



Faculty of Science



BASICs in R

Introduction

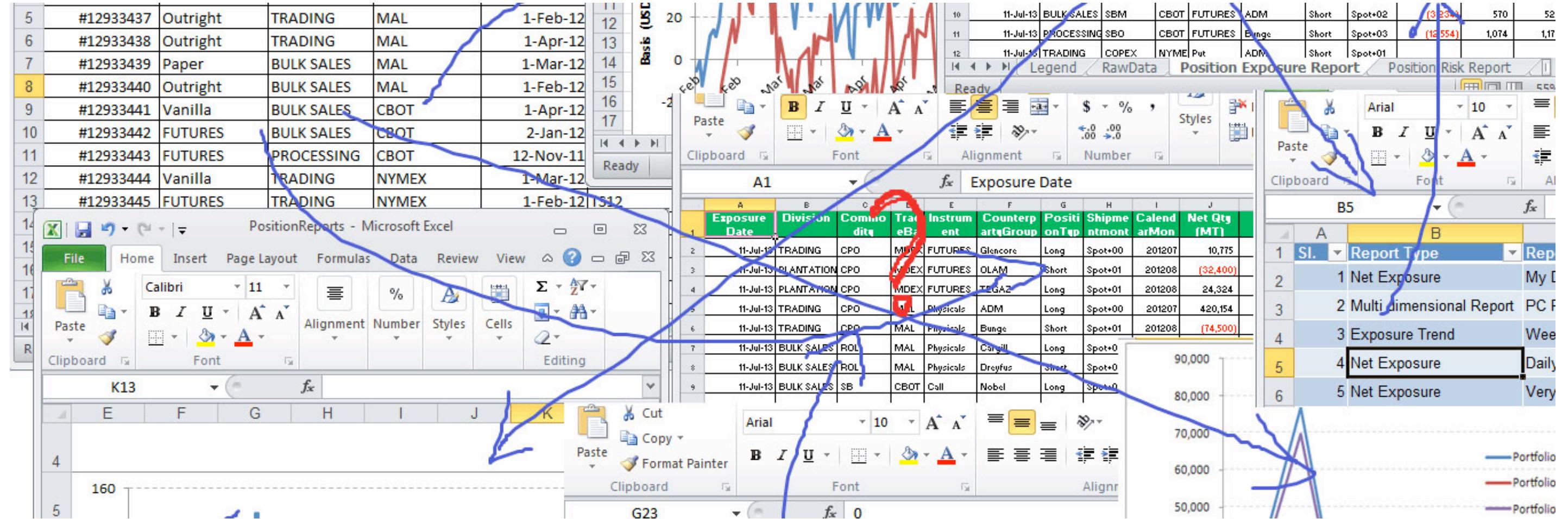
Tuesday 22nd January 2019

Anne Bech Risum
Peter Skou
Morten Arendt Rasmussen

MOTIVATION

- Digitalization on KU
- Reproducibility
- Get quicker insight and More knowledge out of your data





Dias 3



The image shows a screenshot of the R Studio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Window, Help, and a series of system icons. The title bar says "RStudio". The main area has two tabs open: "Untitled14*" and "Untitled15*". The "Untitled15*" tab shows the R console output:

```
R version 3.5.1 (2018-07-02) -- "Feather Spray"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.
```

Below the console is a "Do it!" button in large red text. To the right of the console is a "Where it is executed" section, which displays the RStudio environment pane showing an empty global environment and a file browser pane listing various files and folders in the current directory.

The RStudio interface includes a top bar with system icons like battery level, signal strength, and volume, and a status bar at the bottom showing the date and time (Mon 10.38) and the user's name (Morten Arendt Rasmussen).

What you want to do

The output of your effort

Plots, Files, Packages,...

Packages

R comes with some functionality, but not everything is covered.

Therefor we need additional functions.

There comes in the form of ***packages*** which is installed directly from within R.

```
> install.packages('ggplot2')
```

From
CRAN (<http://cran.r-project.org/>)
(13713 available packages)

```
> devtools::install_github('vqv/ggbiplot')
```

From
github (<https://github.com/>)



Data import

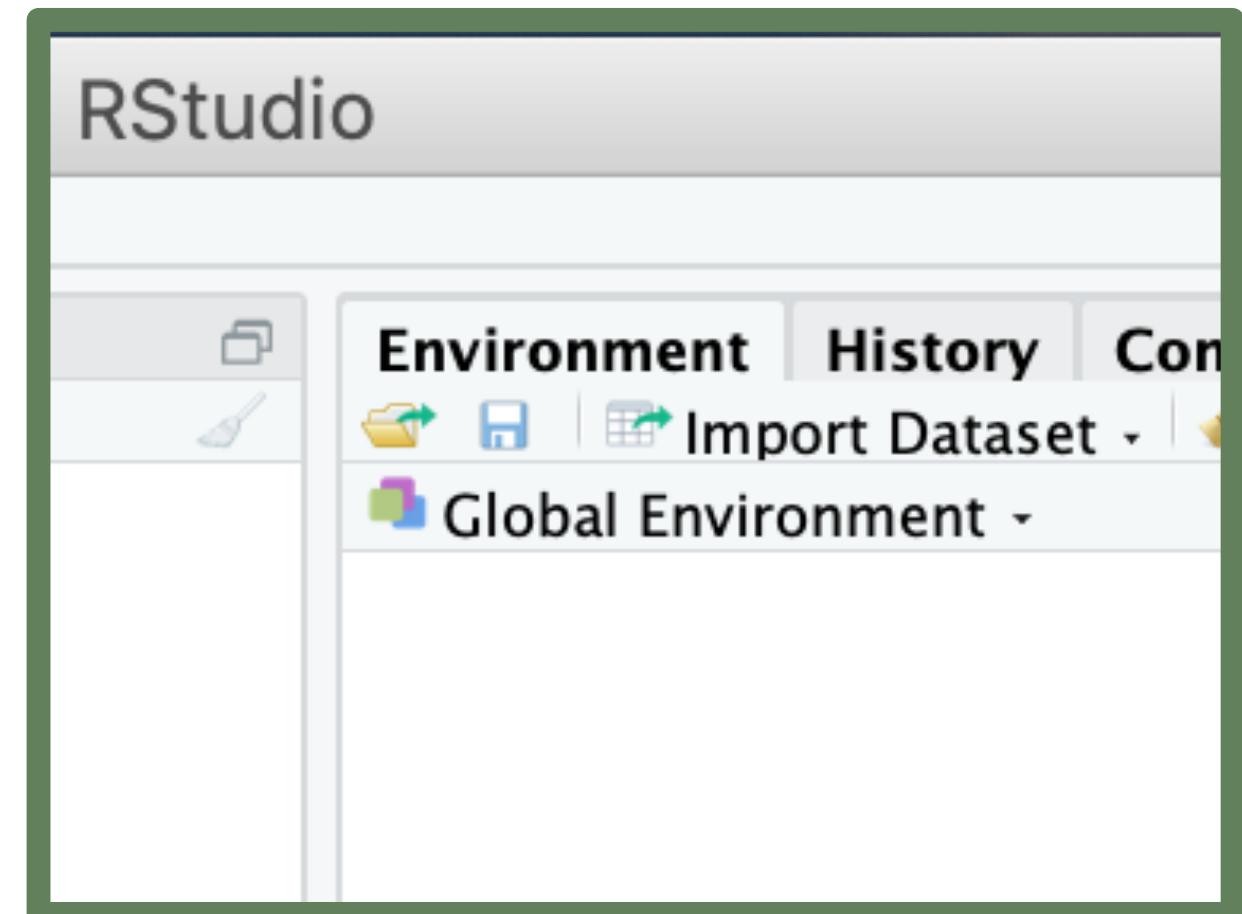
Use the *point-n-click* method in Rstudio

.... Eventually copy paste the commands produced into the script

Or

Use the **rio** package

```
> library(rio)  
> X <-  
import('myExcelFile.xlsx')
```



data.frame()

You should store your data in a data.frame

A data frame is as an excel sheet with first row being variable names.

