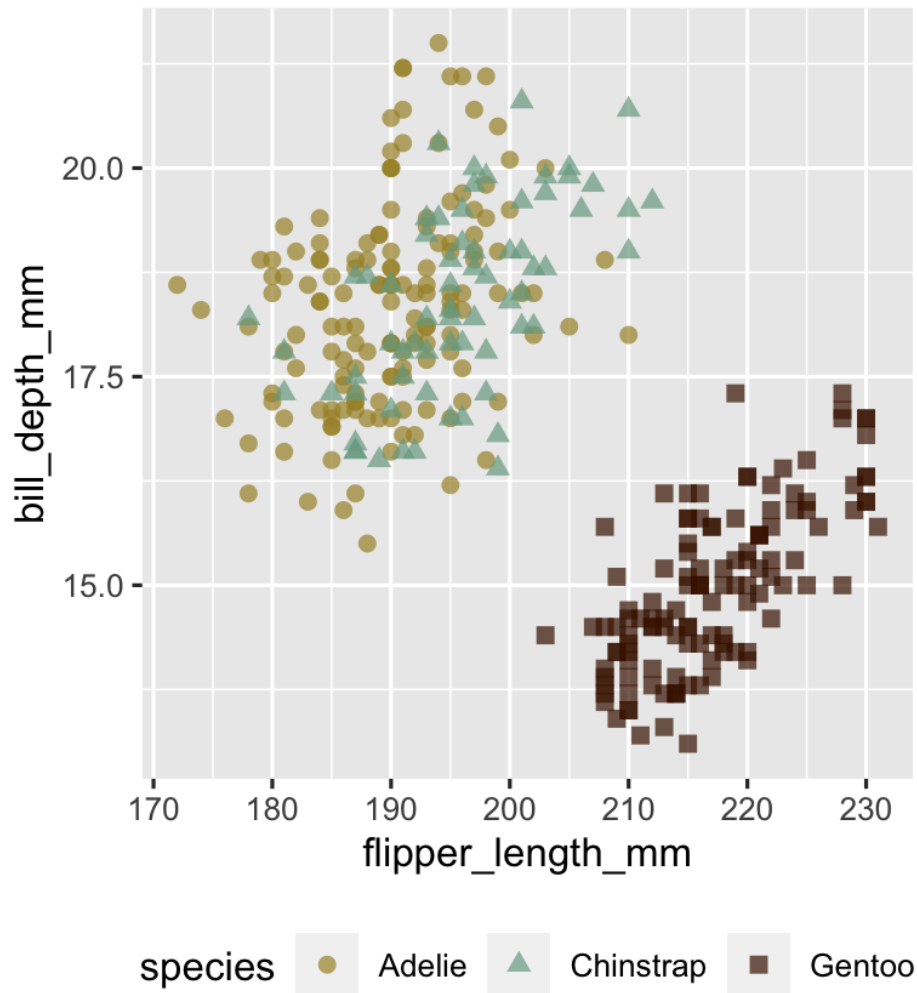
Reading axes - interpretation

Length and direction of axes relative to the pattern of interest

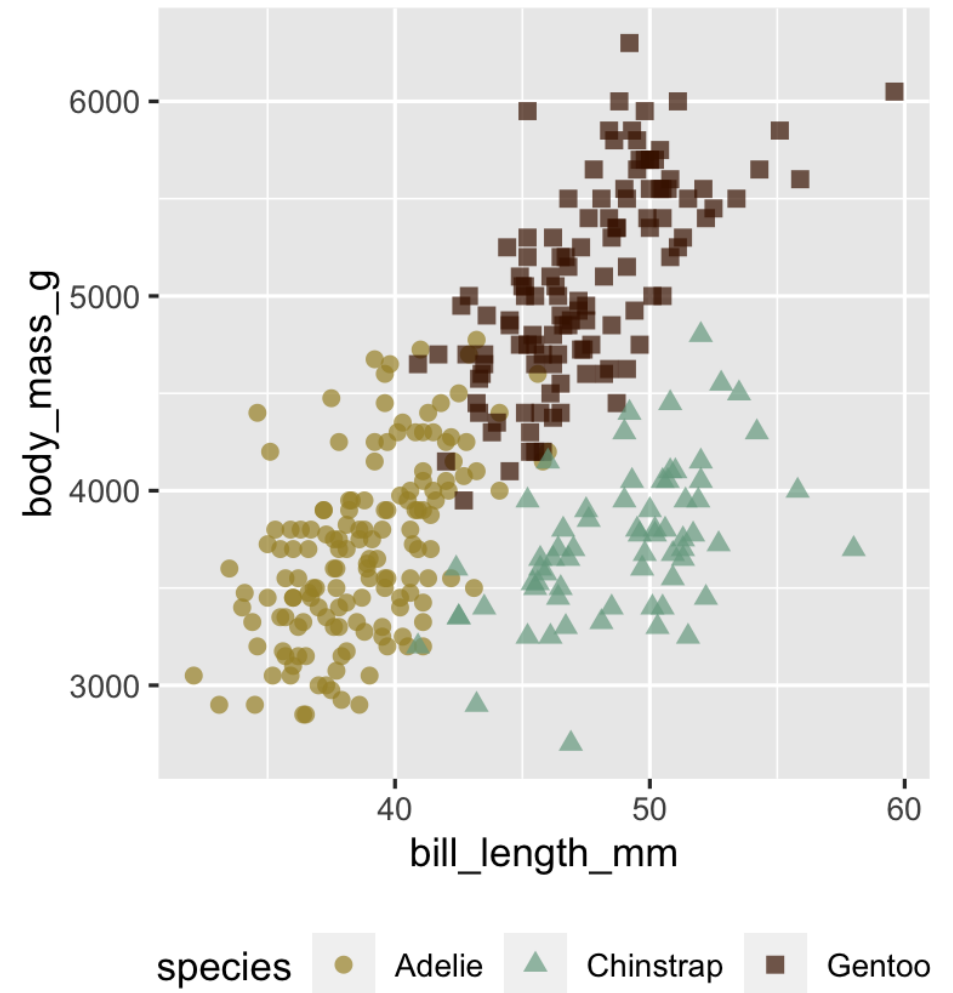
**Brown cluster
separated from
others in NW to
SE direction**



Case study 4 Penguins



Gentoo from others in contrast of fl, bd



Chinstrap from others in contrast of bl, bm

There may be multiple and different combinations of variables that reveal similar structure. 😞

The tour can help to discover these, too. 😂

Other tour types

- 🗨 **guided**: follows the optimisation path for a projection pursuit index.
- 🗨 **little**: interpolates between all variables.
- 🗨 **local**: rocks back and forth from a given projection, so shows all possible projections within a radius.
- 🗨 **dependence**: two independent 1D tours
- 🗨 **frozen**: fixes some variable coefficients, others vary freely.
- 🗨 **manual**: control coefficient of one variable, to examine the sensitivity of structure this variable. (In the **spinifex** package)
- 🗨 **slice**: use a section instead of a projection.

Guided tour

New target bases are chosen using a projection pursuit index function

maximize $g(F'x)$ subject to F being orthonormal
 F

🗨️ holes: This is an inverse Gaussian filter, which is optimised when there is not much data in the center of the projection, i.e. a "hole" or donut shape in 2D.

