

Introduction à ggplot

Denis Puthier

21 August 2017

Présentation de ggplot2

Lorsqu'on fait des statistiques descriptives on souhaite souvent partitionner la fenêtre graphique en fonction des différents niveaux pris par une variable catégorielle (*i.e* qualitative) ou ordinale. Réaliser de tels graphiques se révèle vite assez compliqué avec les librairies de base (*graphics*, *lattice*...). Dans le but de faciliter la réalisation de tels graphiques *Hadley Wickham* a développé la librairie *ggplot2* qui est rapidement devenue populaire dans le monde de la bioinformatique (ici les variables catégorielles peuvent être des gènes, groupe de gènes, chromosomes, voies de signalisation, marque de chromatine...) et les variables ordinales des classes d'expression par exemple. L'une des particularité de la librairie *ggplot2* est que son développement est basé sur un modèle proposé par Leland Wilkison dans son ouvrage "The Grammar of Graphics". Dans ce modèle le graphique est vu comme une entité composé de couches successives (*layers*), d'échelles (*scales*), d'un système de coordonnées et de facettes. Il faut donc créer un graphique et venir ajouter les différents éléments à l'aide de l'opérateur '+'. Le principe est assez déconcertant pour les utilisateurs des librairies basiques de R. Cependant, avec le temps on mesure l'intérêt de cette solution car elle nécessite moins de manipulation pour réaliser des graphiques complexes.