To be aware of:

Avoid using space, leading numerics and repetitions in names.

Some useful functions to look at the data.frame

> head() / tail()	> View()	> rownames()
> str()	> colnames()	



Descriptive stats

Calculate descriptive stats – a *single number describing a distribution of numbers* – splitted according to grouping information.

Functions

- > mean(), sd(), min(), max(), range(), length()
- > aggregate()

Overlook of a design

How many observations are there for a combination of (design) variables

- > table()



Plot everything with ggplot2



PLOTS

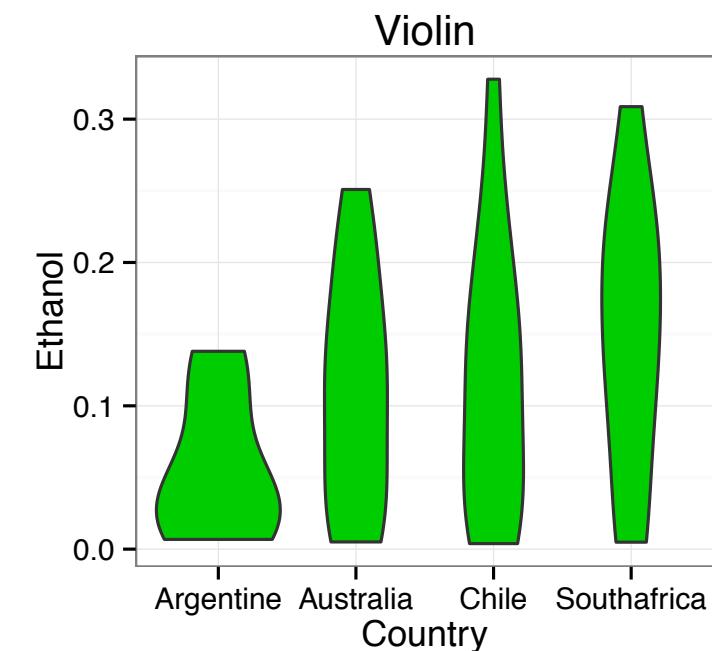
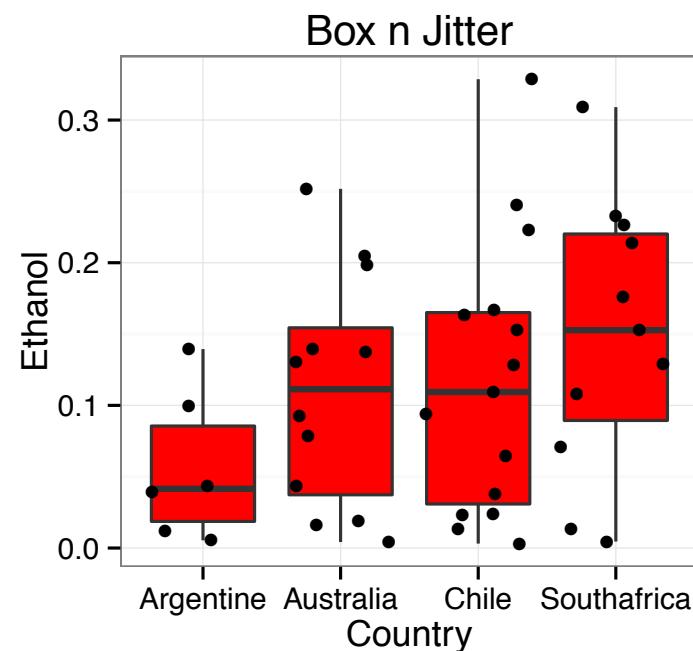
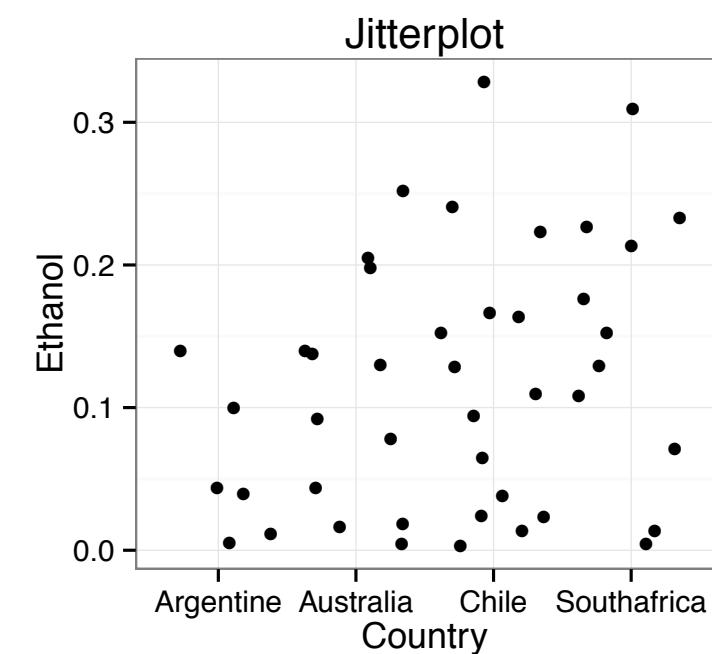
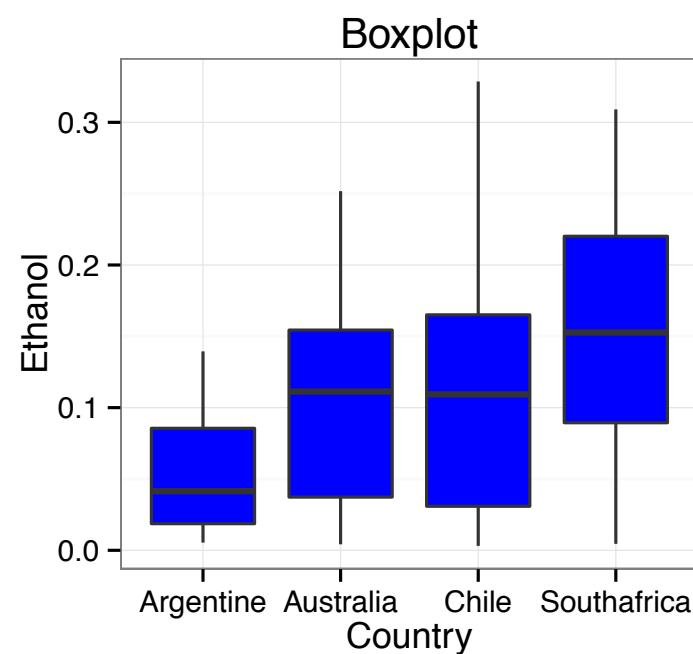
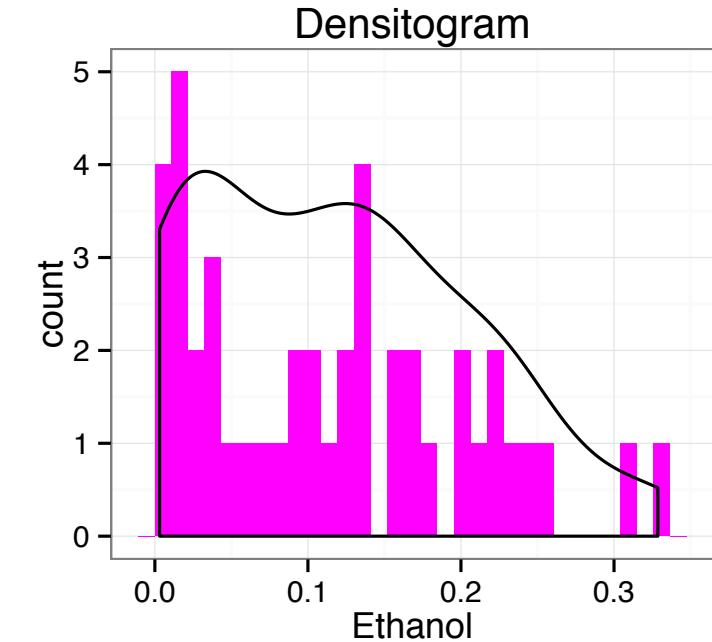
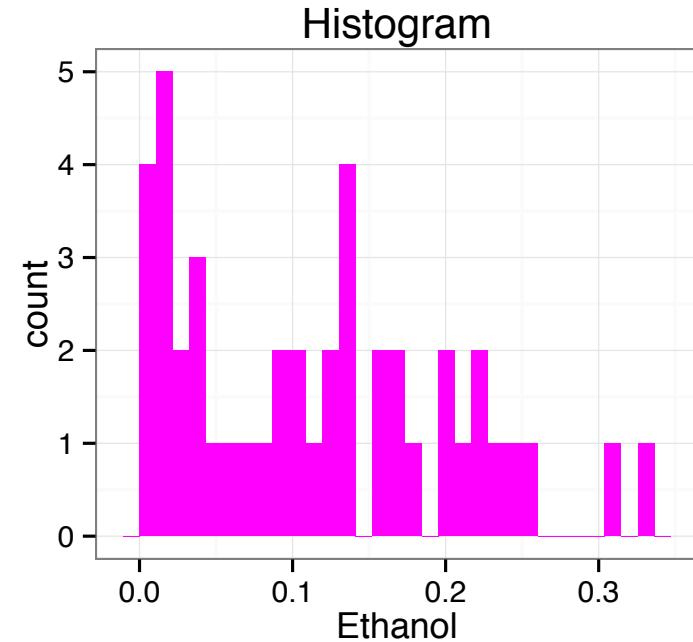
- Histogram, Densitogram, Boxplot, Jitterplot, Violinplot
 - Continuous data
 - Describes distribution
- Lineplot, Stem plot
 - Continuous data
 - Ordinal in e.g. time
- Bar chart, pie chart
 - Categorical data or compositional data
- Scatterplot
 - Bi variate continuous data
- Spline plot
 - Smooth relation between two (or more) variables



