# MBP intro to statistics bootcamp

## Day 1

Jason Lerch

2018/09/10

# Hello World

The three challenges of statistical inference are[1]:

[1] From Andrew Gelman

# Hello World

The three challenges of statistical inference are[1]:

1. Generalizing from sample to population

[1] From Andrew Gelman

# Hello World

The three challenges of statistical inference are[1]:

1. Generalizing from sample to population

2. Generalizing from control to treatment group

[1] From Andrew Gelman

# Hello World

The three challenges of statistical inference are[1]:

1. Generalizing from sample to population

2. Generalizing from control to treatment group

3. Generalizing from observed measurements to underlying constructs of interest

[1] From Andrew Gelman

# Three laws of statistics

Arthur C. Clarke's three laws[1]:

1. When a distinguished but elderly scientist states that something is possible, he is almost certainly right. When he states that something is impossible, he is very probably wrong.

2. The only way of discovering the limits of the possible is to venture a little way past them into the impossible.

3. Any sufficiently advanced technology is indistinguishable from magic.

# Three laws of statistics

Arthur C. Clarke's three laws[1]:

1. When a distinguished but elderly scientist states that something is possible, he is almost certainly right. When he states that something is impossible, he is very probably wrong.

2. The only way of discovering the limits of the possible is to venture a little way past them into the impossible.

3. Any sufficiently advanced technology is indistinguishable from magic.

Andrew Gelman's updates[2]:

1. When a distinguished but elderly scientist states that "You have no choice but to accept that the major conclusions of these studies are true," don't believe him.

2. The only way of discovering the limits of the reasonable is to venture a little way past them into the unreasonable.

3. Any sufficiently crappy research is indistinguishable from fraud.

# Three laws of statistics

Arthur C. Clarke's three laws[1]:

1. When a distinguished but elderly scientist states that something is possible, he is almost certainly right. When he states that something is impossible, he is very probably wrong.

2. The only way of discovering the limits of the possible is to venture a little way past them into the impossible.

3. Any sufficiently advanced technology is indistinguishable from magic.

Andrew Gelman's updates[2]:

1. When a distinguished but elderly scientist states that "You have no choice but to accept that the major conclusions of these studies are true," don't believe him.

2. The only way of discovering the limits of the reasonable is to venture a little way past them into the unreasonable.

3. Any sufficiently crappy research is indistinguishable from fraud.

[1] https://en.wikipedia.org/wiki/Clarke%27s_three_laws

[2] http://andrewgelman.com/2016/06/20/clarkes-law-of-research/

# The MBP statistics bootcamp

Goals of this week:

1. Teach the theory and practice of statistics

2. Applied data analysis problem solving using R

3. Think hard about truth and replicability in science

# The MBP statistics bootcamp

Goals of this week:

1. Teach the theory and practice of statistics

2. Applied data analysis problem solving using R

3. Think hard about truth and replicability in science

| Hour | Monday | Tuesday | Wednesday | Thursday | Friday |
|------|--------|---------|-----------|----------|--------|
| 9-12 | Introduction, R, visualization, data munging | Linear models, testing proportions, hypothesis tests | | | Presentations, exam |
| 12-1 | | | | | |
| 1-3 | Group assignment #1 | Group assignment #2 | Machine learning, Bayesian statistics | Truth and replicabity. Group assignment #3 | |
| 3-4 | | | | | |

# Grading

Exams (concepts only, no R):

| What | When | How much |
| --- | --- | --- |
| Short exam | Tuesday | 5% |
| Short exam | Wednesday | 5% |
| Short exam | Thursday | 5% |
| Final exam | Friday | 35% |

Group assignments and presentations (R analyses and concepts):

| What | Due when | How much |
| --- | --- | --- |
| Group assignment #1 | Tuesday | 10% |
| Group assignment #2 | Wednesday | 10% |
| Group assignment #3 | Friday | 10% |
| Final presentation | Friday | 20% |

# Exams

- true/false, multiple choice, and short paragraphs.

- each class begins with ~ 10 minute, short exam covering previous day.
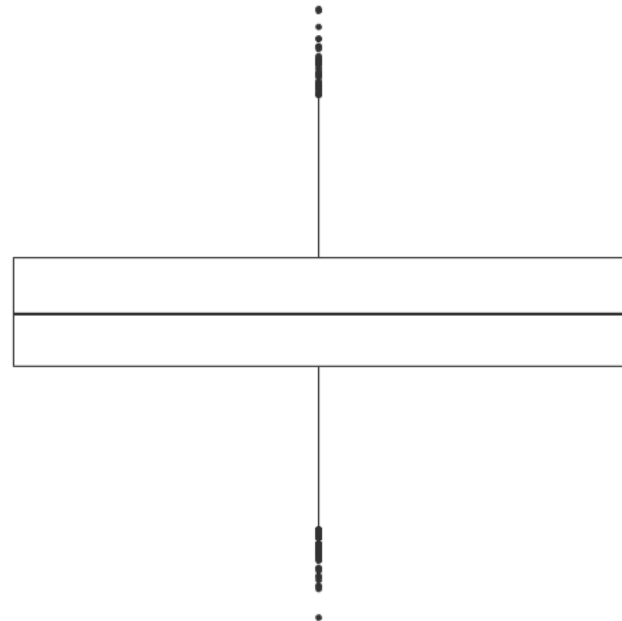
- final exam 30-60 minutes.

# Exams

- true/false, multiple choice, and short paragraphs.

- each class begins with ~ 10 minute, short exam covering previous day.

- final exam 30-60 minutes.

Sample questions:

*Describe the null hypothesis*

*Identify elements of a box and whiskers plot (on a drawing)*

*Discuss analysis pre-registration advantages and disadvantages*

*TRUE/FALSE: if you compute a 95% confidence interval, you have a 95% chance of it containing the true value*

# Group assignments

- split into small groups of 3-4.

- we will assign groups.

- will try to mix groups by R and programming expertise.

- each group will be graded as a unit.

- final presentation given by a member of the group with least R/programming expertise.

# Let's get started

# Statistical software

Common software

1. Excel

2. SPSS

3. SAS

4. matlab

5. python

6. R

# Statistical software

Common software

1. Excel

2. SPSS

3. SAS

4. matlab

5. python

6. R

Ups and downs of R

1. Open source, free, and powerful.

2. If a statistical test exists, it likely exists in R.

3. Literate programming/self documenting analyses.

4. Very strong in bioinformatics.
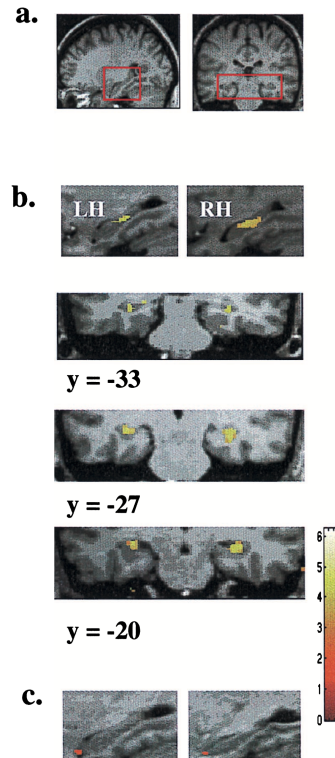
5. Steeper learning curve.

# Intro to R

Over to Mehran

# Reading and summarizing our data

# Intro to our dataset

How do our brains change as we learn or undergo new experiences?

Earliest evidence that our brains are *plastic* at larger, or *mesoscopic*, scales came from a study of taxi drivers in London, UK.

Mechanism of how that happens is unclear.

# Mouse models

We can create taxi driving mice.

Use high-field MRI to get similar readout as in humans.

Use genetic models to test hypotheses of implicated pathways.

Use RNA sequencing to assess what changes per genotype or experimental group.

# The dataset

# The dataset

There are 283 mice in this dataset, with MRI scans acquired at 6 timepoints.

# The dataset

There are 283 mice in this dataset, with MRI scans acquired at 6 timepoints.