Réaliser un graphique basique

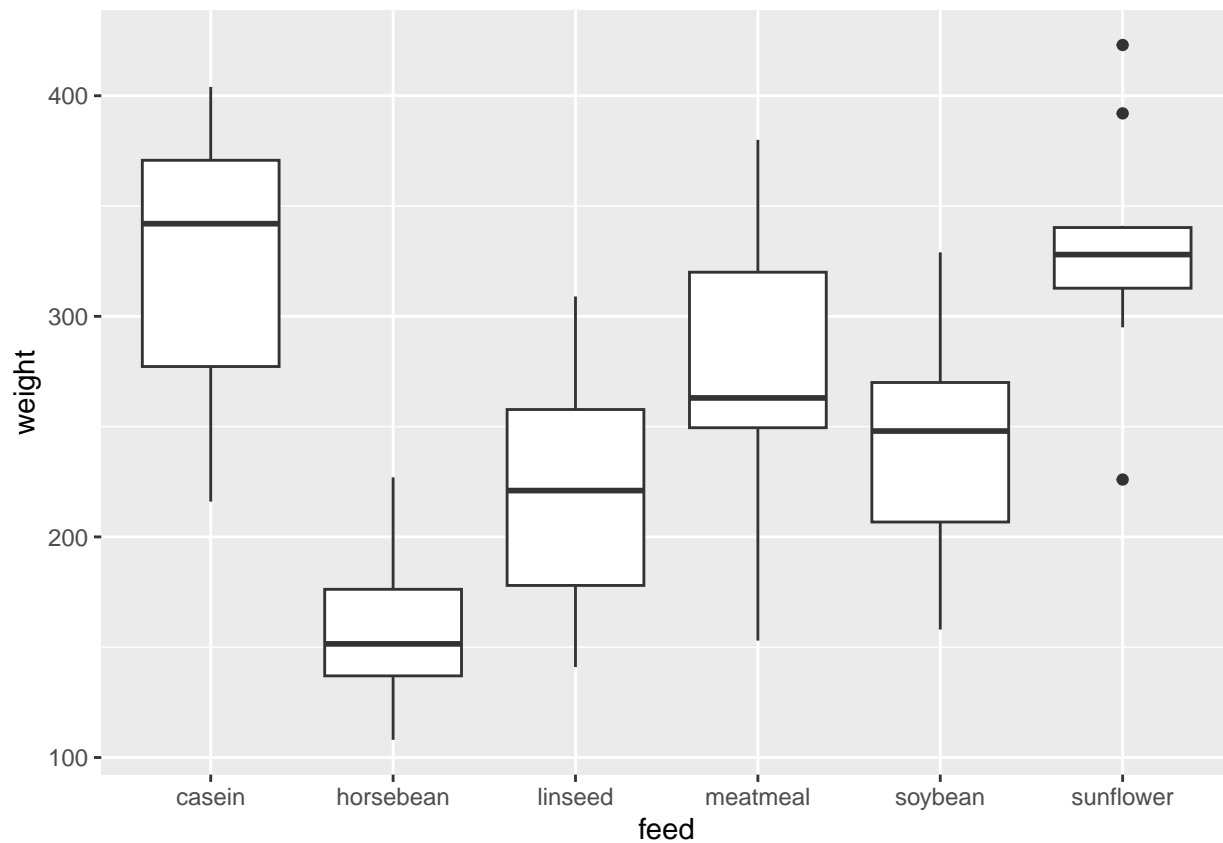
boxplot et violin plot

Les boîtes à moustaches (*boxplots*) et diagramme en violon (violin plots) peuvent être utilisés pour représenter les distributions associées à un jeu de données. On donne ci-dessous quelques exemples.

```
## loading ggplot2 package
library(ggplot2)
## Then we can load a demonstration dataset
data(chickwts)
View(chickwts)

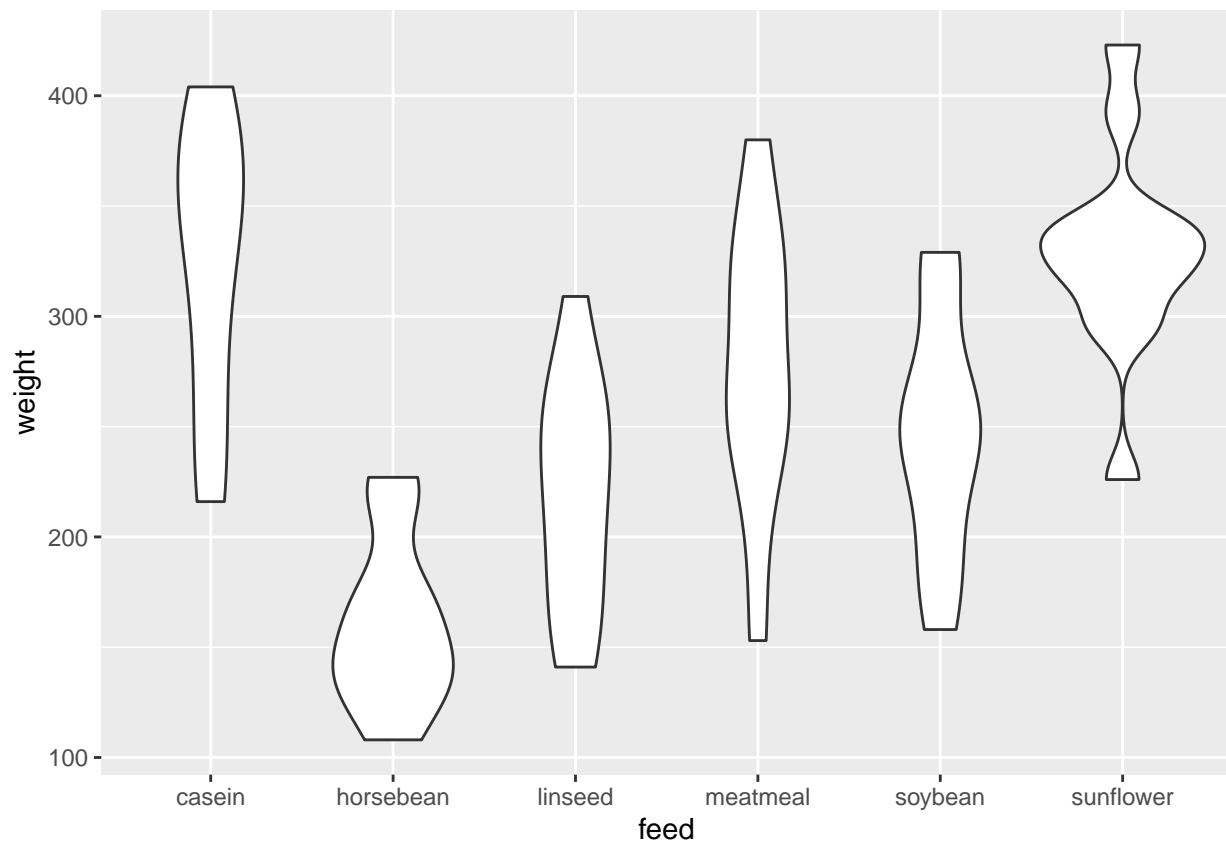
## Then we declare a new graphics and associate
## a dataset. Here the aes (aesthetic) argument is set
## to feed and len that correspond to Insectfeeds dataset column names and
## will be the x and y axes respectively.
p <- ggplot(data=chickwts, aes( x=feed, y=weight))

## We have to indicate the type of requested graphics
p.bp <- p + geom_boxplot()
print(p.bp)
```