Histogram, Densitogram, Boxplot, Jitterplot, Violin plot,

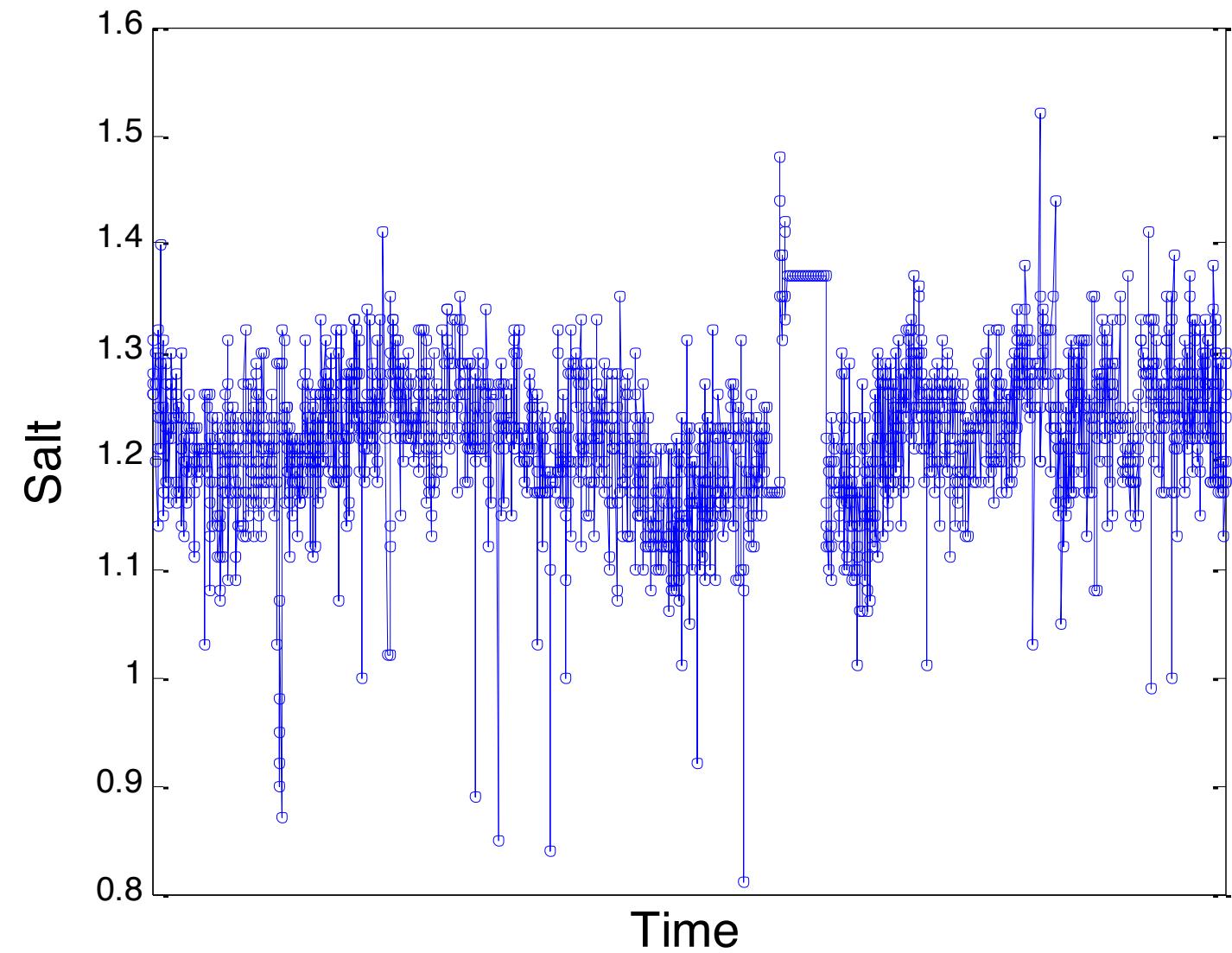
- Jitter shows raw data
 - Boxplot and violin can handle many observations
- Densitogram and Violin may cheat in the representation of raw data

Continuous data
Describes distribution



Lineplot

- Is used in e.g. Statistical Process Control (SPC)
- Excellent for capturing drift (systematical change over time) in the laboratory or in the production.