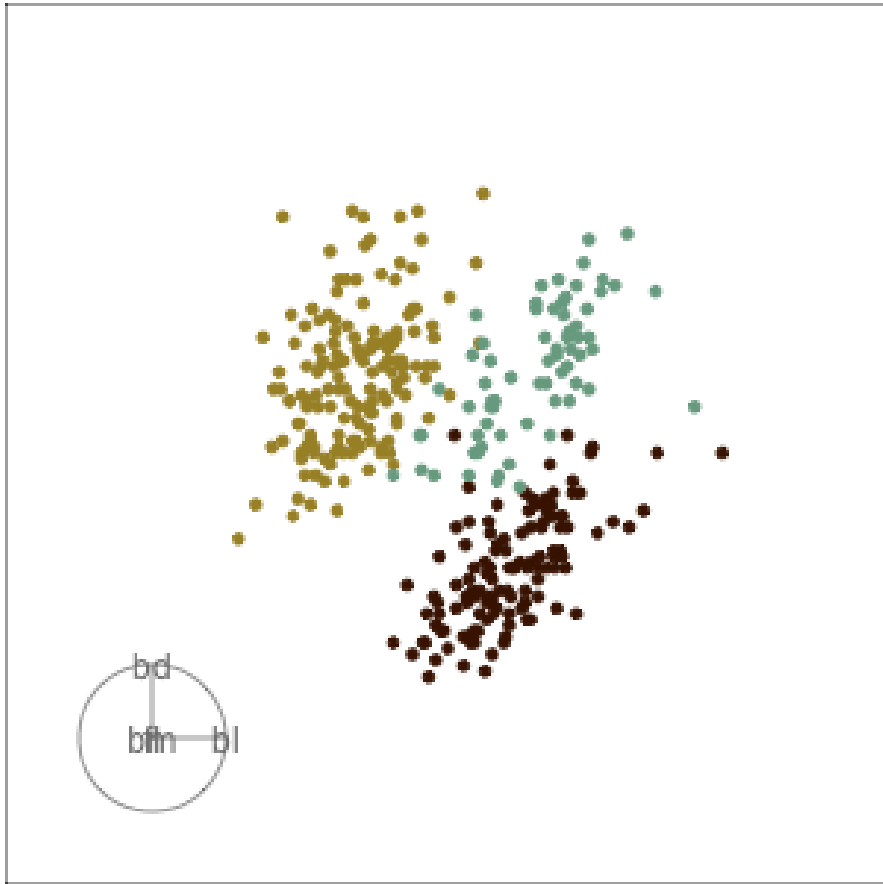
🗨️ central mass: The opposite of holes, high density in the centre of the projection, and often "outliers" on the edges.

🗨️ LDA/PDA: An index based on the linear discriminant dimension reduction (and penalised), optimised by projections where the named classes are most separated.