We can also easily produce a violin plot using the following instructions

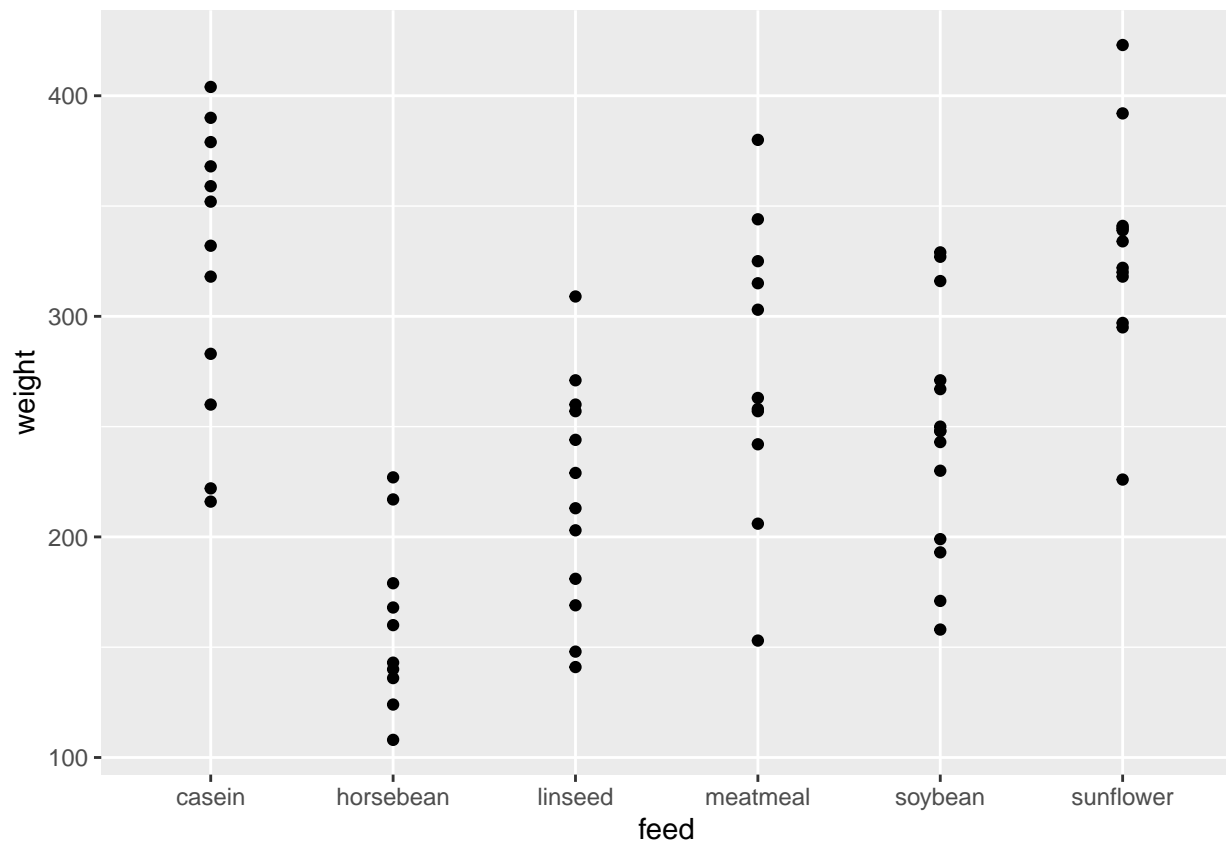
```
p.vp <- p + geom_violin()
print(p.vp)
```