Continuous data
Ordinal in e.g. time

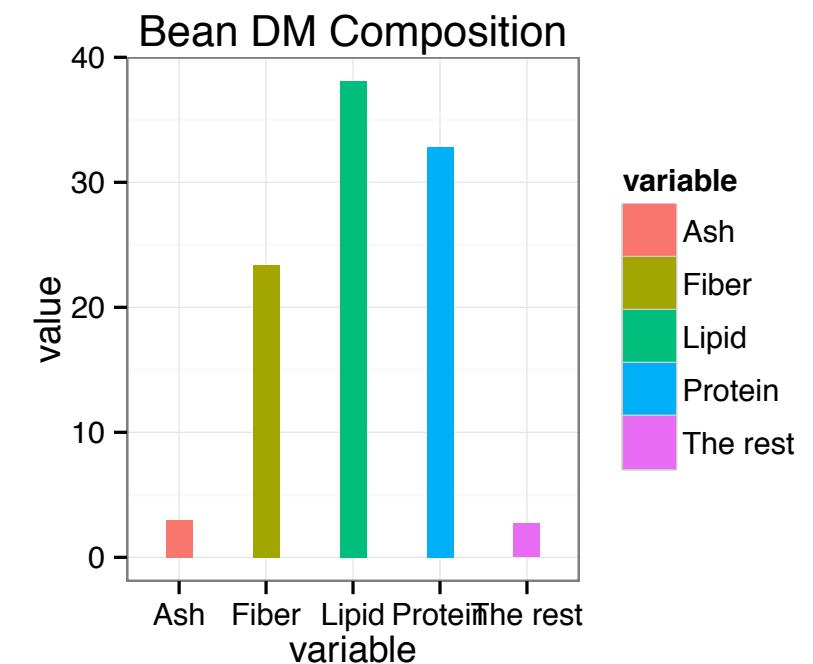
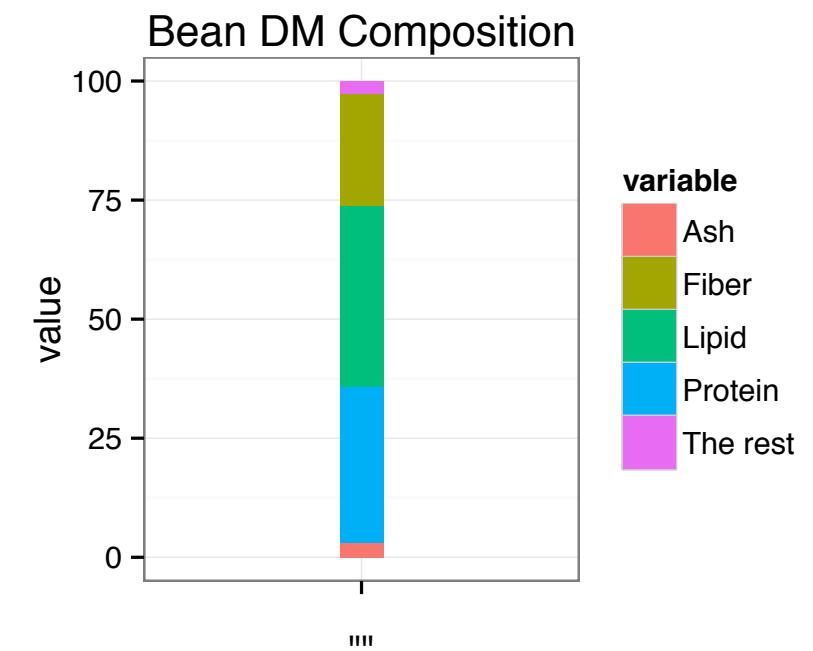
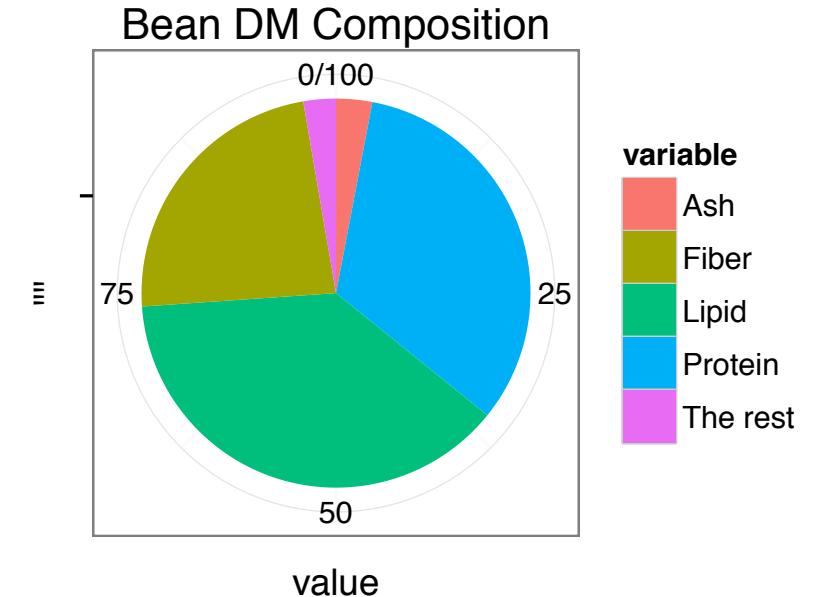


Bar charts

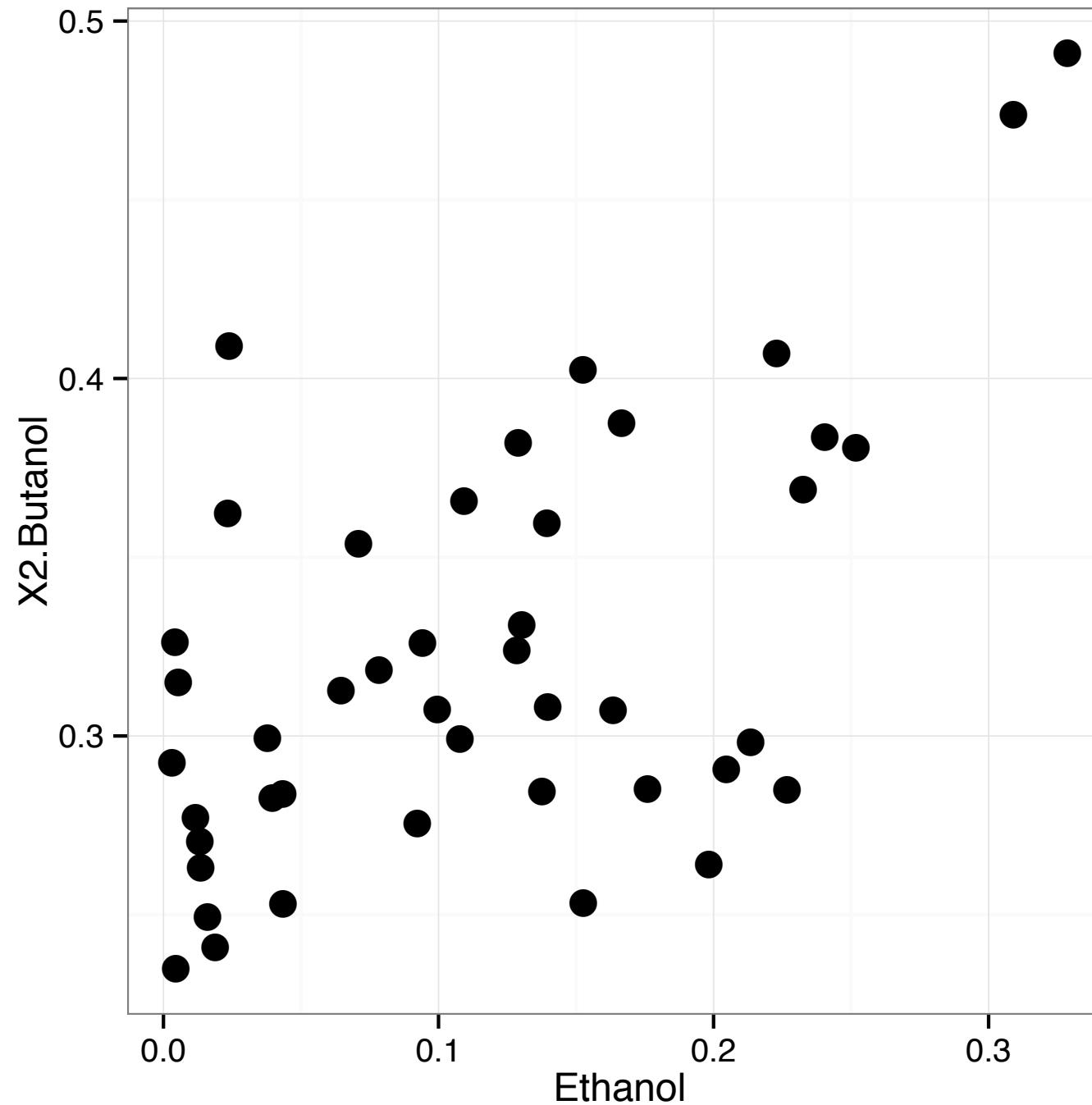
Pie charts

- OBS: Do not replace a boxplot with a bar chart!