We have 3 genotypes: CREB -/-, CREB +/-, CREB +/+

# The dataset

There are 283 mice in this dataset, with MRI scans acquired at 6 timepoints.

We have 3 genotypes: CREB -/-, CREB +/-, CREB +/+

There are 4 environmental conditions: Enriched, Exercise, Isolated Standard, Standard

# The dataset

There are 283 mice in this dataset, with MRI scans acquired at 6 timepoints.

We have 3 genotypes: CREB -/-, CREB +/-, CREB +/+

There are 4 environmental conditions: Enriched, Exercise, Isolated Standard, Standard

MRIs were acquired at every timepoint, and the brains automatically segmented into regions.

# The dataset

There are 283 mice in this dataset, with MRI scans acquired at 6 timepoints.

We have 3 genotypes: CREB -/-, CREB +/-, CREB +/+

There are 4 environmental conditions: Enriched, Exercise, Isolated Standard, Standard

MRIs were acquired at every timepoint, and the brains automatically segmented into regions.

There are good reasons to believe that the hippocampus and the dentate gyrus of the hippocampus will be the most affected by the environmental interventions.

# The dataset

There are 283 mice in this dataset, with MRI scans acquired at 6 timepoints.

We have 3 genotypes: CREB -/-, CREB +/-, CREB +/+

There are 4 environmental conditions: Enriched, Exercise, Isolated Standard, Standard

MRIs were acquired at every timepoint, and the brains automatically segmented into regions.

There are good reasons to believe that the hippocampus and the dentate gyrus of the hippocampus will be the most affected by the environmental interventions.

The effect of the three genotypes alone is interesting.

# The dataset

There are 283 mice in this dataset, with MRI scans acquired at 6 timepoints.

We have 3 genotypes: CREB -/-, CREB +/-, CREB +/+

There are 4 environmental conditions: Enriched, Exercise, Isolated Standard, Standard

MRIs were acquired at every timepoint, and the brains automatically segmented into regions.

There are good reasons to believe that the hippocampus and the dentate gyrus of the hippocampus will be the most affected by the environmental interventions.

The effect of the three genotypes alone is interesting.

A separate cohort of mice was used for the RNA-seq experiment (but we'll get to that later in the course).

# Enrichment

# Reading data

A surprising amount of time in data analysis is spent in prepping data for
visualization and analysis.

```r
library(tidyverse)
library(forcats)

mice <- read_csv("mice.csv")
```

```
## Parsed with column specification:
## cols(
##   Age = col_double(),
##   Sex = col_character(),
##   Condition = col_character(),
##   Mouse.Genotyping = col_character(),
##   ID = col_integer(),
##   Timepoint = col_character(),
##   Genotype = col_character(),
##   DaysOfEE = col_integer(),
##   DaysOfEE0 = col_integer()
## )
```

# Meet the mice

```
str(mice, give.attr=FALSE)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':    1392 obs. of  9 variables:
##  $ Age            : num  8.5 8.5 8.5 9.5 9.5 8.5 8.5 9.5 8.5 9.5 ...
##  $ Sex            : chr  "M" "M" "M" "M" ...
##  $ Condition      : chr  "Enriched" "Standard" "Standard" "Enriched" ...
##  $ Mouse.Genotyping: chr  "Heterozygous" "Heterozygous" "Heterozygous" "Wi
##  $ ID             : int  901 899 898 891 893 901 899 889 898 895 ...
##  $ Timepoint      : chr  "Pre1" "Pre1" "Pre1" "Pre1" ...
##  $ Genotype       : chr  "CREB +/-" "CREB +/-" "CREB +/-" "CREB +/+" ...
##  $ DaysOfEE       : int  -4 -4 -4 -4 -4 -3 -3 -3 -3 -3 ...
##  $ DaysOfEE0      : int  0 0 0 0 0 0 0 0 0 0 ...
```

# Numeric variable: age

```
mice %>%
  summarise(mean=mean(Age),
            min=min(Age),
            max=max(Age))
```

| mean | min | max |
|------|-----|-----|
| 6.58 | 3.1 | 10.1 |

# Intro to pipes?

# Factors: Sex, Condition, Genotype

```
mice %>%
  group_by(Sex) %>%
  summarise(n=n())
```

```
mice %>%
  group_by(Genotype) %>%
  summarize(n=n())
```

| Sex | n |
|-----|-----|
| F | 543 |
| M | 849 |

| Genotype | n |
|----------|-----|
| CREB -/- | 426 |
| CREB +/- | 486 |
| CREB +/+ | 480 |

# Subject descriptors: ID and Timepoint

```
mice %>%
  select(ID, Timepoint) %>%
  head
```

| ID | Timepoint |
|----|-----------|
| 901 | Pre1 |
| 899 | Pre1 |
| 898 | Pre1 |
| 891 | Pre1 |
| 893 | Pre1 |
| 901 | Pre2 |

# Alternate encodings: Genotype

```
mice %>%
  select(Genotype, Mouse.Genotyping) %>%
  head
```

| Genotype | Mouse.Genotyping |
|----------|------------------|
| CREB +/- | Heterozygous |
| CREB +/- | Heterozygous |
| CREB +/- | Heterozygous |
| CREB +/+ | Wildtype |
| CREB +/+ | Wildtype |
| CREB +/- | Heterozygous |

# Alternate encodings: Days of EE, DaysofEE0

```
mice %>%
  filter(ID == 901) %>%
  select(Timepoint, DaysOfEE, DaysOfEE0) %>%
  head
```

| Timepoint | DaysOfEE | DaysOfEE0 |
|-----------|---------:|----------:|
| Pre1 | -4 | 0 |
| Pre2 | -3 | 0 |
| 24h | 1 | 1 |
| 48h | 2 | 2 |
| 1 week | 8 | 8 |
| 2 week | 16 | 16 |

# Overview of subject numbers

```
with(mice,
     ftable(Condition, Genotype, Timepoint))
```

```
##                           Timepoint 1 week 2 week 24h 48h Pre1 Pre2
## Condition        Genotype
## Enriched         CREB -/-             24     25   7  24   24   24
##                  CREB +/-             30     33  12  34   30   33
##                  CREB +/+             27     30   8  30   27   28
## Exercise         CREB -/-             22     21   0  21   20   18
##                  CREB +/-             17     18   0  18   17   15
##                  CREB +/+             19     19   0  19   19   16
## Isolated Standard CREB -/-            14     14   0  14   13   14
##                  CREB +/-             12     12   0  12   11   12
##                  CREB +/+             17     17   0  17   17   14
## Standard         CREB -/-             23     26   4  26   25   23
##                  CREB +/-             29     34   9  34   32   32
##                  CREB +/+             28     31   6  31   31   29
```

# Factors, revisited

The Timepoint order makes no sense. Let's reorder

```
mice <- mice %>%
  mutate(Timepoint=fct_relevel(Timepoint, "Pre1", "Pre2", "24h",
                               "48h", "1 week", "2 week"))
with(mice, ftable(Condition, Genotype, Timepoint))
```

```
##                            Timepoint Pre1 Pre2 24h 48h 1 week 2 week
## Condition         Genotype
## Enriched          CREB -/-            24   24   7  24     24     25
##                   CREB +/-            30   33  12  34     30     33
##                   CREB +/+            27   28   8  30     27     30
## Exercise          CREB -/-            20   18   0  21     22     21
##                   CREB +/-            17   15   0  18     17     18
##                   CREB +/+            19   16   0  19     19     19
## Isolated Standard CREB -/-            13   14   0  14     14     14
##                   CREB +/-            11   12   0  12     12     12
##                   CREB +/+            17   14   0  17     17     17
## Standard          CREB -/-            25   23   4  26     23     26
##                   CREB +/-            32   32   9  34     29     34
##                   CREB +/+            31   29   6  31     28     31
```