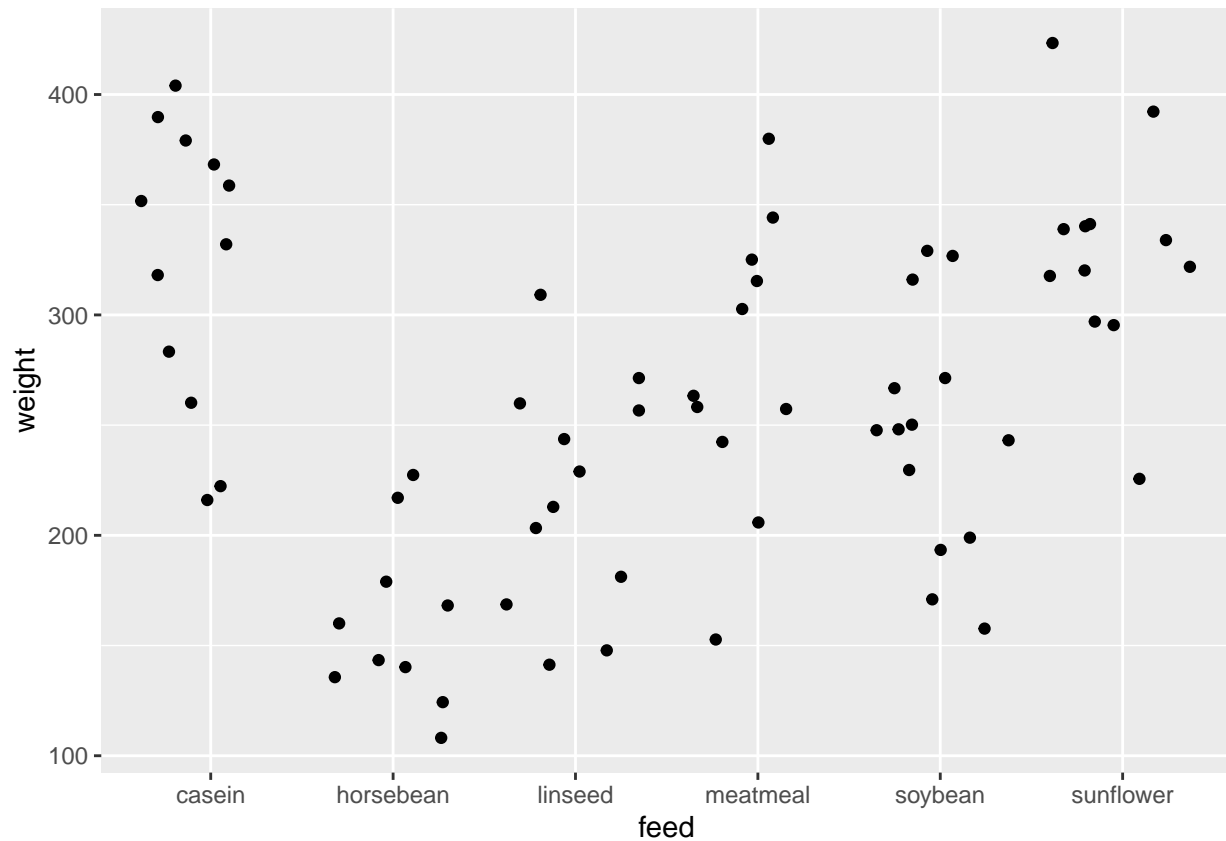
Nuages de points

Il y a environ 40 types de graphiques disponibles. Quelques exemples présentent ci-dessous les instructions pour réaliser des nuages de points.

```
## We can for instance show the values associated to each feed  
p.pt <- p + geom_point()  
print(p.pt)
```



```
## However as there are some ties it may be advised to
## use the jitter option that will add some randomness to the value of the x axis (that here are categories)
p.jt <- p + geom_jitter()
print(p.jt)
```