Compositional data
Categorical data



Scatterplot

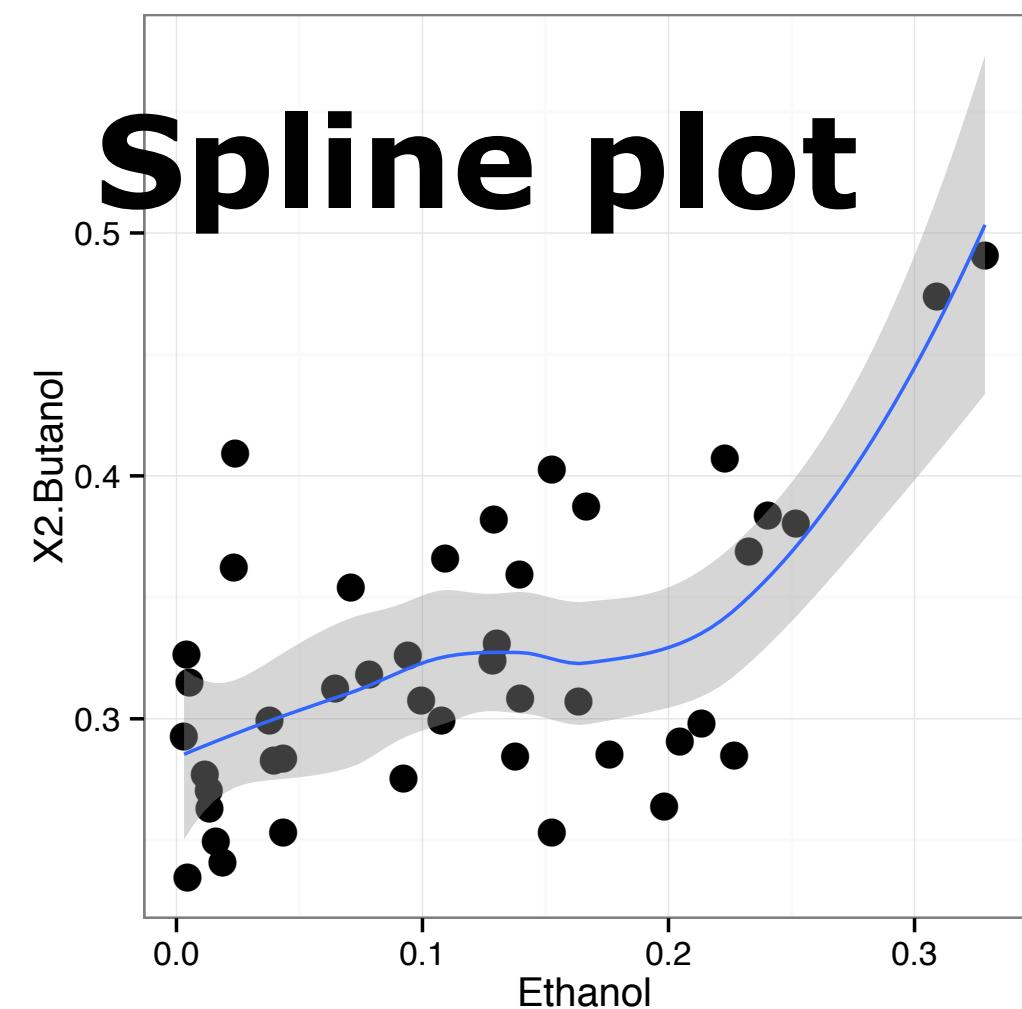
