

# Class 4: High dimensional visualisation

Andrew Parnell  
andrew.parnell@mu.ie



PRESS RECORD

[https://andrewcparnell.github.io/dataviz\\_course](https://andrewcparnell.github.io/dataviz_course)

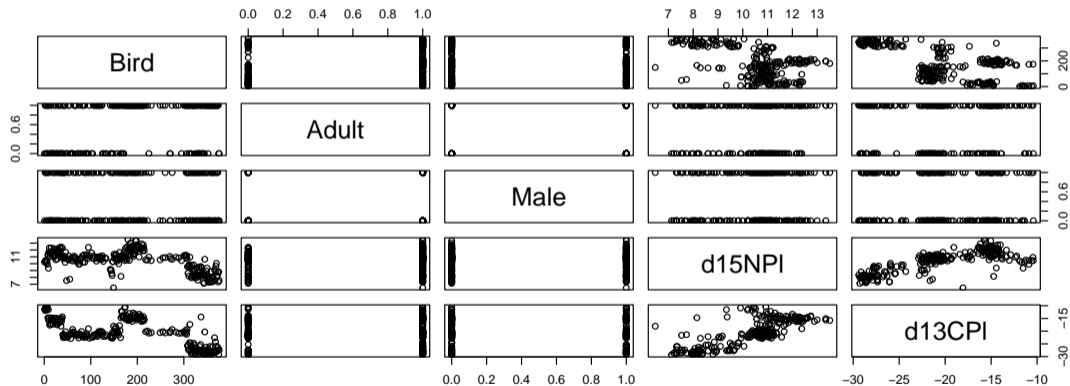
## Learning outcomes

- ▶ Learn how to create multiple scatter plots
- ▶ Plot some map data
- ▶ Plot some time series data
- ▶ Use some of the 3D geoms
- ▶ Our last session where we mostly use ggplot



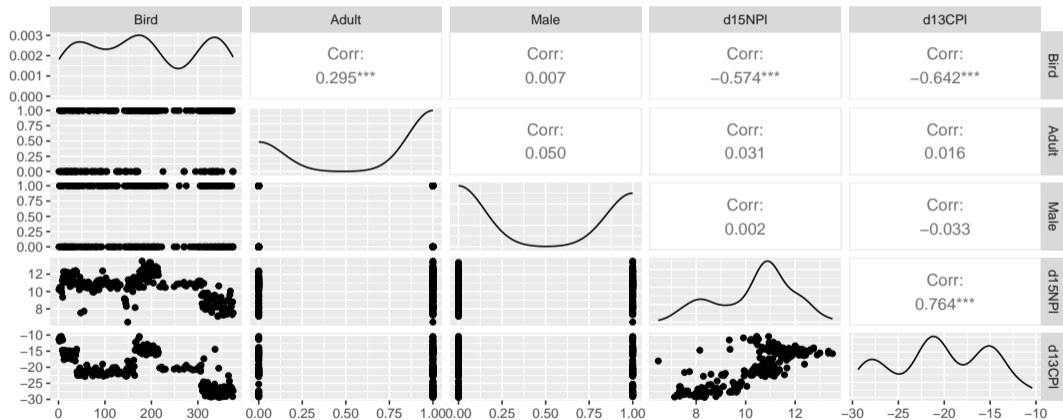
# Multiple scatter plots: pairs

```
pairs(geese[,1:5])
```