# Redo in tidyverse

```
mice %>%
  group_by(Condition, Genotype, Timepoint) %>%
  summarise(n=n()) %>% spread(Timepoint, value=n)
```

```
## # A tibble: 12 x 8
## # Groups:   Condition, Genotype [12]
##    Condition         Genotype  Pre1  Pre2 `24h` `48h` `1 week` `2 week`
##    <chr>             <chr>    <int> <int> <int> <int>    <int>    <int>
##  1 Enriched          CREB -/-    24    24     7    24       24       25
##  2 Enriched          CREB +/-    30    33    12    34       30       33
##  3 Enriched          CREB +/+    27    28     8    30       27       30
##  4 Exercise          CREB -/-    20    18    NA    21       22       21
##  5 Exercise          CREB +/-    17    15    NA    18       17       18
##  6 Exercise          CREB +/+    19    16    NA    19       19       19
##  7 Isolated Standard CREB -/-    13    14    NA    14       14       14
##  8 Isolated Standard CREB +/-    11    12    NA    12       12       12
##  9 Isolated Standard CREB +/+    17    14    NA    17       17       17
## 10 Standard          CREB -/-    25    23     4    26       23       26
## 11 Standard          CREB +/-    32    32     9    34       29       34
## 12 Standard          CREB +/+    31    29     6    31       28       31
```

# Reading more data

```
volumes <- read_csv("volumes.csv")
```

```
## Parsed with column specification:
## cols(
##   .default = col_double(),
##   ID = col_integer(),
##   Timepoint = col_character()
## )

## See spec(...) for full column specifications.
```

# Inspecting the new data

```
str(volumes)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':    1392 obs. of  161 variables:
##  $ amygdala                        : num  9.84 10.3 10.5
##  $ anterior commissure: pars anterior  : num  1.42 1.48 1.5
##  $ anterior commissure: pars posterior : num  0.392 0.428 0.
##  $ basal forebrain                 : num  4.72 4.96 4.93
##  $ bed nucleus of stria terminalis : num  1.24 1.31 1.28
##  $ cerebellar peduncle: inferior   : num  0.908 0.967 0.
##  $ cerebellar peduncle: middle     : num  1.23 1.31 1.26
##  $ cerebellar peduncle: superior   : num  0.991 0.848 0.
##  $ cerebral aqueduct               : num  0.373 0.44 0.4
##  $ cerebral peduncle               : num  2.58 2.54 2.6
##  $ colliculus: inferior            : num  5.46 5.6 5.34
##  $ colliculus: superior            : num  9.52 9.83 9.37
##  $ corpus callosum                 : num  14.1 14.3 14 1
##  $ corticospinal tract/pyramids    : num  1.59 1.59 1.56
##  $ cuneate nucleus                 : num  0.27 0.254 0.2
##  $ dentate gyrus of hippocampus    : num  3.96 3.92 3.95
##  $ facial nerve (cranial nerve 7)  : num  0.221 0.233 0.
##  $ fasciculus retroflexus          : num  0.214 0.212 0.
##  $ fimbria                         : num  2.67 2.9 2.73
```

# Linking data

```
volumes %>%
  select(ID, Timepoint) %>%
  head
```

```
mice %>%
  select(ID, Timepoint) %>%
  head
```

| ID | Timepoint |
|----|-----------|
| 901 | Pre1 |
| 899 | Pre1 |
| 898 | Pre1 |
| 891 | Pre1 |
| 893 | Pre1 |
| 901 | Pre2 |

| ID | Timepoint |
|----|-----------|
| 901 | Pre1 |
| 899 | Pre1 |
| 898 | Pre1 |
| 891 | Pre1 |
| 893 | Pre1 |
| 901 | Pre2 |

# Joining data

```
mice <- mice %>%
  inner_join(volumes)
```

```
## Joining, by = c("ID", "Timepoint")

## Warning: Column `Timepoint` joining factor and character vector, coercing
## into character vector
```

```
str(mice)
```

```
## Classes 'tbl_df', 'tbl' and 'data.frame':    1392 obs. of  168 variables:
##  $ Age               : num  8.5 8.5 8.5 9.
##  $ Sex               : chr  "M" "M" "M" "M
##  $ Condition         : chr  "Enriched" "St
##  $ Mouse.Genotyping  : chr  "Heterozygous"
##  $ ID                : int  901 899 898 89
##  $ Timepoint         : chr  "Pre1" "Pre1"
##  $ Genotype          : chr  "CREB +/-" "CR
##  $ DaysOfEE          : int  -4 -4 -4 -4 -4
##  $ DaysOfEE0         : int  0 0 0 0 0 0 0
##  $ amygdala          : num  9.84 10.3 10.5
```

# Data visualization

# Data visualization

Data visualization communicates your data to your audience - and can be how your data communicates with you.

# Data visualization

Data visualization communicates your data to your audience - and can be how your data communicates with you.

Excellent guide to visualization:

https://www.data-to-viz.com

# Data visualization

Data visualization communicates your data to your audience - and can be how your data communicates with you.

Excellent guide to visualization:

https://www.data-to-viz.com

Your task for later will be to look at the interesting variables in this dataset. For now, we will look at sex and the brain instead.

# Histogram

```
ggplot(mice) +
  aes(x=`bed nucleus of stria terminalis`) +
  geom_histogram()
```

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.

# Make it prettier

```
ggplot(mice) +
  aes(x=`bed nucleus of stria terminalis`) +
  geom_histogram() +
  xlab(bquote(Volume ~ (mm^3))) +
  ggtitle("Bed nucleus of stria terminalis") +
  theme_gray(16)
```
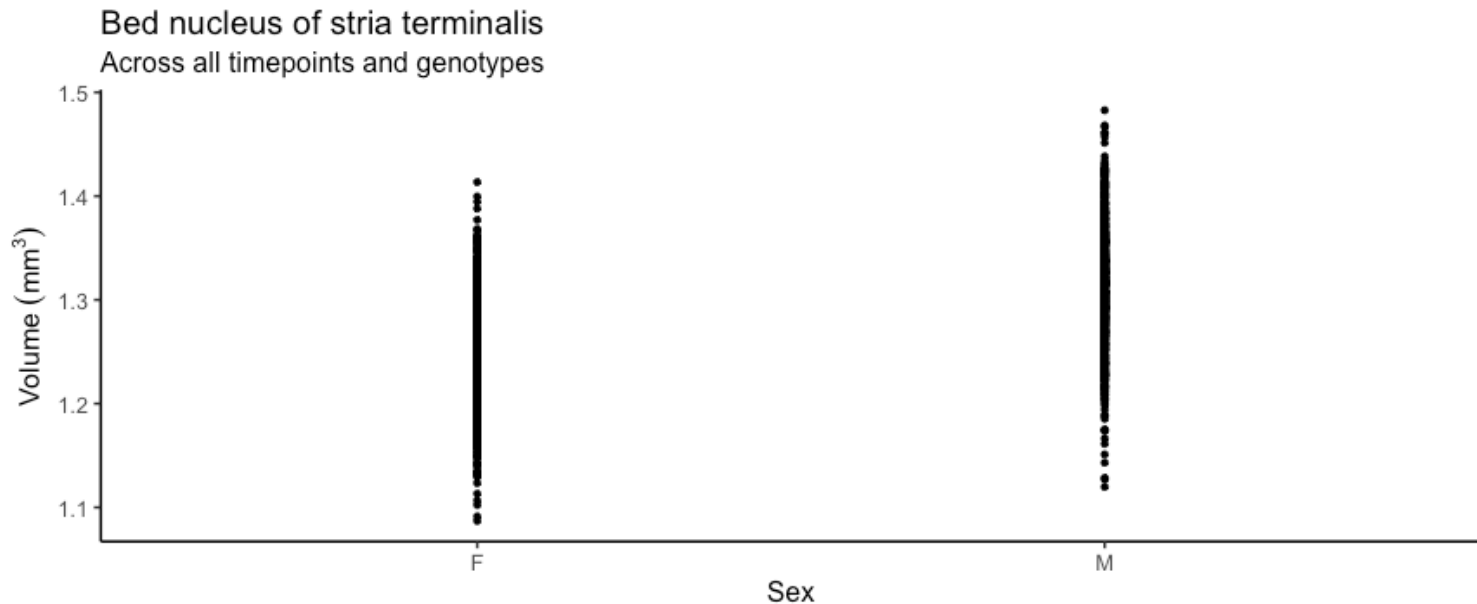
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