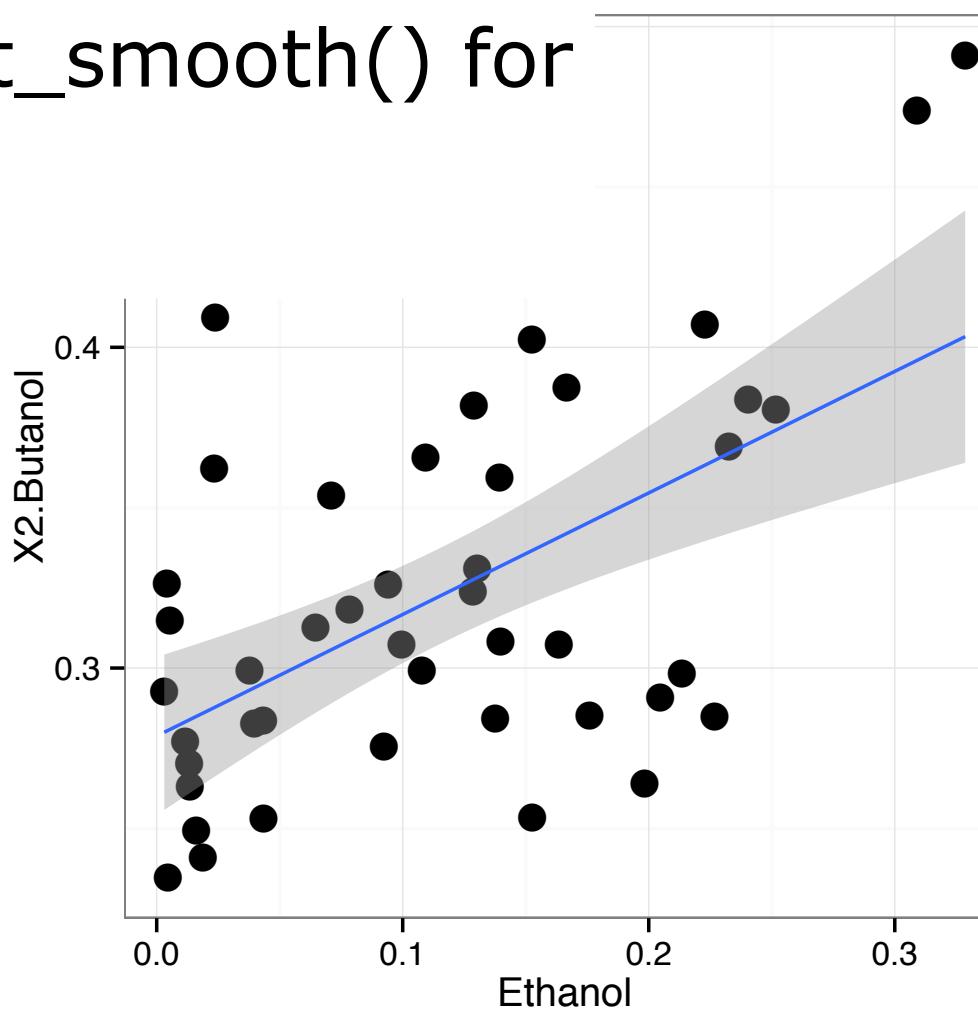
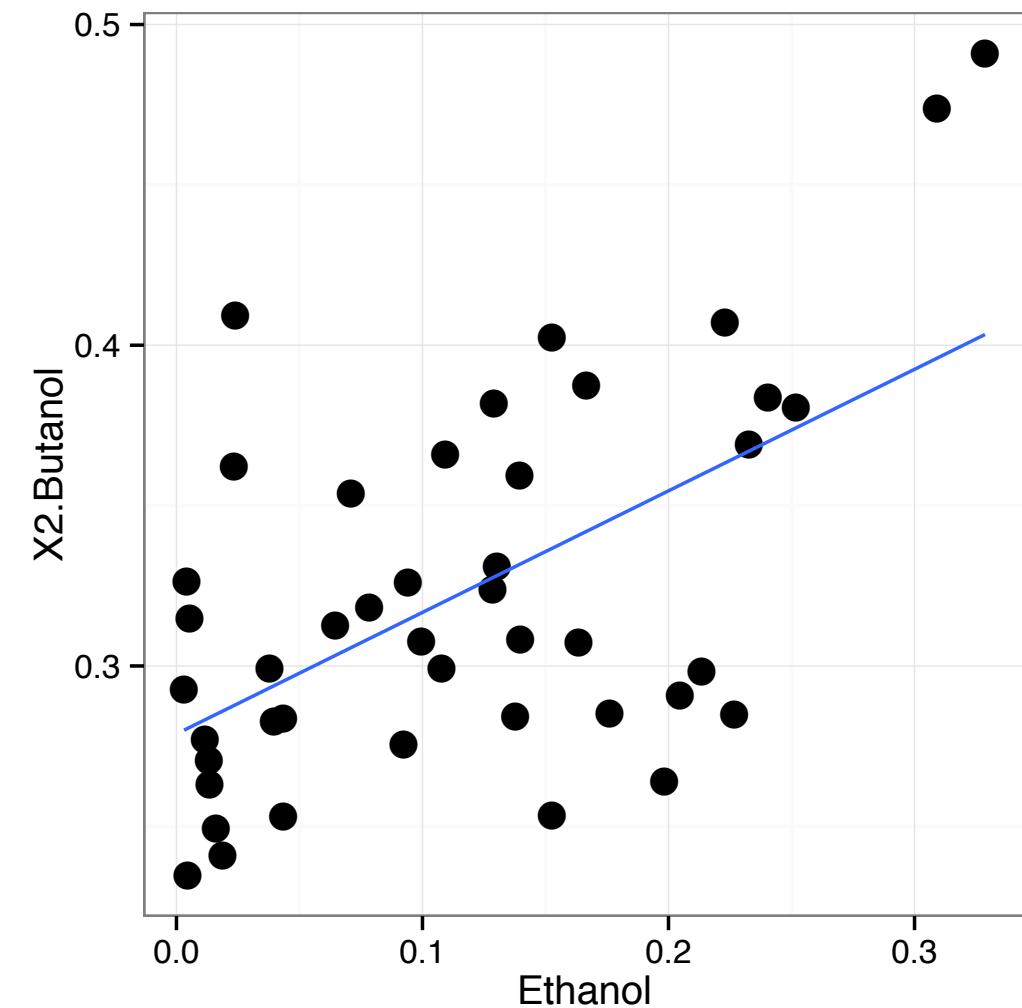
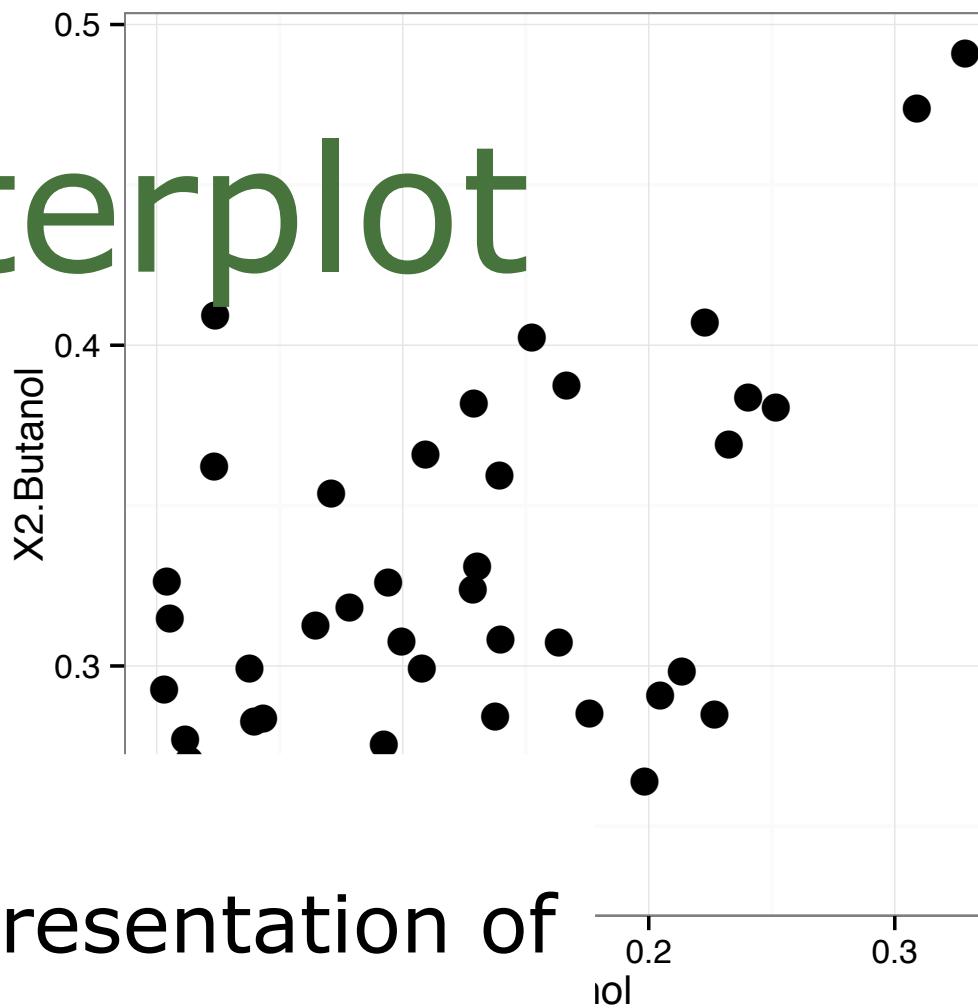