## ggplot version using package GGally

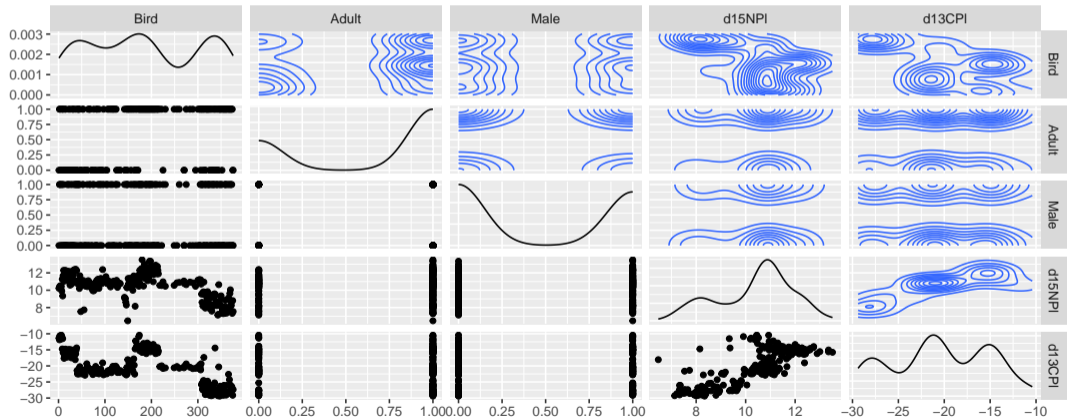
```
library(GGally)
ggpairs(geese[,1:5])
```



Very customisable too - can change what the upper, lower and diagonal panels contain or add colours to different groups

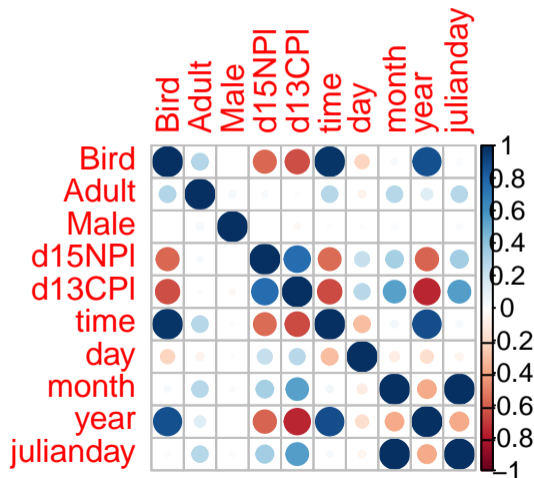
# Extended ggpairs plot

```
ggpairs(geese[,1:5],  
        upper = list(continuous = "density"),  
        lower = list(continuous = "points"))
```



## Visualise correlations

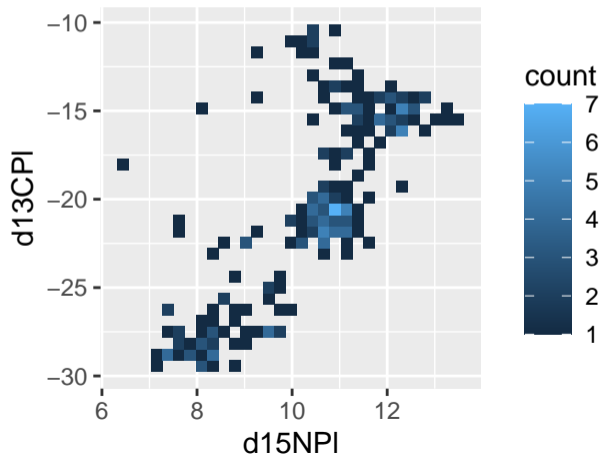
```
library(corrplot)
corrplot(cor(geese)) # Note: need a correlation matrix!
```



## 2D density plots

ggplot2 has specific functions for plotting 3D density data

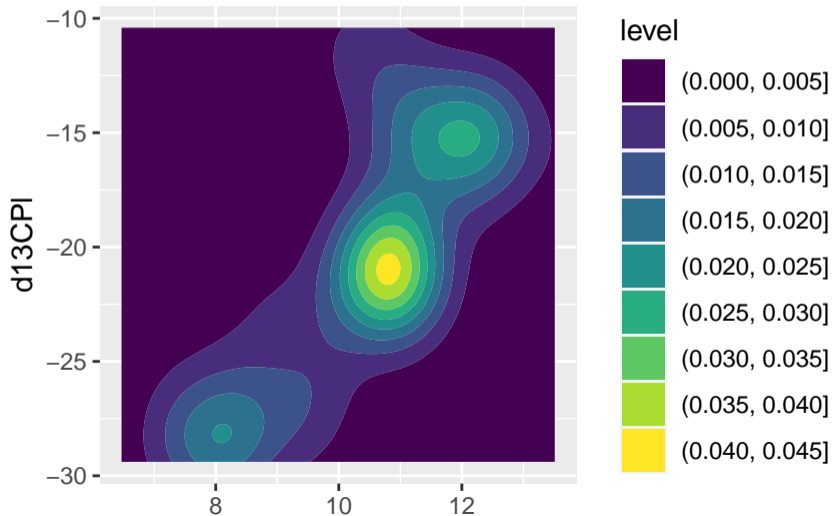
```
ggplot(geese, aes(x = d15NP1, y = d13CP1)) +  
  geom_bin_2d() + theme(aspect.ratio = 1)
```