# Histogram bins

```
ggplot(mice) +
  aes(x=`bed nucleus of stria terminalis`) +
  geom_histogram(binwidth = 0.01) +
  xlab(bquote(Volume ~ (mm^3))) +
  ggtitle("Bed nucleus of stria terminalis") +
  theme_gray(16)
```

# Facets

```
ggplot(mice) +
  aes(x=`bed nucleus of stria terminalis`) +
  geom_histogram(binwidth = 0.01) +
  xlab(bquote(Volume ~ (mm^3))) +
  ggtitle("Bed nucleus of stria terminalis") +
  theme_gray(16) +
  facet_grid(Sex ~ .)
```

# Colours

```
ggplot(mice) +
  aes(x=`bed nucleus of stria terminalis`, fill=Sex) +
  geom_histogram(binwidth = 0.01) +
  xlab(bquote(Volume ~ (mm^3))) +
  ggtitle("Bed nucleus of stria terminalis") +
  theme_gray(16)
```

Bed nucleus of stria terminalis

# Points

```
ggplot(mice) +
  aes(x=Sex, y=`bed nucleus of stria terminalis`) +
  geom_point() +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all timepoints and genotypes") +
  ylab(bquote(Volume ~ (mm^3))) +
  theme_classic(16)
```
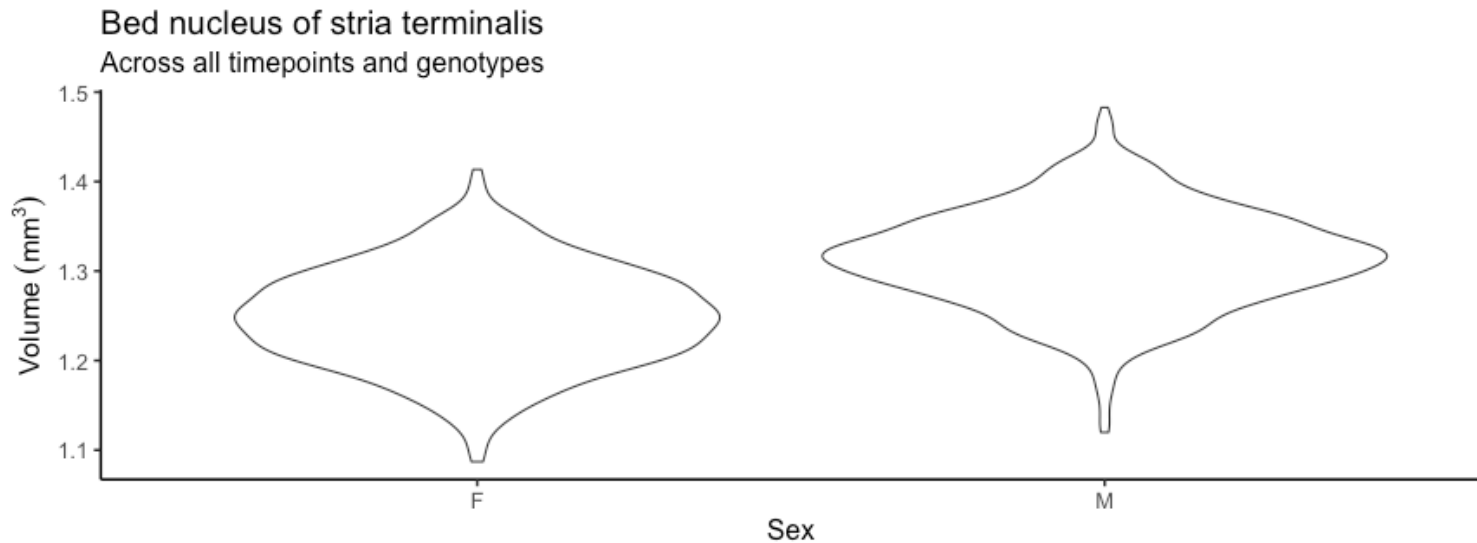
# Points

That's not very useful - too many points to see separation.

```
ggplot(mice) +
  aes(x=Sex, y=`bed nucleus of stria terminalis`) +
  geom_jitter() +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all timepoints and genotypes") +
  ylab(bquote(Volume ~ (mm^3))) +
  theme_classic(16)
```



Bed nucleus of stria terminalis
Across all timepoints and genotypes
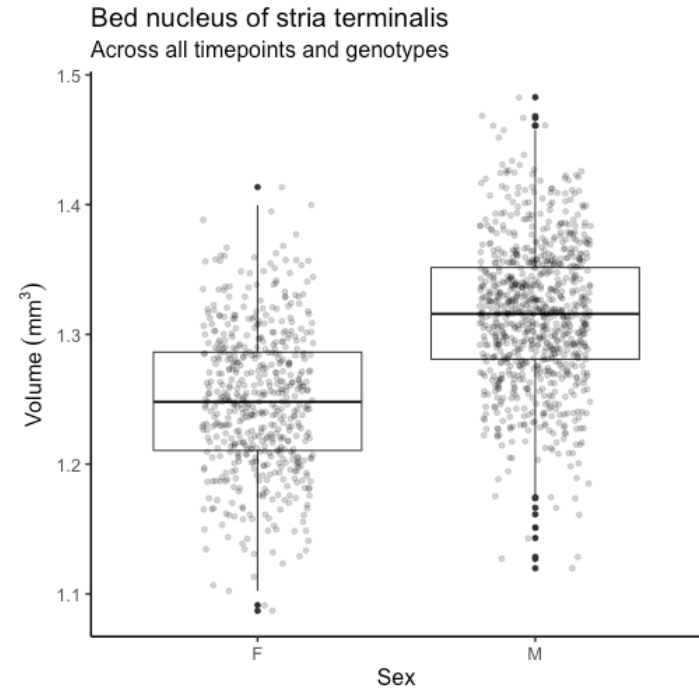
# Boxplot

Good view of data distribution

```
ggplot(mice) +
  aes(x=Sex, y=`bed nucleus of stria terminalis`) +
  geom_boxplot() +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all timepoints and genotypes") +
  ylab(bquote(Volume ~ (mm^3))) +
  theme_classic(16)
```

# Ridge lines

```
suppressMessages(library(ggridges))
ggplot(mice) +
  aes(y=Sex, x=`bed nucleus of stria terminalis`) +
  geom_density_ridges() +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all timepoints and genotypes") +
  xlab(bquote(Volume ~ (mm^3))) +
  theme_classic(16)
```

```
## Picking joint bandwidth of 0.0132
```