Case study 4 Penguins

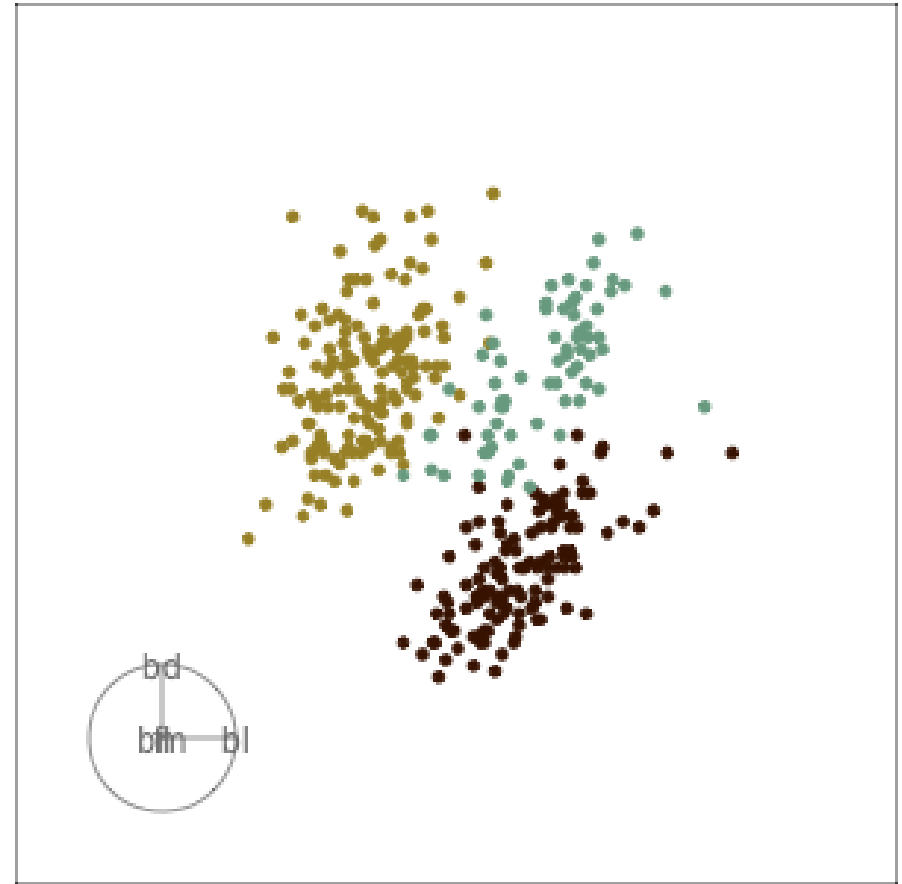


Grand



Might accidentally see best separation

Guided, using LDA index



Moves to the best separation

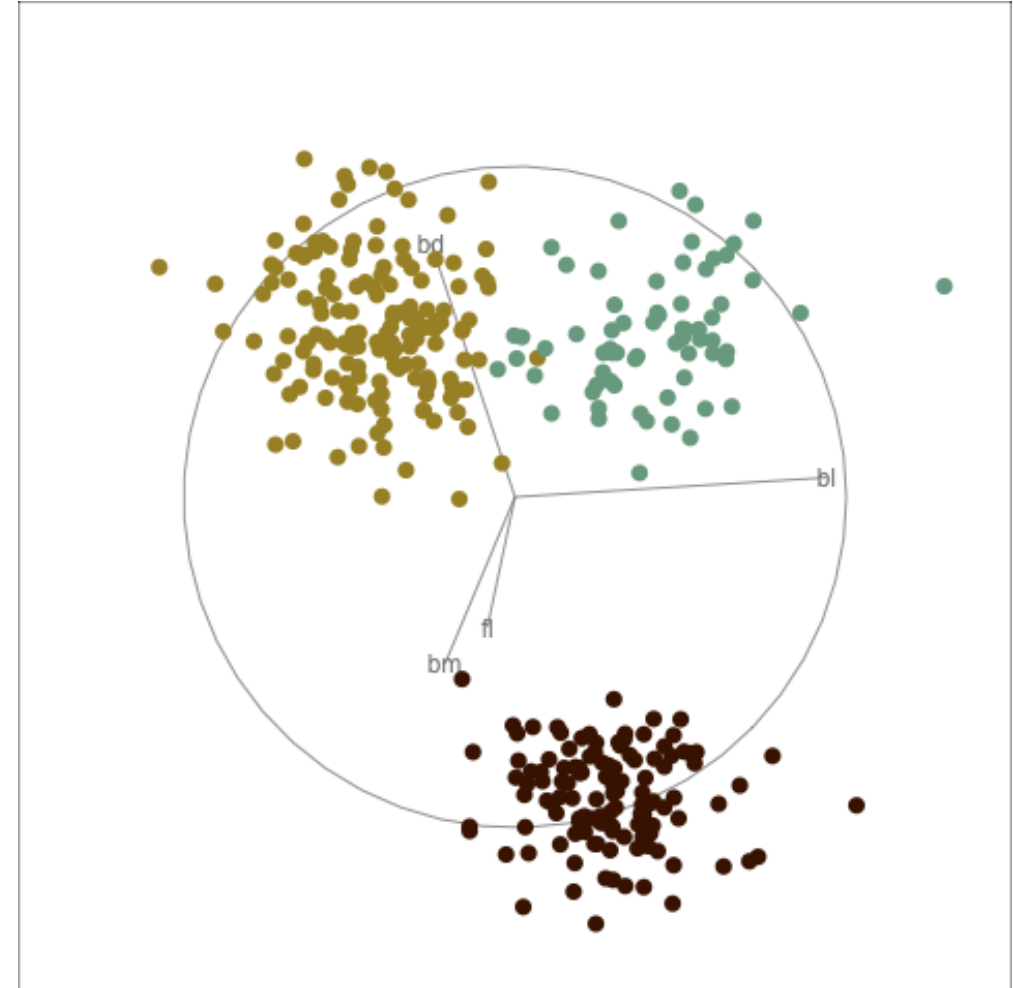
Manual tour

Control the coefficient of one variable, reduce it to zero, then increase it to 1, maintaining orthonormality