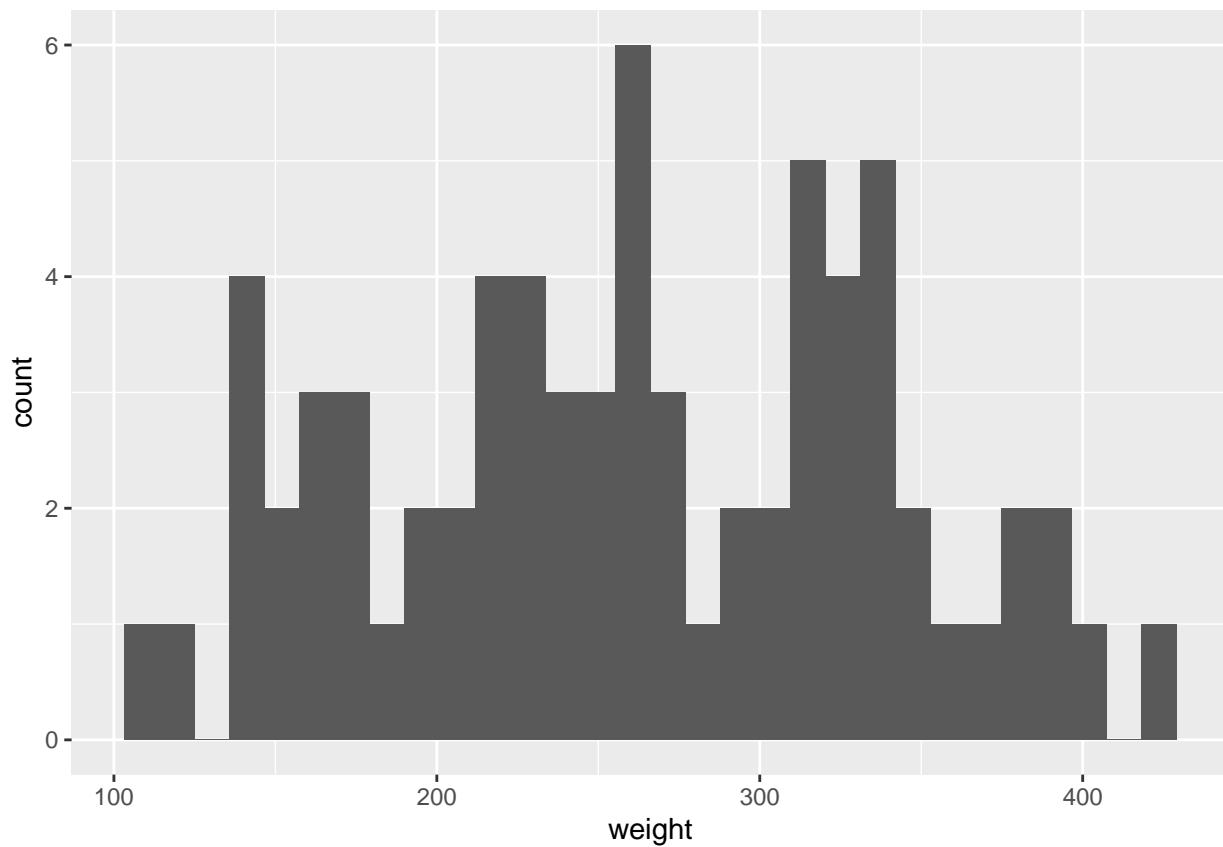


Histogrammes et densités

Dans le cas de l'histogramme, l'axe des x correspond à des interval (*bins*) et l'axe des y au nombre de fois on les valeurs de comptage sont observés dans ces interval. Il n'y aura donc qu'une seule variable à fournir pour la fonction *aes*.

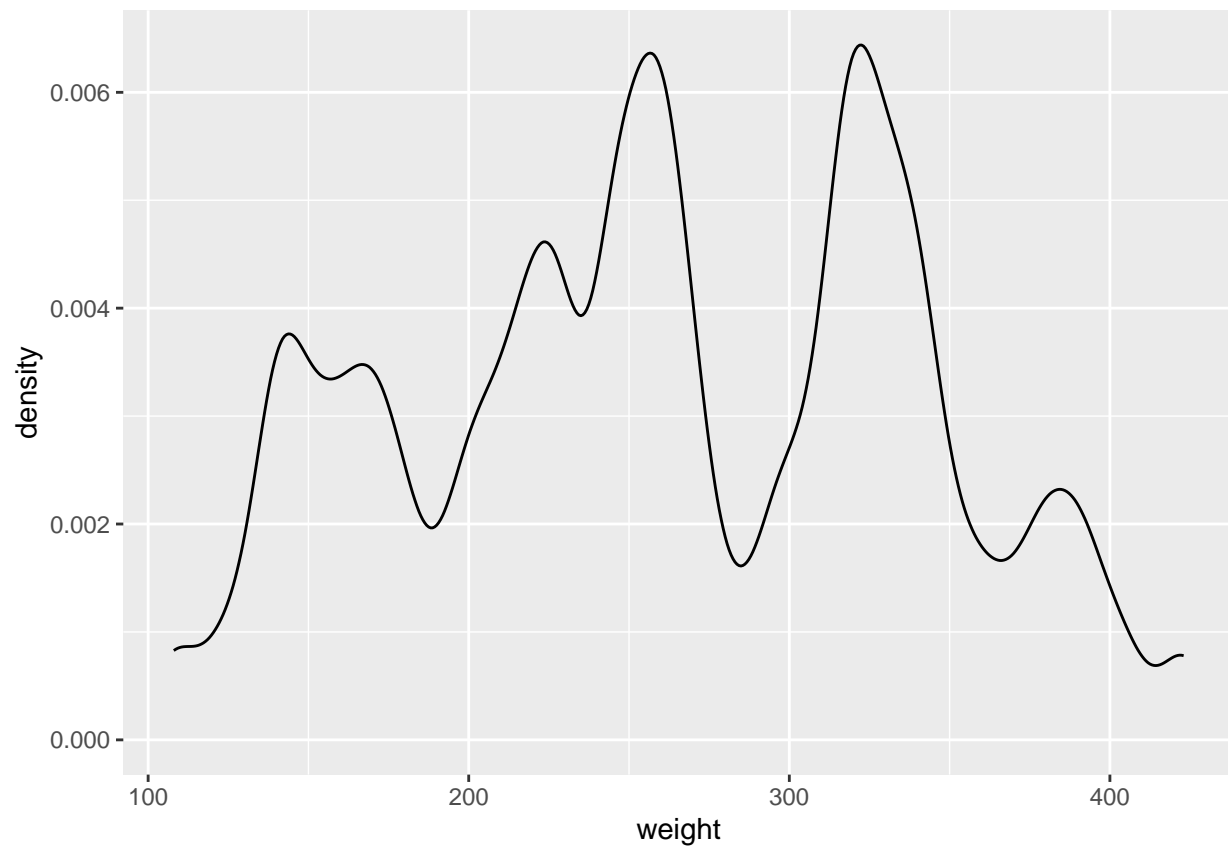
```
## Then we declare a new graphics and associate
ggplot(data=chickwts, aes(x=weight)) + geom_histogram()
```