# Violins

```
ggplot(mice) +
  aes(x=Sex, y=`bed nucleus of stria terminalis`) +
  geom_violin() +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all timepoints and genotypes") +
  ylab(bquote(Volume ~ (mm^3))) +
  theme_classic(16)
```
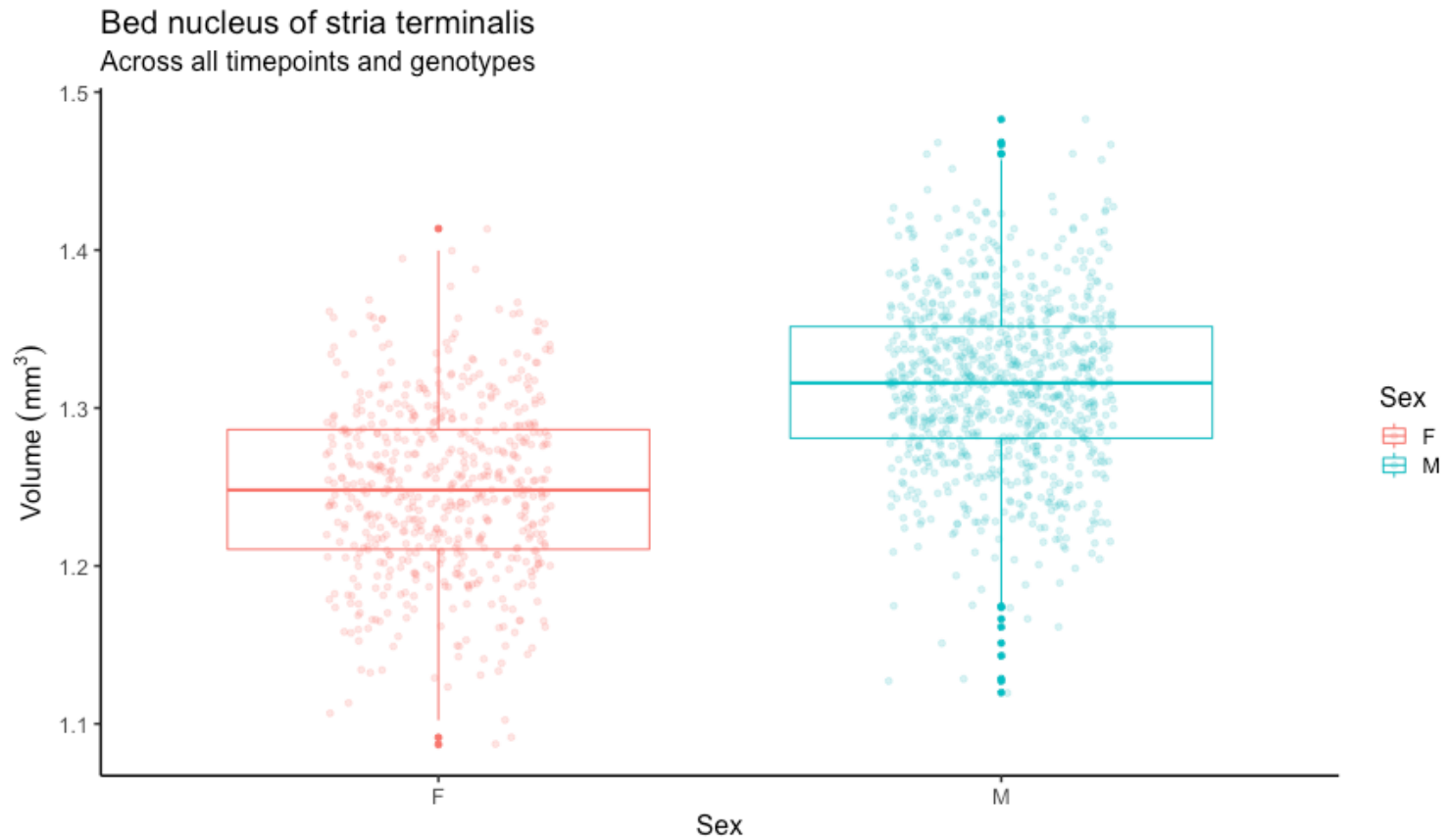
# Combining plot types

```r
ggplot(mice) +
  aes(x=Sex,
      y=`bed nucleus of stria terminalis`
      ) +
  geom_boxplot() +
  geom_jitter(width=0.2,
              alpha=0.2) +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all timepoints and genotypes")
  ylab(bquote(
    Volume ~ (mm^3))) +
  theme_classic(16)
```
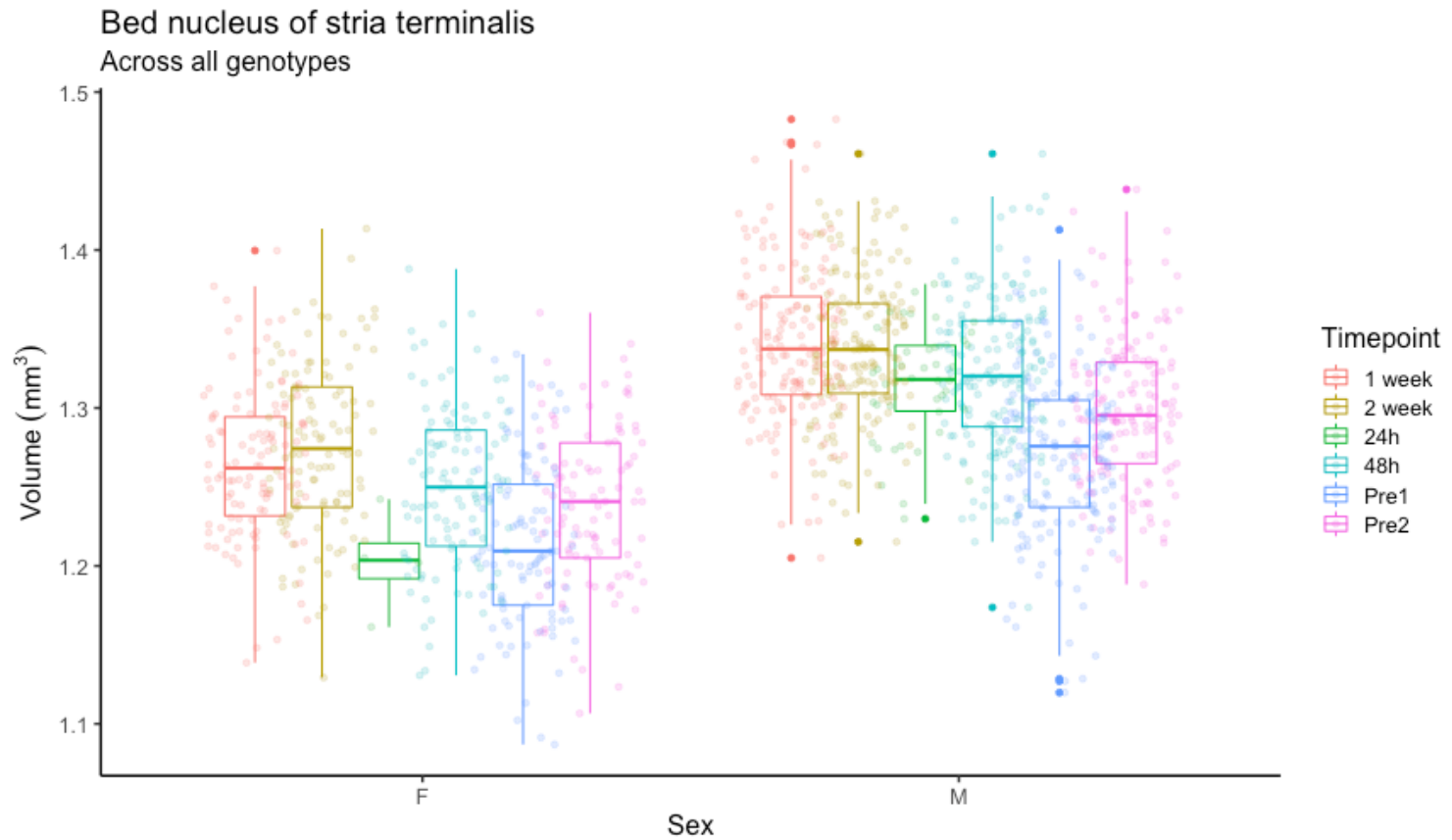


Bed nucleus of stria terminalis
Across all timepoints and genotypes

# Adding colour

```r
ggplot(mice) +
  aes(x=Sex,
      y=`bed nucleus of stria terminalis`,
      colour=Sex) +
  geom_boxplot() +
  geom_jitter(width=0.2,
              alpha=0.2) +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all timepoints and genotypes") +
  ylab(bquote(
    Volume ~ (mm^3))) +
  theme_classic(16)
```