- Two response variables versus each other
- Very usefull for Multivariate data



Scatterplot

- Same data!
- Different representation of the relation
- Check: `+stat_smooth()` for details



t.test Linear models and ANOVA

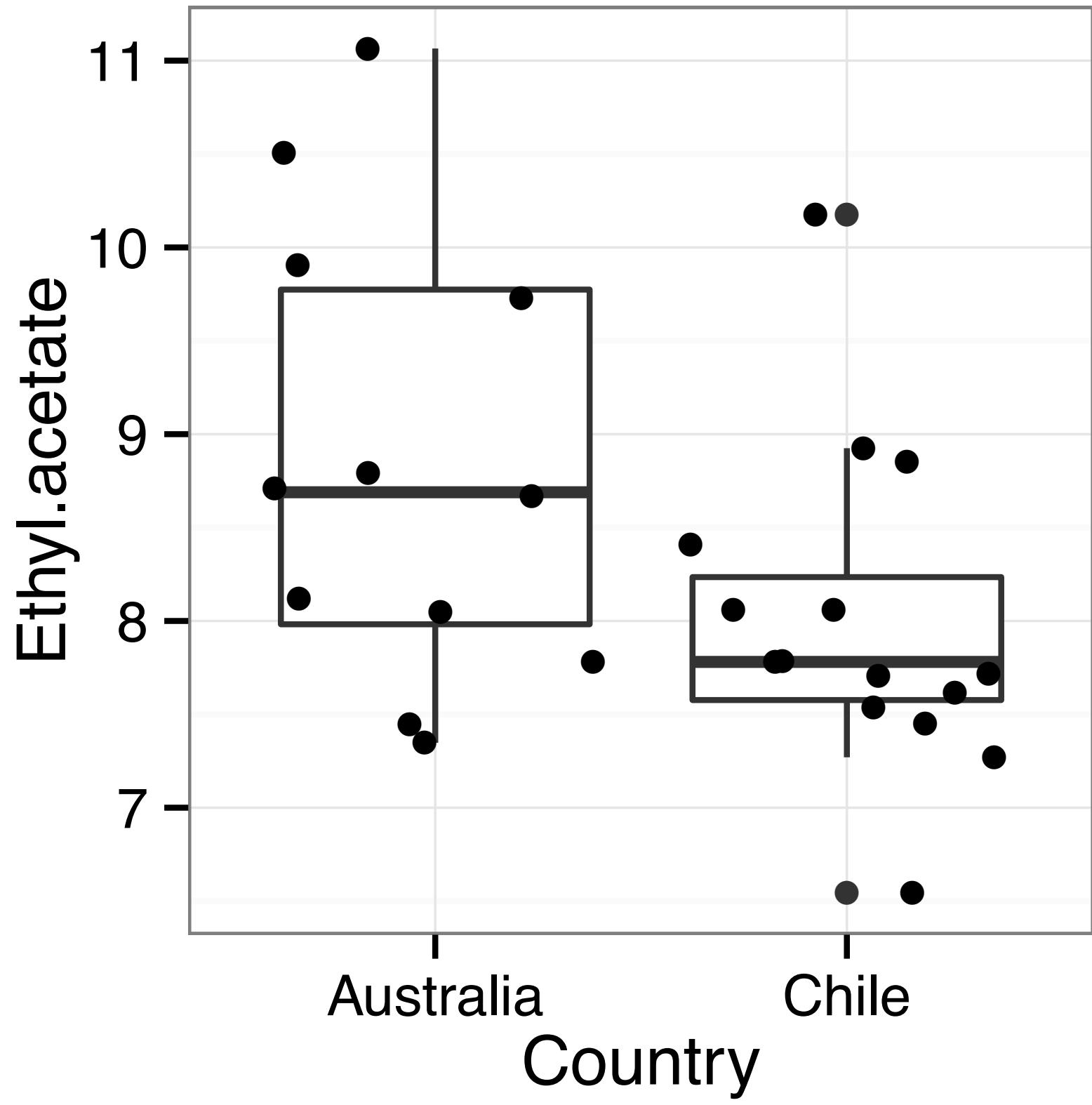


T-test

Above ALL
statistical software should
be able to do a **t.test**

... to compare the means of
two distributions

Function
`> t.test()`



“Account” of the Variation ANOVA table

Source	$SumSq$	Df	$MeanSq$	F_{obs}	$Pr(F_{df_{Source}, df_e} > F_{obs})$
Factor(A)	SS_A	$df_A = k_A - 1$	$MS_A = SS_A / df_A$	$F_A = MS_A / MS_e$	p_A
Factor(B)	SS_B	$df_B = k_B - 1$	$MS_B = SS_B / df_B$	$F_B = MS_B / MS_e$	p_B
Residuals	SS_e	$df_e = n - df_A - df_B - 1$	$MS_e = SS_e / df_e$		
Total	SS_{tot}	$df_{tot} = n - 1$			

k_A and k_B describes the number of levels within factor A and B respectively



Model assumptions

```
in R  
m <- aov()  
plot(m)
```

$$X_i = \alpha(A_i) + \beta(B_i) + \gamma(A_i, B_i) + e_i$$

where $e_i \sim \mathcal{N}(0, \sigma^2)$ and independent

for $j = 1, \dots, n$

We assume that the residuals:

- 1) Comes from the **same** Normal distribution
- 2) That they are **independent**

This can be checked via inspection of the residuals.

Fix for **same** distribution: transformation of the response variable

Fix for **Independence**: Subject for another day...



ANOVA- *Cook book*

1. **Visualize data** (w. two factors: connect the points).
2. Consider data **transformation**
3. Specify a model and some relevant hypotheses.
4. Calculate ANOVA and maybe adjust the model (remove terms)
5. Check the model assumptions (`plot(aov())`) – maybe repeat from pt 2
6. Compare contrasts via **confidence intervals** for differences and maybe a **test** for whether the difference could be =0
7. **Comment on the results**



Useful aov() and lm() functions

```
model<-aov()  
# calculates model  
  
anova(model)  
= summary(model)  
drop1(model)  
# tests each factor versus  
the full model  
  
model.tables(model)  
# returns estimates or  
means  
  
coef(model)  
# display the model  
parameter estimates  
  
confint(model)  
# calculates confidence  
intervals for the model  
parameters  
  
plot(m)  
# makes plots for validating  
the model assumptions
```



PCA



Correlated Multivariate Data

GC-Aroma profiles, Sensorical attribute data, Spectral data,... are all examples of **Multivariate data**.

PCA might be the most powerful tool to initially look at such data, to gain insight on outliers, grouping, structure between variables etc.

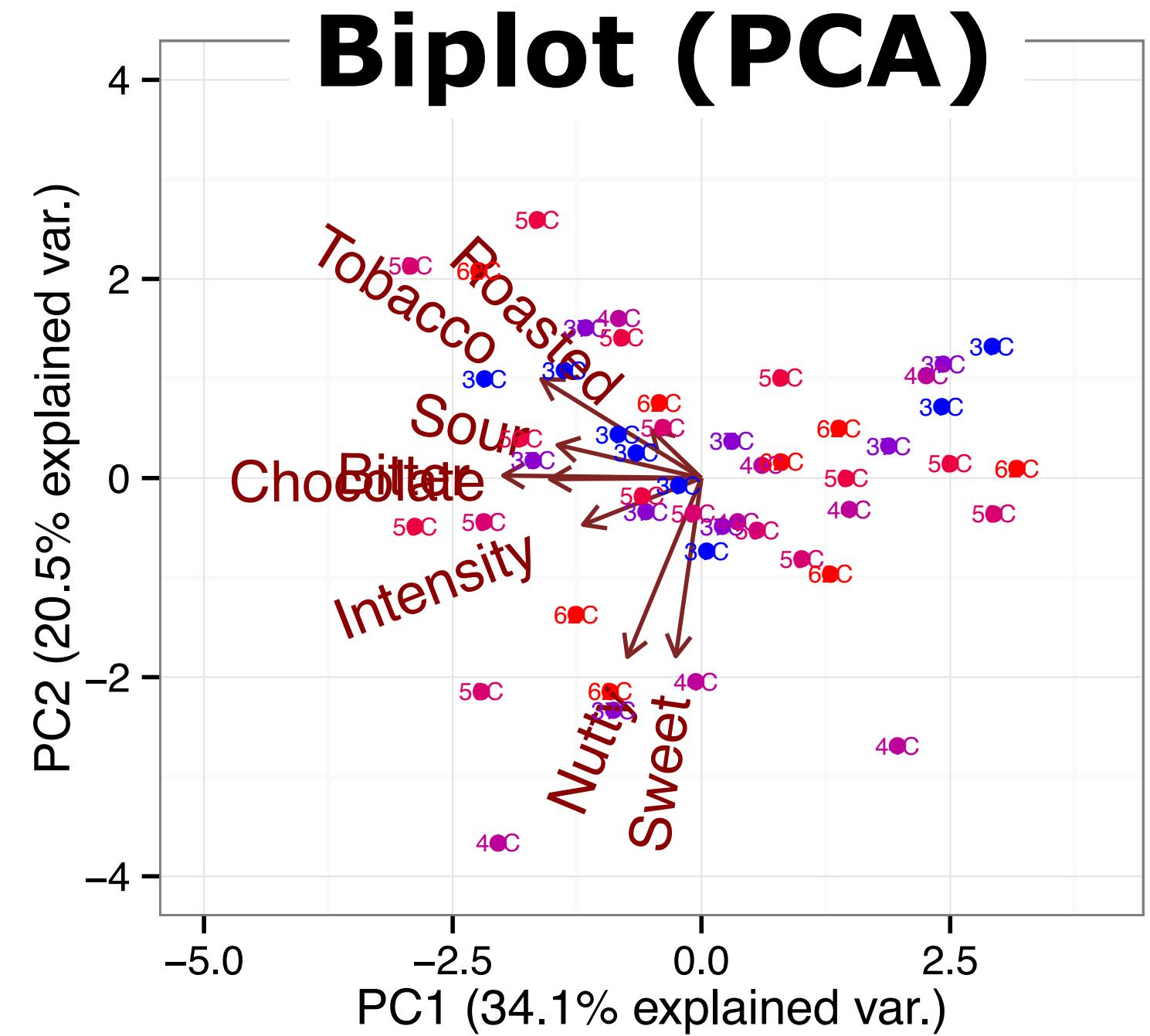
In R it is easily estimated and plotted.