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



On peut aussi réaliser un profil de densité de probabilité en utilisant la fonction `geom_density()`

```
ggplot(data=chickwts, aes(x=weight)) + geom_density(adjust = 1/4)
```

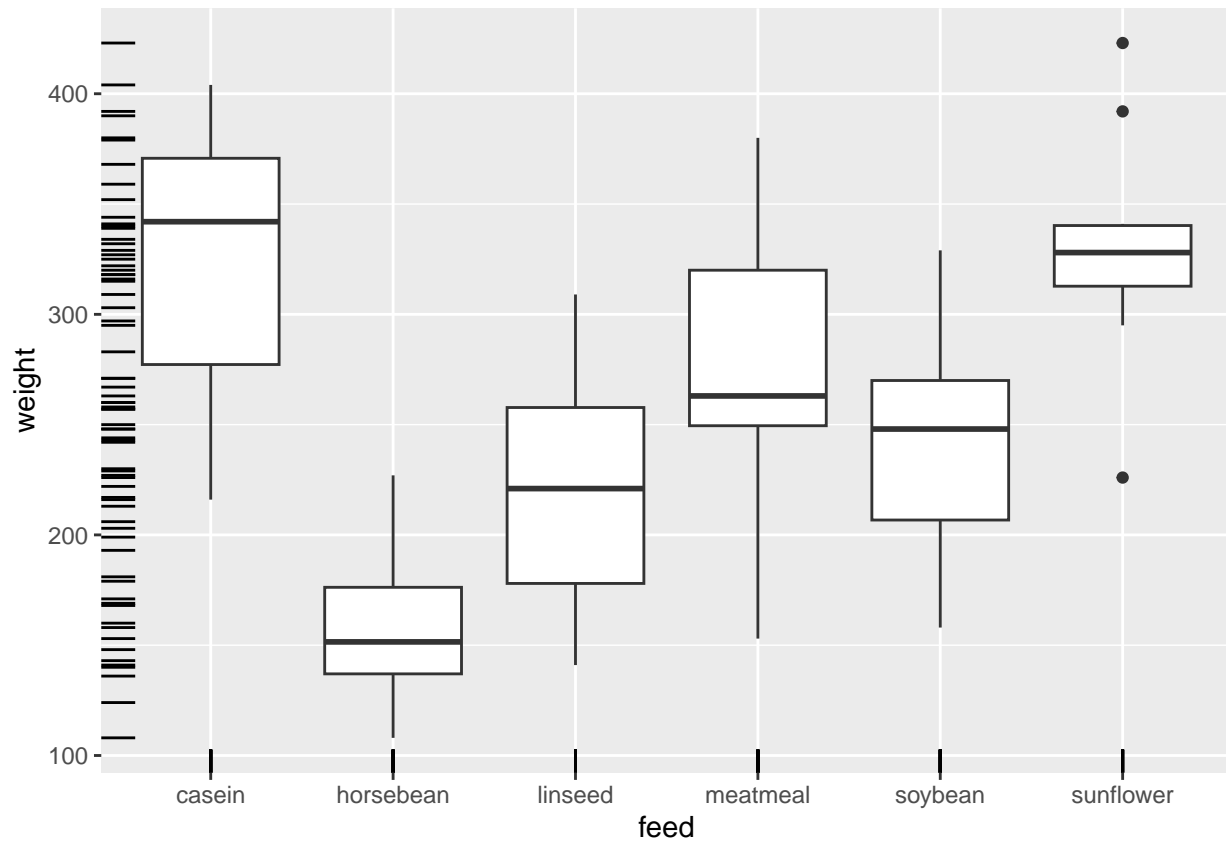


Superposer des éléments graphiques

Exemple autour du boxplot

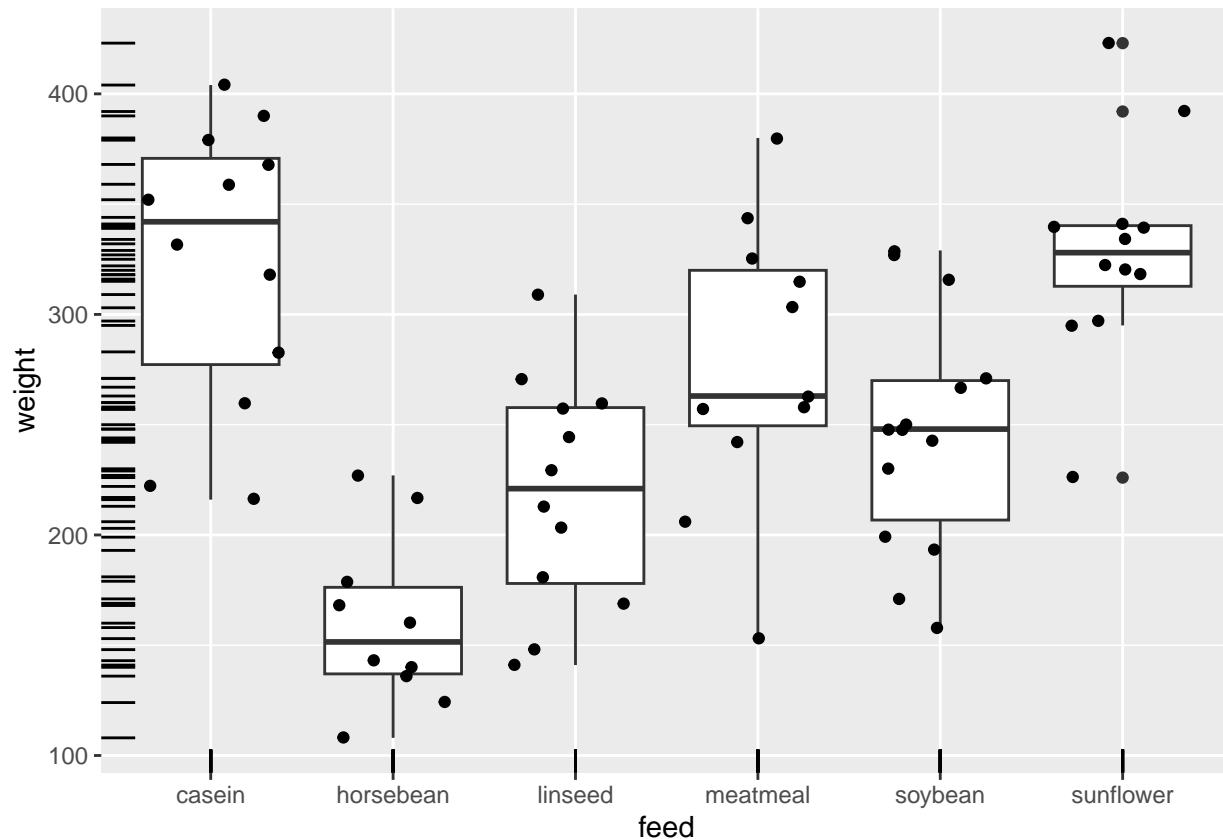
Le modèle sous-jacent à ggplot permet de superposer relativement facilement des couches graphiques.

```
## One can add a rug to the boxplot graphics  
p.bp + geom_rug()
```



One can display the scattered values on the boxplot

```
p.bp + geom_jitter() + geom_rug()
```

->

Facettes

L'utilisation des facettes permet d'explorer les données en fonction d'un facteur ou d'un groupe de facteurs donnés. Pour l'exemple suivant nous allons créer une matrice contenant les résultats d'un test ELISA fictif dans lequel on mesure à 2 temps différents (jours) les expériences réalisés par quatres opérateurs différents.

```
url <- "https://github.com/dputhier/jgb53d-bd-prog_github/blob/gh-pages/data/elisa/elisa_artificial.txt"
elisa <- read.table(url, sep="\t", header=TRUE)
```

```
## Warning in scan(file = file, what = what, sep = sep, quote = quote, dec = dec, :
## EOF within quoted string
```

```
x.1 <- matrix(round(rnorm(96,20,1),0.5),nr=8,ncol=12)
row.names(x.1) <- c('cont',letters[1:7])
colnames(x.1) <- LETTERS[1:12]
x.1[2,] <- x.1[2,] + 4
x.1[3,] <- x.1[3,] + 3
x.1[4,] <- x.1[4,] + 2
x.1[5,] <- x.1[5,] + 1
x.1[6,] <- x.1[6,] + 0.5
x.1[7,] <- x.1[7,] + 0.25
x.1[8,] <- x.1[8,] + 0.125
```