# Adding colour

# Using colour for additional information

```
ggplot(mice) +
  aes(x=Sex,
      y=`bed nucleus of stria terminalis`,
      colour=Timepoint) +
  geom_boxplot() +
  geom_jitter(alpha=0.2,
              position = position_jitterdodge(jitter.width = 0.2)) +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all genotypes") +
  ylab(bquote(Volume ~ (mm^3))) +
  theme_classic(16)
```
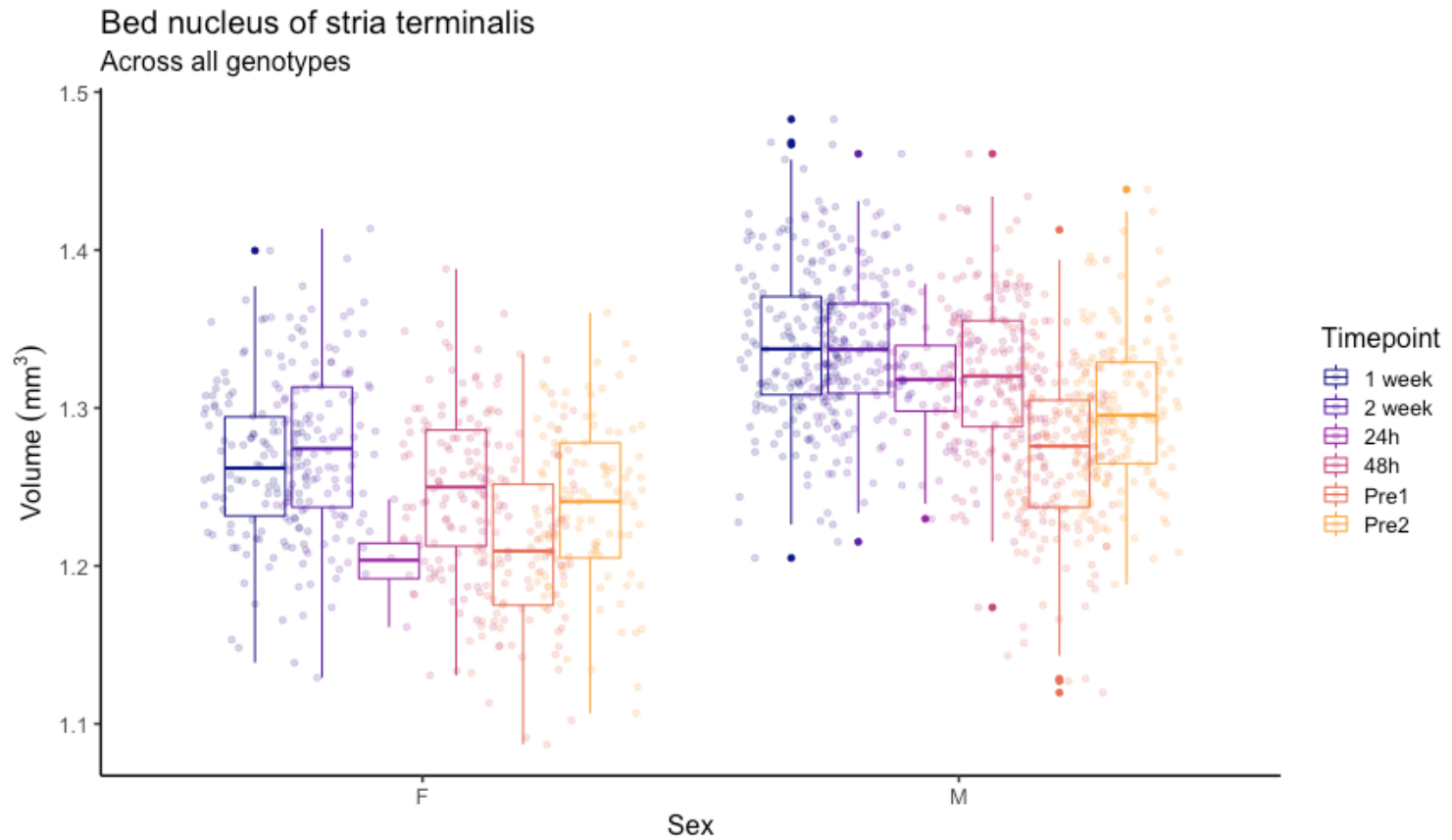
# Using colour for additional information



Bed nucleus of stria terminalis
Across all genotypes

# Using colour for additional information

```
ggplot(mice) +
  aes(x=Sex,
      y=`bed nucleus of stria terminalis`,
      colour=Timepoint) +
  geom_boxplot() +
  geom_jitter(alpha=0.2,
              position = position_jitterdodge(jitter.width = 0.2)) +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all genotypes") +
  ylab(bquote(Volume ~ (mm^3))) +
  scale_colour_viridis_d(option="C", end=0.8) +
  theme_classic(16)
```

# Using colour for additional information



Bed nucleus of stria terminalis
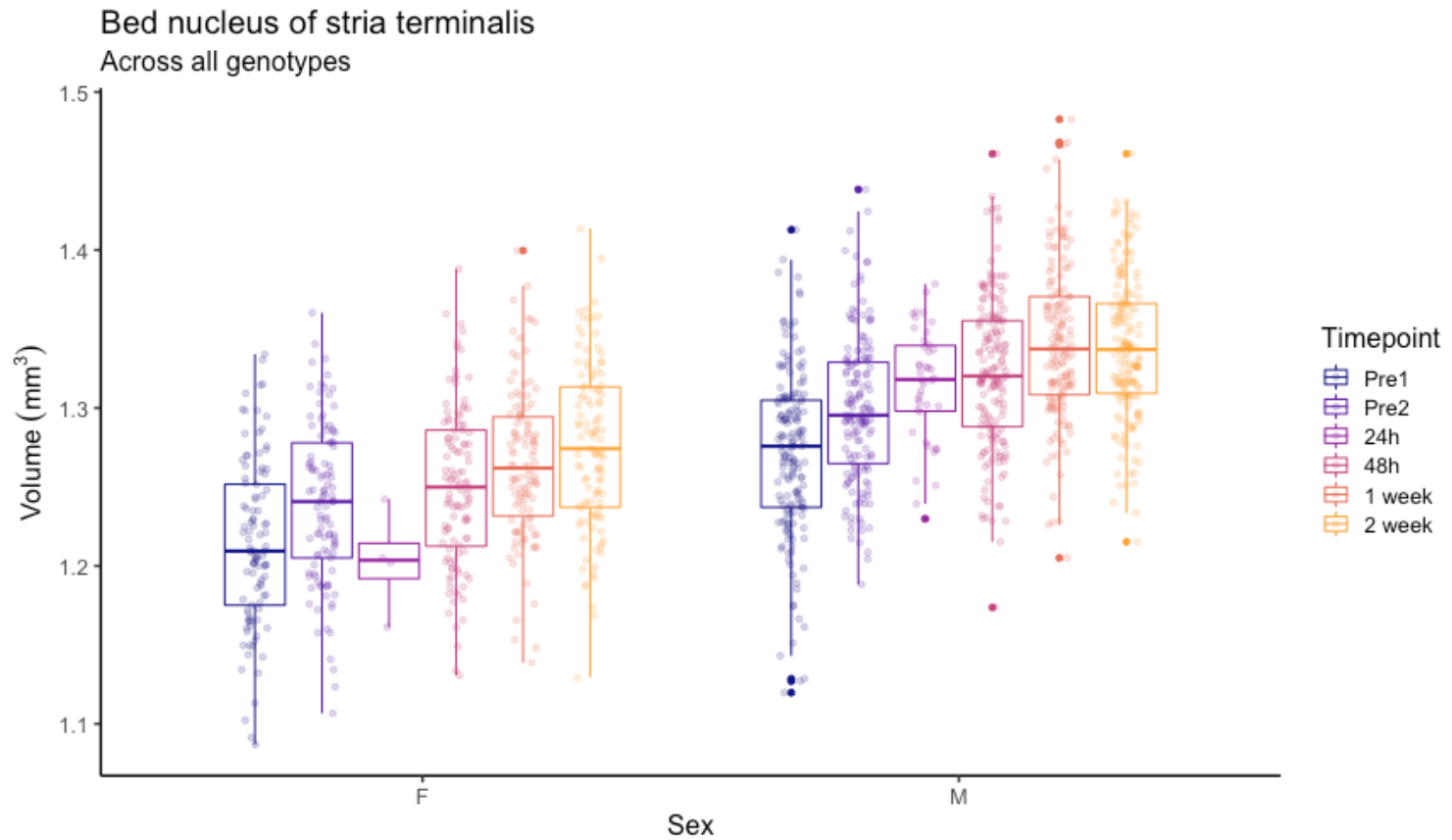Across all genotypes