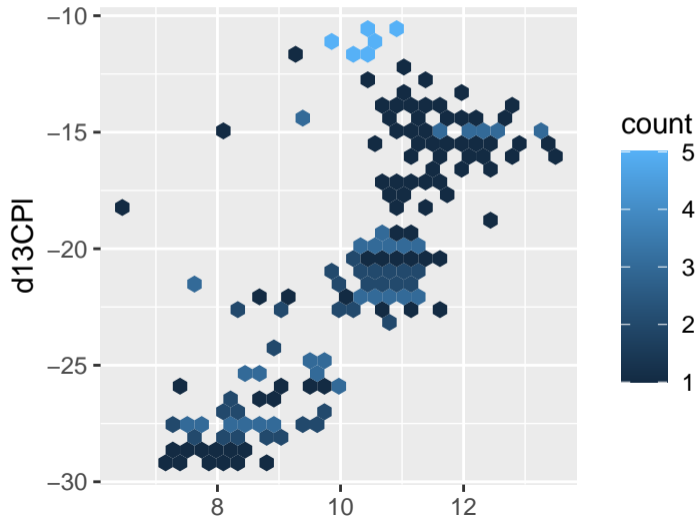
## Contour plots

```
ggplot(geese, aes(x = d15NP1, y = d13CPI)) +  
  geom_density_2d_filled() + theme(aspect.ratio = 1)
```



## Hexagonal plots

```
ggplot(geese, aes(x = d15NP1, y = d13CP1)) +  
  geom_hex() + theme(aspect.ratio = 1)
```

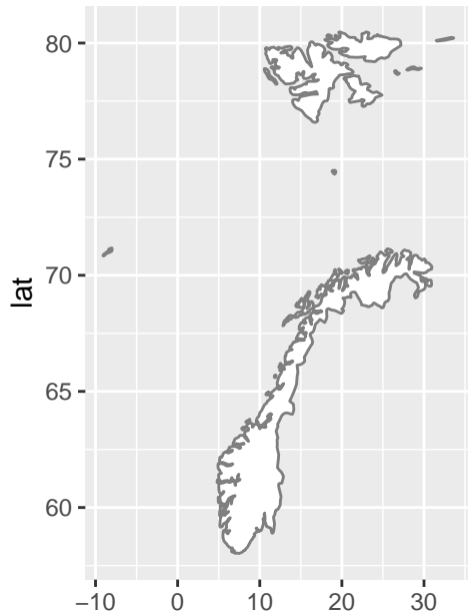


## Plotting data with maps

- ▶ Most base R plots of maps are ugly. `ggplot` has an interface to the `maps` package but these are quite ugly and out of date
- ▶ Better package are `ggmap`, `ggspatial` and `sf` (simple features) which have loads of examples and interfaces for both base R and `ggplots`
- ▶ Three common use cases:
  - ▶ You want to take a map of an area and add information to it by filling in the areas or placing points on the map
  - ▶ You want to take an e.g. Google or OpenStreetMap (raster) map and add information on top
  - ▶ You have your own shape file and want to create a plot on top of that (not covered – see `read_sf`)