Model

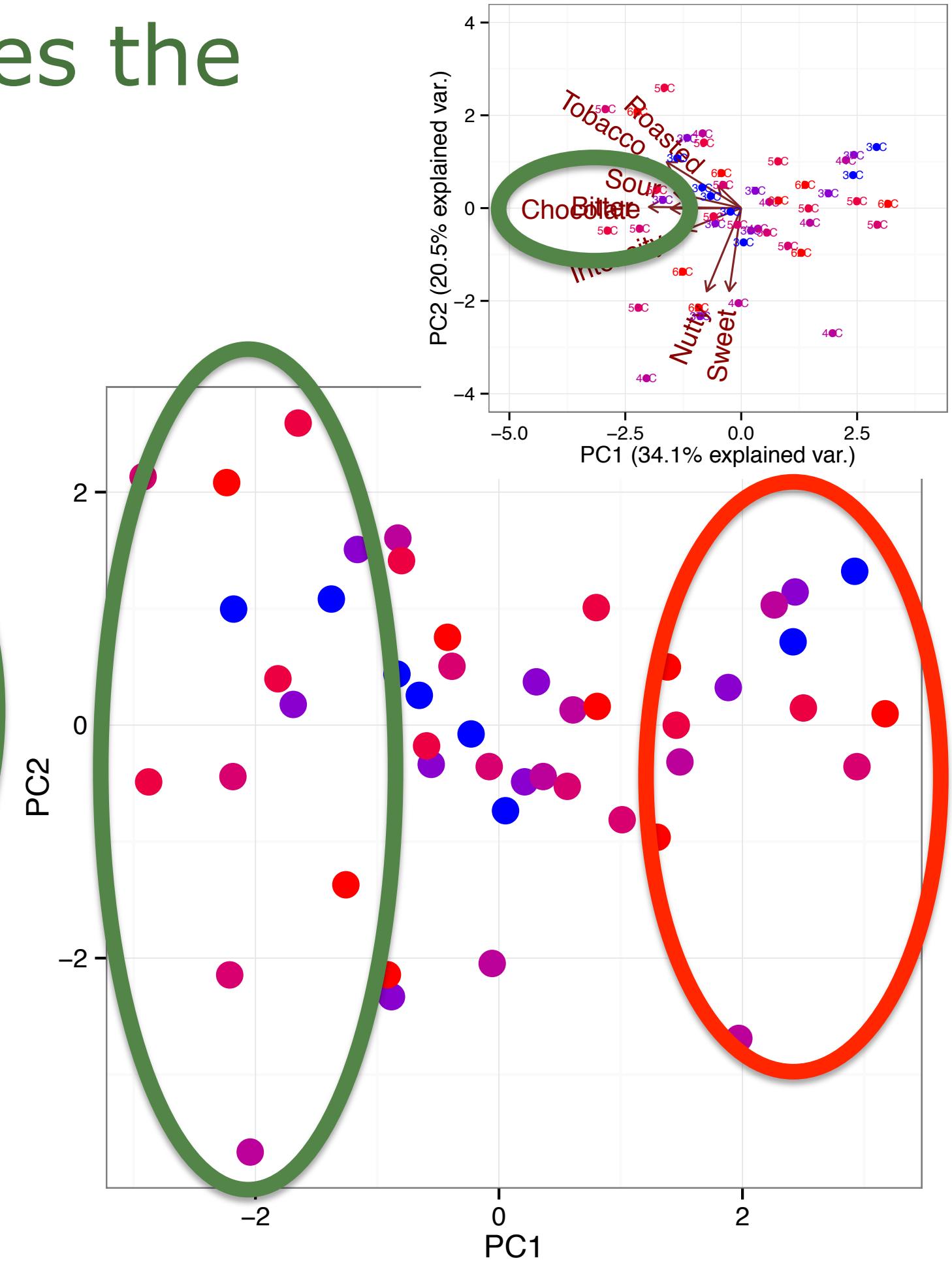
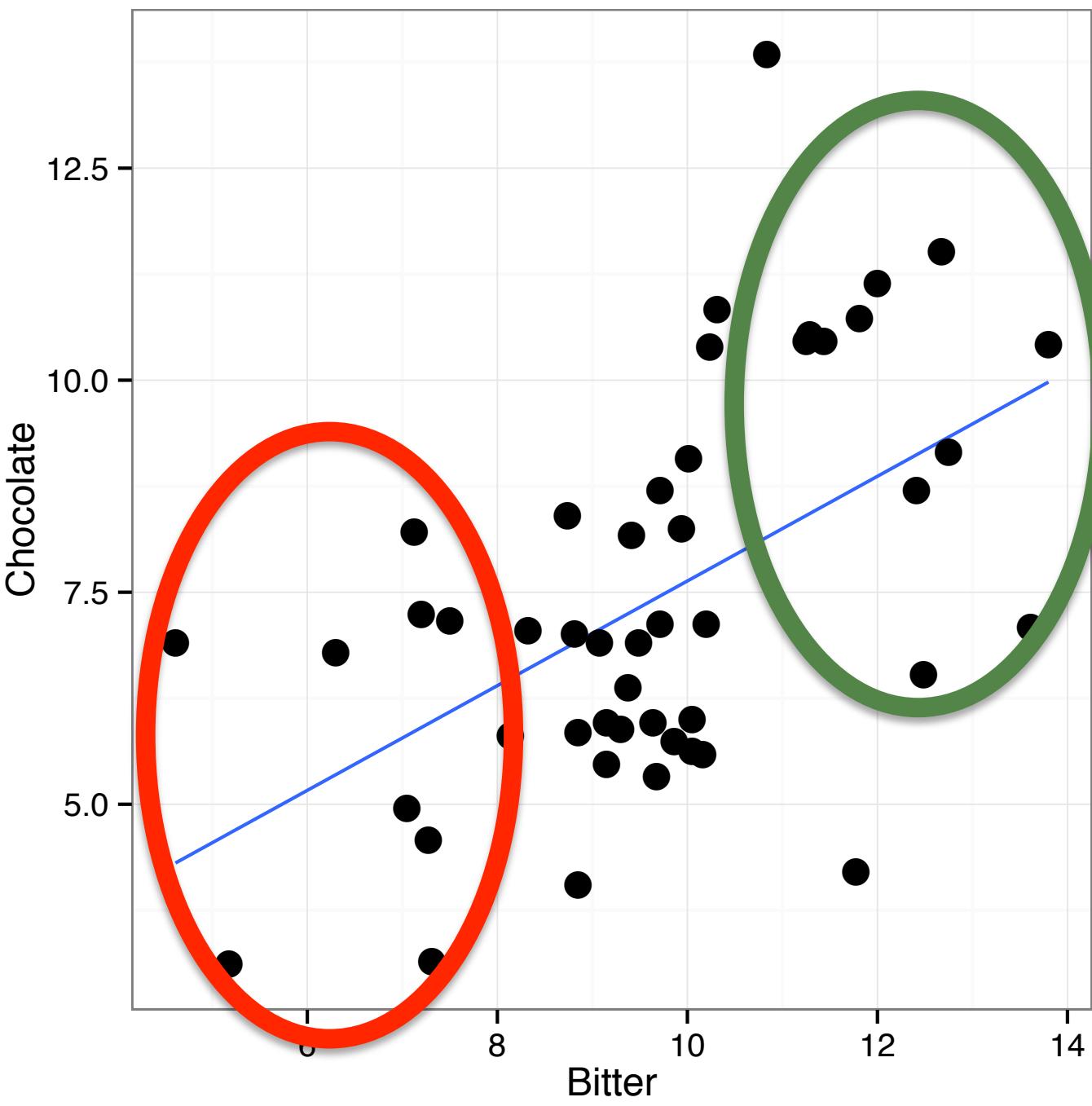
```
> prcomp() [remember scaling]
```

Plot

```
> ggbioplot() [from the ggbioplot package]
```



Which samples drives the association?



PCA versus Correlation analyses