# Factor order, again

Apparently the factor ordering was lost in data joining?

```
mice <- mice %>%
  mutate(Timepoint=fct_relevel(Timepoint, "Pre1", "Pre2", "24h",
                               "48h", "1 week", "2 week"))

ggplot(mice) +
  aes(x=Sex,
      y=`bed nucleus of stria terminalis`,
      colour=Timepoint) +
  geom_boxplot() +
  geom_jitter(alpha=0.2,
              position = position_jitterdodge(jitter.width = 0.2)) +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle="Across all genotypes") +
  ylab(bquote(Volume ~ (mm^3))) +
  scale_colour_viridis_d(option="C", end=0.8) +
  theme_classic(16)
```
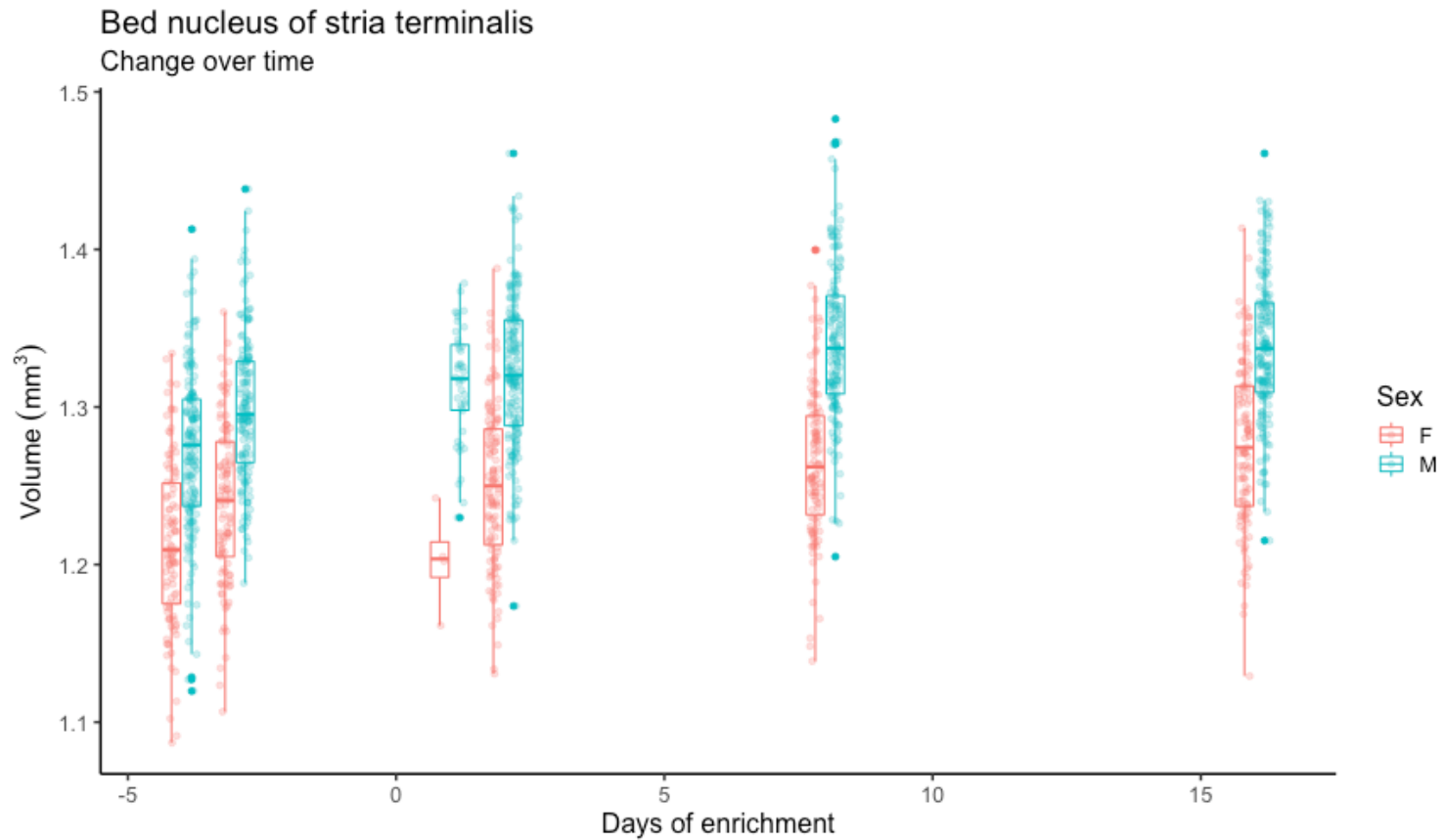
# Factor ordering, again

# Better encoding of time

```r
ggplot(mice) +
  aes(x=DaysOfEE,
      y=`bed nucleus of stria terminalis`,
      colour=Sex) +
  geom_boxplot(aes(group=interaction(Timepoint, Sex))) +
  geom_jitter(alpha=0.25, position =
                position_jitterdodge(jitter.width = 0.2)) +
  ylab(bquote(Volume ~ (mm^3))) +
  xlab("Days of enrichment") +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle = "Change over time") +
  theme_classic(16)
```
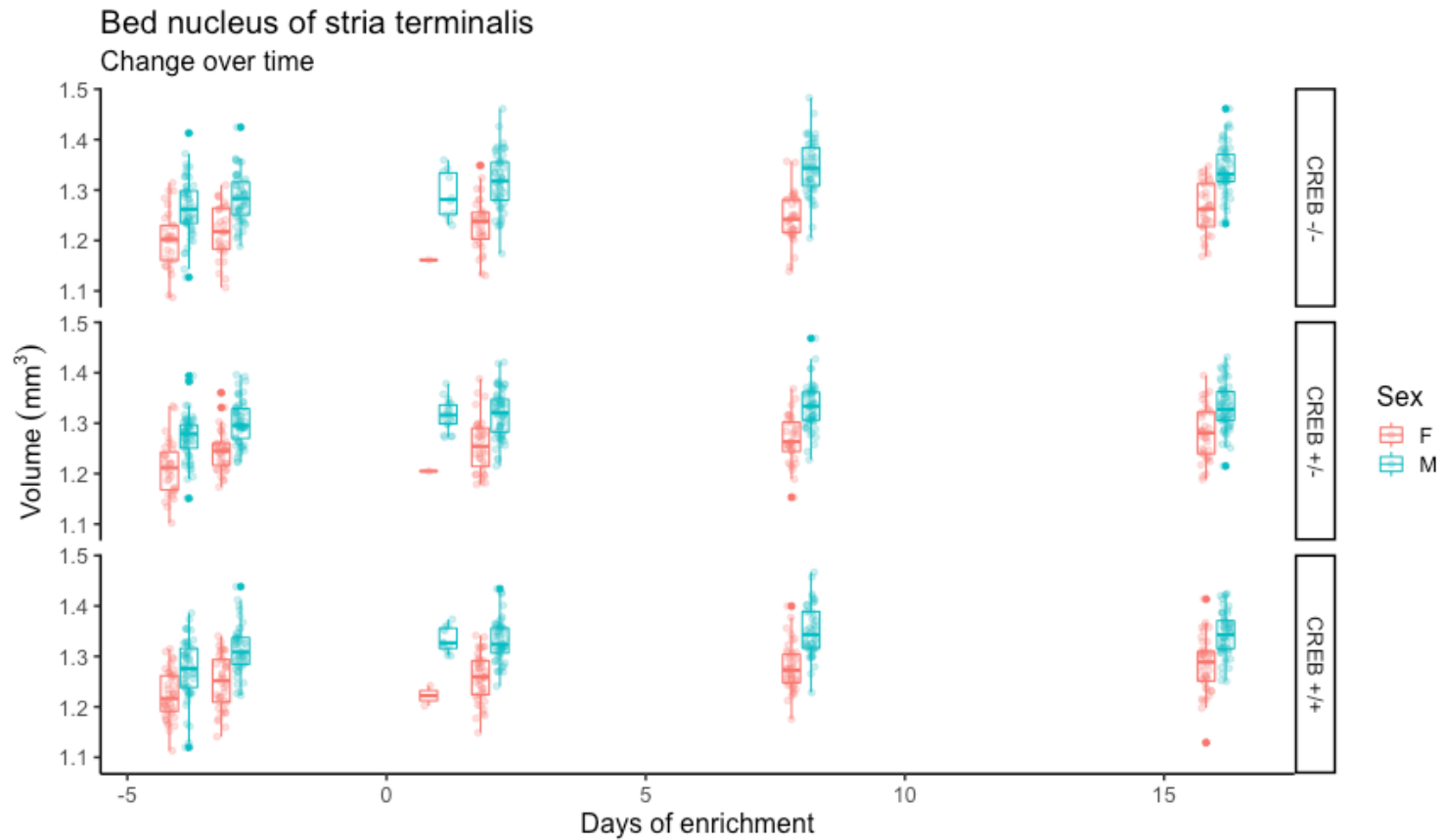
# Better encoding of time