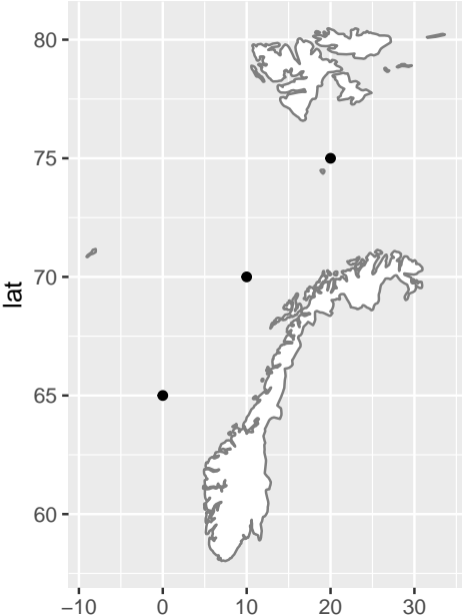
## Scenario 1: plot



## Adding information

```
points <- tibble(  
  long = c(0, 10, 20),  
  lat = c(65, 70, 75),  
  group = 1 # Necessary to match the fitting command  
)  
ggplot(norway, aes(long, lat, group = group)) +  
  geom_polygon(fill = "white", colour = "grey50") +  
  coord_quickmap() +  
  geom_point(data = points) # Add colours etc here
```

# Plot with added information

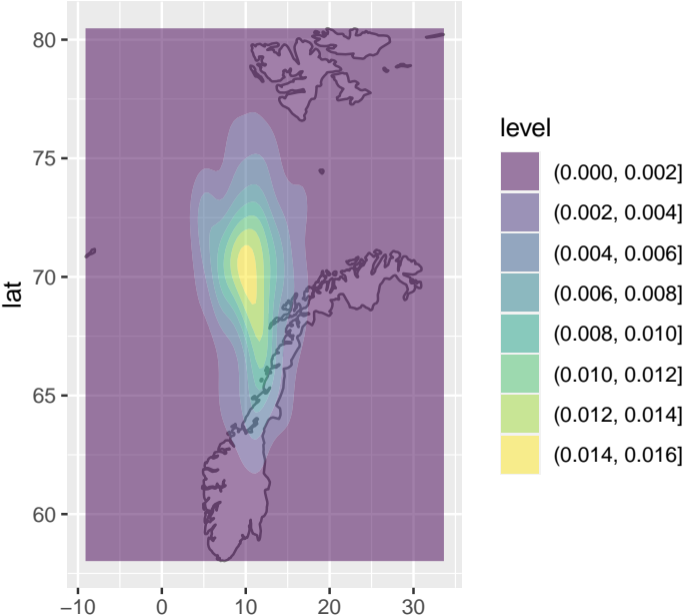


## Scenario 1 (cont): with e.g. contours

```
new_data <- MASS::mvrnorm(n = 100,  
                          mu = c(10, 70),  
                          Sigma = diag(10, 2)) %>%  
  as_tibble() %>%  
  rename(long = V1, lat = V2) %>%  
  mutate(group = 1)  
ggplot(norway, aes(long, lat, group = group)) +  
  geom_polygon(fill = "white", colour = "grey50") +  
  coord_quickmap() +  
  geom_density_2d_filled(data = new_data, alpha = 0.5)
```



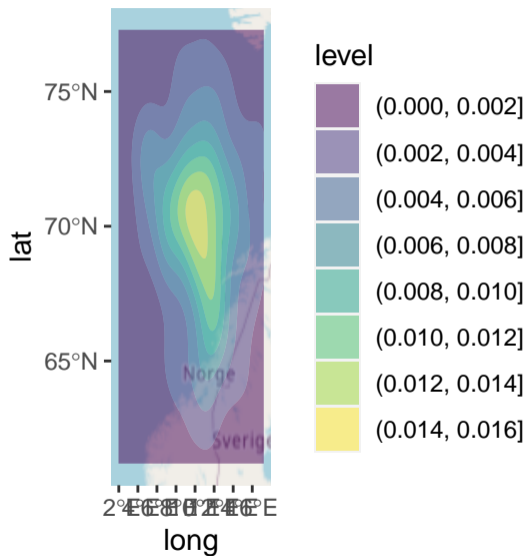
# Scenario 1: plot



## Scenario 2: overlaying an existing web map

```
library(ggspatial)
ggplot(new_data, aes(x = long, y = lat)) +
  annotation_map_tile(type = "osm") +
  geom_density_2d_filled(alpha = 0.5) +
  coord_sf(crs = 4326)
```