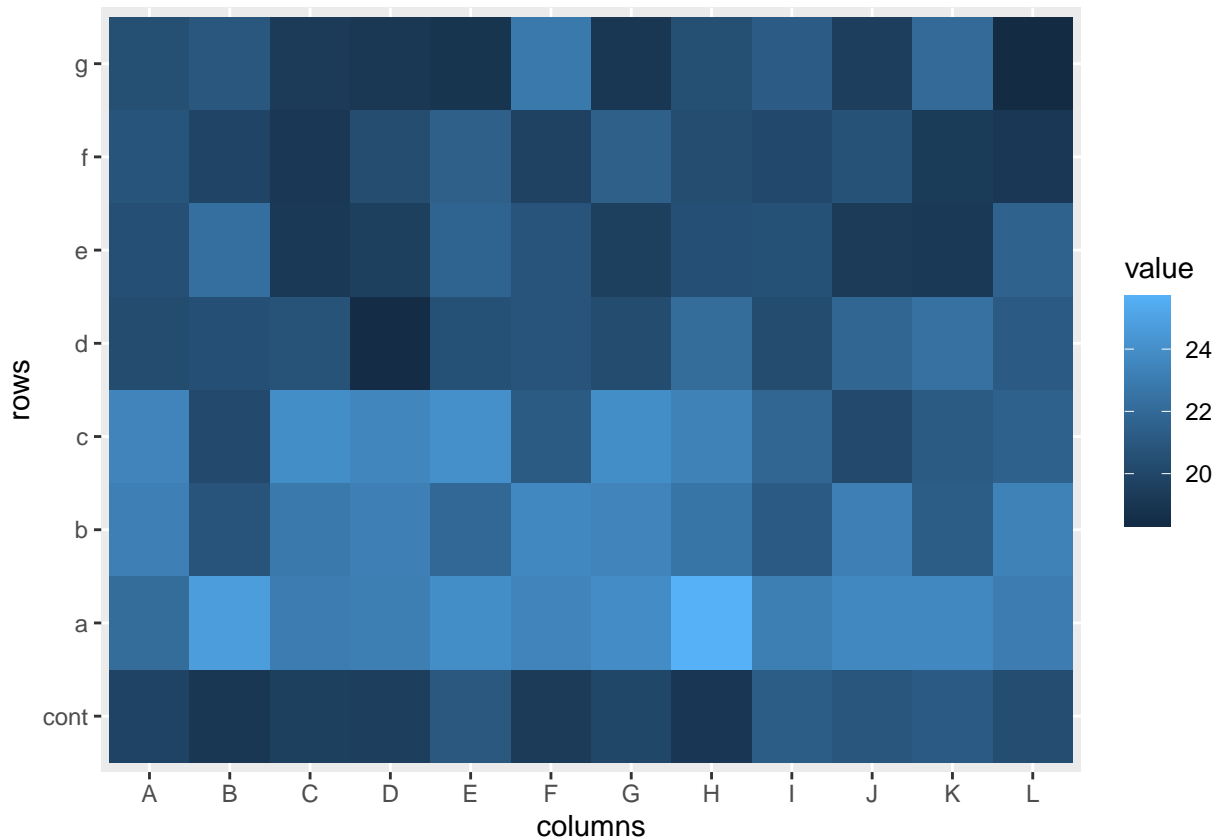
```
library(reshape2)
x.1.long <- melt(x.1, varnames=c("rows","columns"))
```

```

x.1.long <- data.frame(day=c(rep("Mon",96),rep("Fri",96)), rbind(x.1.long,x.1.long))
x.2.long <- x.1.long
x.2.long$value <- x.2.long$value + rnorm(nrow(x.1.long))
x.3.long <- x.1.long
x.3.long$value <- x.3.long$value + rnorm(nrow(x.1.long))
x.4.long <- x.1.long
x.4.long$value <- x.4.long$value + rnorm(nrow(x.1.long))
x.long <- rbind(x.1.long, x.2.long, x.3.long, x.4.long)
x.long <- data.frame(x.long, user=c(rep("Alain", nrow(x.1.long)),
                                   rep("Mathilde", nrow(x.1.long)),
                                   rep("Yvan", nrow(x.1.long)),
                                   rep("Sophie", nrow(x.1.long)))))

library(ggplot2)
p <- ggplot(data=x.long[x.long$user=="Alain" & x.long$day=="Mon",], mapping=aes(x = columns, y = rows, fill = value))
p <- p + geom_raster()
print(p)

```



```

x.1.long$rows <- factor(x.1.long$rows, levels = rev(levels(x.1.long$rows)))
p <- ggplot(data=x.1.long, mapping=aes(x = columns, y = rows, fill = value))
p <- p + geom_raster()
print(p)

```