# Combining colours and facets

```
ggplot(mice) +
  aes(x=DaysOfEE,
      y=`bed nucleus of stria terminalis`,
      colour=Sex) +
  geom_boxplot(aes(group=interaction(Timepoint, Sex))) +
  geom_jitter(alpha=0.25, position =
                position_jitterdodge(jitter.width = 0.2)) +
  ylab(bquote(Volume ~ (mm^3))) +
  xlab("Days of enrichment") +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle = "Change over time") +
  facet_grid(Genotype ~ .) +
  theme_classic(16)
```
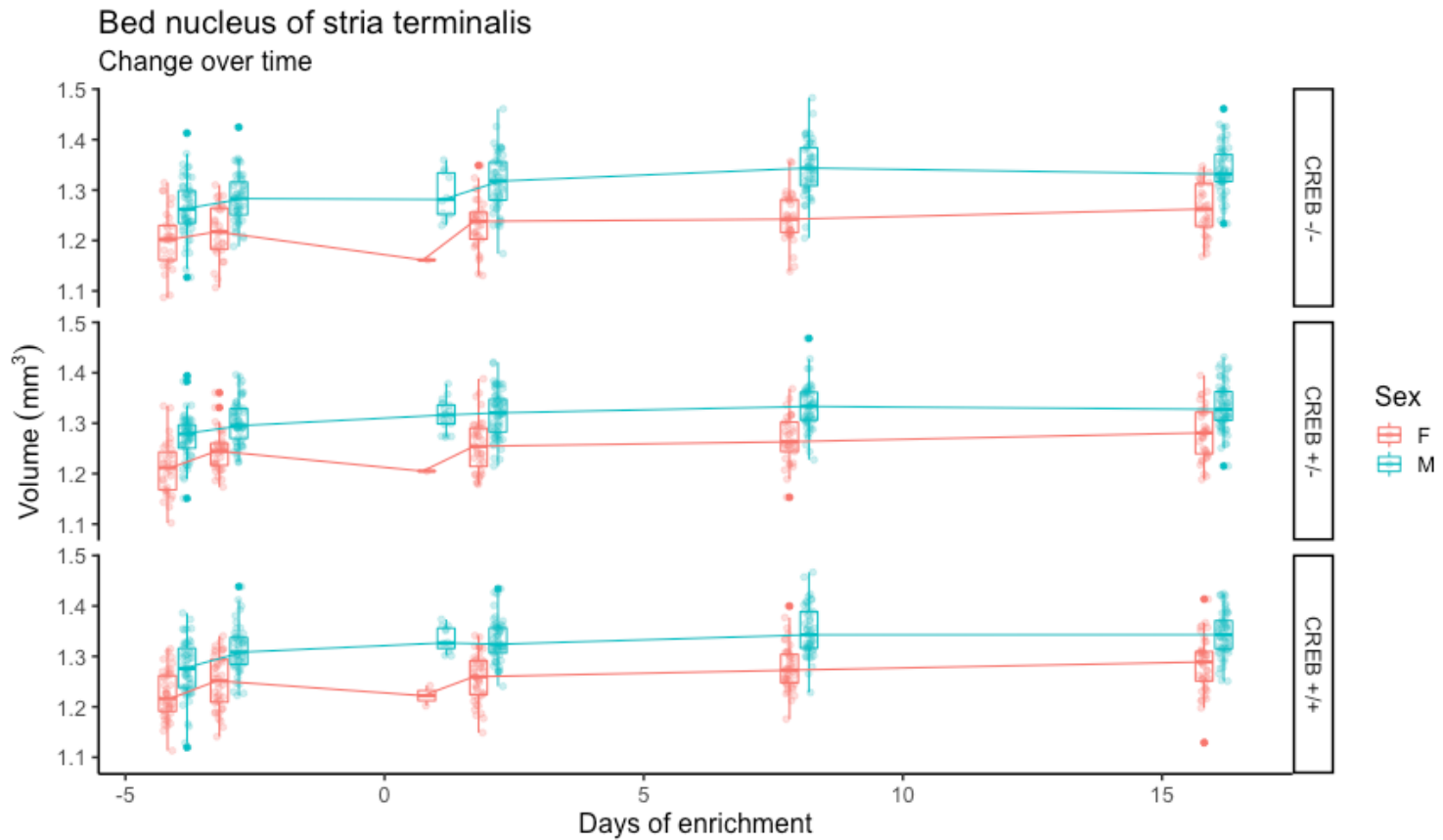
# Combining colours and facets

# Adding lines

```
ggplot(mice) +
  aes(x=DaysOfEE,
      y=`bed nucleus of stria terminalis`,
      colour=Sex) +
  geom_boxplot(aes(group=interaction(Timepoint, Sex))) +
  geom_jitter(alpha=0.25, position =
                position_jitterdodge(jitter.width = 0.2)) +
  stat_summary(fun.y = median, geom="line",
               position =
                 position_jitterdodge(jitter.width = 0.2)) +
  ylab(bquote(Volume ~ (mm^3))) +
  xlab("Days of enrichment") +
  ggtitle("Bed nucleus of stria terminalis",
          subtitle = "Change over time") +
  facet_grid(Genotype ~ .) +
  theme_classic(16)
```

# Adding lines



Bed nucleus of stria terminalis
Change over time

# Literate programming

# Literate programming

## The Idea:

- mix code, text, and figures in one document.

- All analyses and their outputs remain in sync

- Can work as a notebook

# Literate programming

## The Idea:

- mix code, text, and figures in one document.

- All analyses and their outputs remain in sync

- Can work as a notebook

## The Implementation:

- rmarkdown

  - simple markup language for text
  - code embedded in document
  - documents are compiled - or knitted - to produce output html or pdf

- Great alternative: Jupyter

# Assignment

# Group assignment number 1

1. Assemble into your assigned teams.

2. Ensure that RStudio is running and you can load all required libraries.

3. Load the required data

4. Create an rmarkdown document that contains the following:

    1. A summary table of the subject numbers per timepoint, genotype, and condition

    2. Visualization(s) of the difference in hippocampal volume by Genotype at the final timepoint.

    3. Visualization(s) of the difference in hippocampal volume by Condition at the final timepoint.

    4. Visualization(s) of the change over time by Condition and Genotype.

5. Make sure that all team members are listed as authors.

6. Any questions: ask here in person, or email us (jason.lerch@utoronto.ca, mehran.karimzadehreghbati@mail.utoronto.ca) and we promise to answer quickly.