## Scenario 2: plot

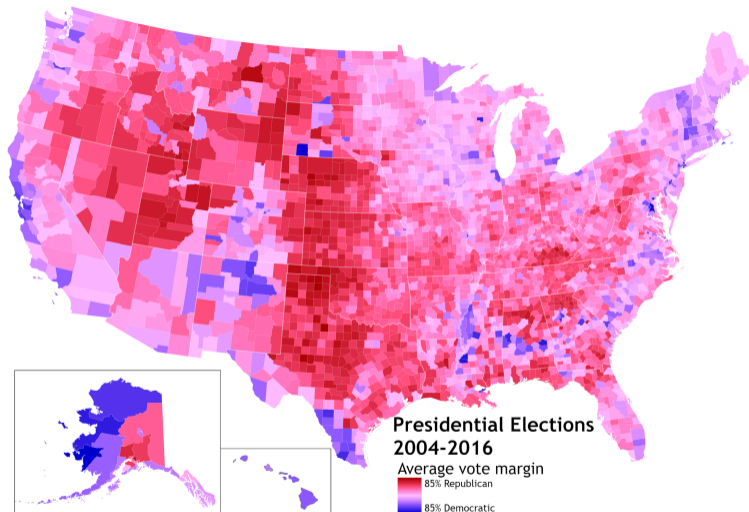


# Leaflet

See demo in Rstudio:

```
library(leaflet)
leaflet() %>%
  addTiles() %>%
  setView(lng = 10, lat = 62, zoom = 5) %>%
  addMarkers(lng = 10, lat = 62) %>%
  addCircles(lng = 11, lat = 63, radius = 5000)
```

## Common problems with choropleth maps

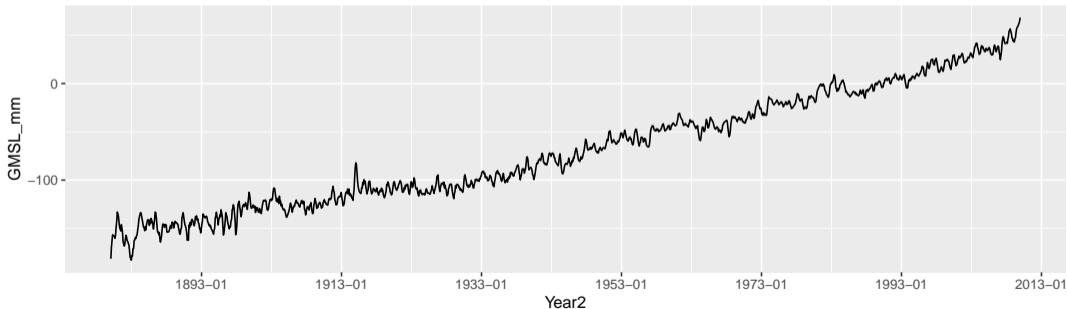


From Wikipedia



## Plots of time series

```
sl_month %>%  
  mutate(Year2 = lubridate::date_decimal(Year)) %>%  
  ggplot(aes(x = Year2, y = GMSL_mm)) +  
  geom_line() +  
  scale_x_datetime(date_breaks = "20 years",  
                  date_labels = "%Y-%m")
```



## dygraphs - interactive time series

See Rstudio demo

```
library(dygraphs)
sl_month %>%
  select(Year, GMSL_mm) %>%
  dygraph() %>%
  dyRangeSelector()
```