Case study 4 Penguins

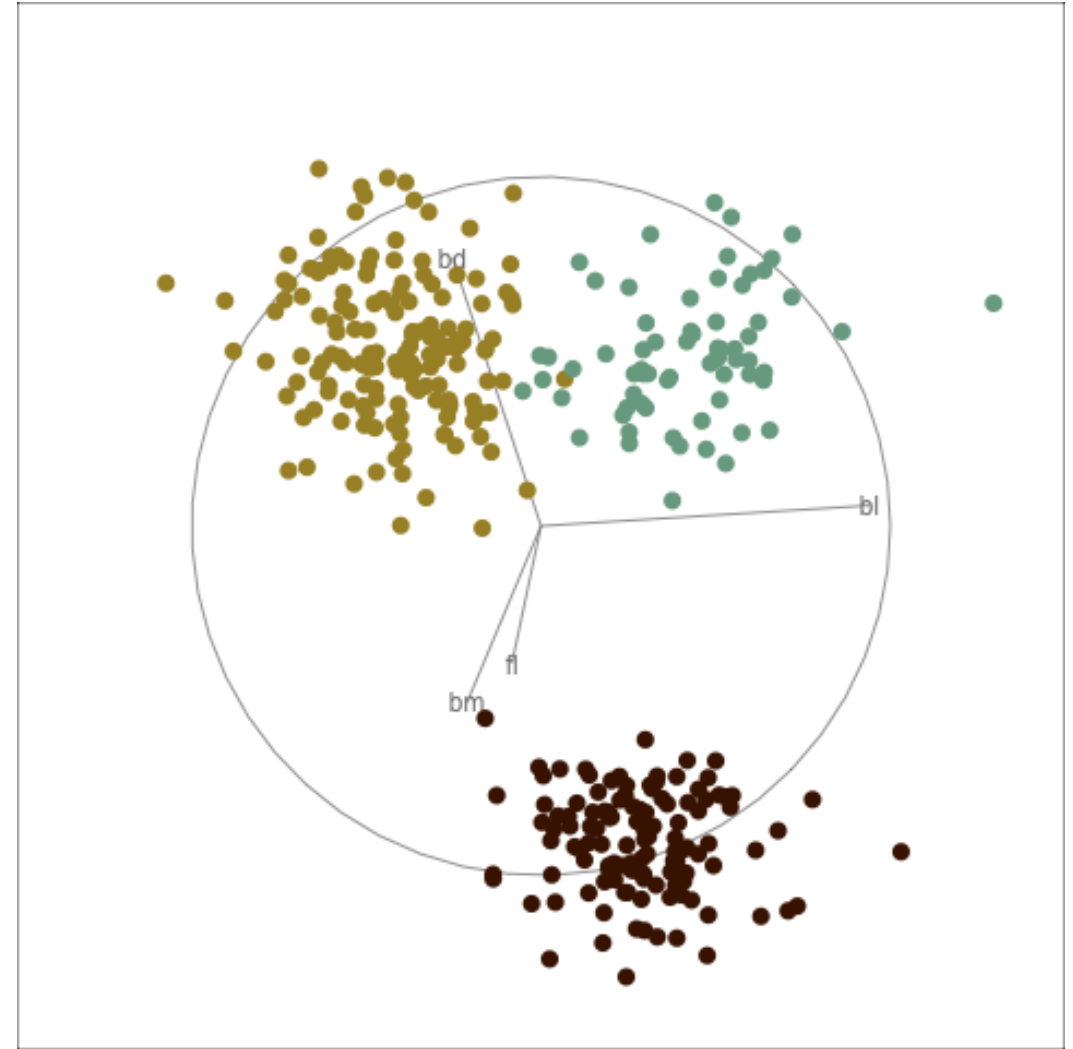


- start from best projection, given by projection pursuit
- b1 contribution controlled
- if b1 is removed from projection, Adelie and chinstrap are mixed
- b1 is important for Adelie



Case study 4 Penguins

- start from best projection, given by projection pursuit
- f1 contribution controlled
- cluster less separated when f1 is fully contributing
- f1 is important, in small amounts, for Gentoo



That's it, for this lecture!



This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/).

Lecturer: Di Cook

Department of Econometrics and Business Statistics

✉ ETC5521.Clayton-x@monash.edu