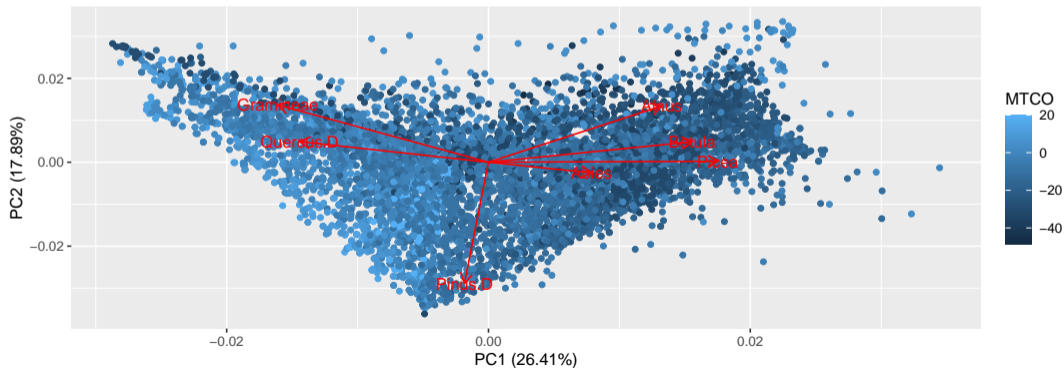


## Automatic plotting with `ggfortify`

- ▶ There is a function in `ggplot2` called `autoplot` that tries to take into account the object class when creating a plot
- ▶ Other packages can use and extend this function, for example the `forecast` package uses `autoplot` to create nice time series plots
- ▶ The `ggfortify` package uses this even more generally to create nice plots for a load of different classes of object

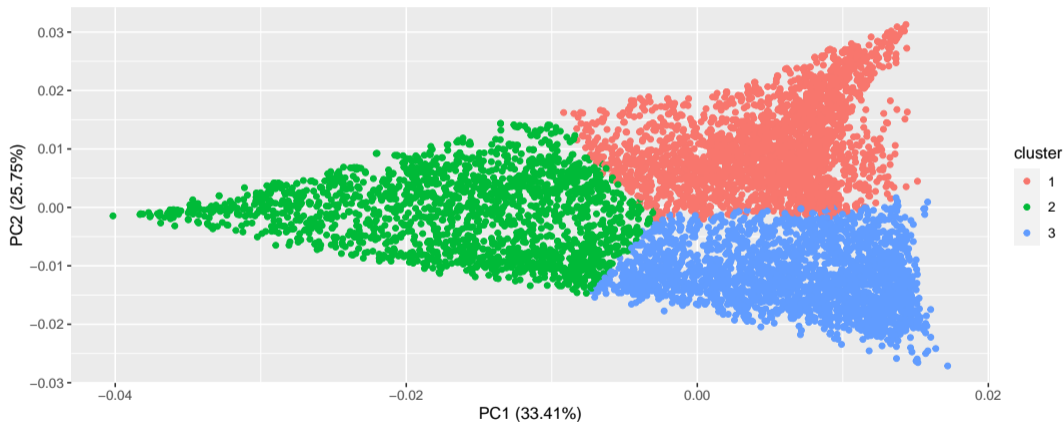
# Principal Components Analysis and ggfortify

```
library(ggfortify)
pollen <- read.csv('../data/pollen.csv')
pca <- prcomp(pollen[,3:9], scale. = TRUE)
autoplot(pca, data = pollen, colour = "MTCO",
         loadings = TRUE, loadings.label = TRUE)
```



## Cluster analysis and ggfortify

```
autoplot(kmeans(pollen[,3:9], 3),  
         data = pollen)
```



- ▶ Can also do generalised linear models, maps, time series, other dimension reduction methods (factor analysis and multidimensional scaling), and survival data

## Exercise

Run some of the above dimension reduction techniques on the geese data above and produce some plots. Try removing different columns (e.g. the time columns. What does the PCA tell you? (Post plots and answers in Slack)

## Summary

- ▶ Lots of options for high dimensional data. Start with matrix scatter plots
- ▶ Plotting maps is a minefield. Pick a package and work through it for what you want to produce. There are many other packages and complications
- ▶ Always a good idea to try some dimension reduction when plotting high dimensional data; `ggfortify` a simple first choice