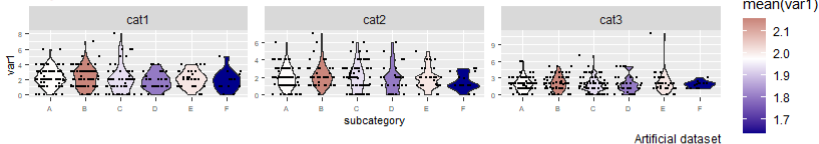


Objective: multiple violin plots

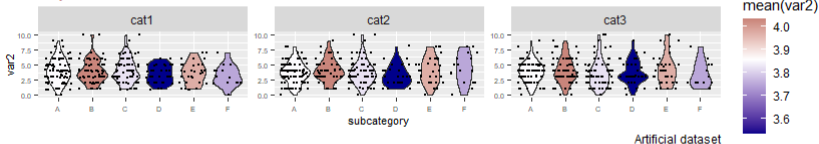
Violin plot for var1 x subcategory, for each category group

Color gradient indicate the mean of the variable "var1"



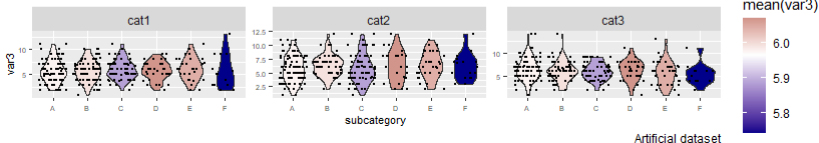
Violin plot for var2 x subcategory, for each category group

Color gradient indicate the mean of the variable "var2"



Violin plot for var3 x subcategory, for each category group

Color gradient indicate the mean of the variable "var3"



Creating fake numerical and categorical data including missing values

```
1 # load libraries and create artificial dataset
2 library(caret)
3 library(missForest)
4 library(tidyverse)
5 library(RANN)
6 library(gridExtra)
7
8 set.seed(2023) # for reproducibility
9 var1 <- rpois(900,2)
10 var2 <- rpois(900,4)
11 var3 <- rpois(900,6)
12 category <- c(rep('cat1', 300), rep('cat2', 300), rep('cat3', 300))
13 subcategory <- sample(LETTERS[1:6], size = 900, replace = TRUE,
14                      prob = c(0.25, 0.3, 0.2, 0.1, 0.1, 0.05))
15 dataset <- data.frame(var1, var2, var3, subcategory, category)
16
17 # prodNA produce 5% missing data
18 dataset.mis = data.frame(prodNA(dataset[,1:3], noNA = 0.05), subcategory,
19                          category)
20
21 # save a copy of the dataset in .csv
22 write.csv(dataset.mis, "path/dataset.mis.csv", row.names = FALSE)
23
24 dataset.mis = read.csv("path/dataset.mis.csv", header = TRUE)
25
26 dataset.mis[298:301, ] # excerpt of the dataset contains NA values
27 #   var1 var2 var3 subcategory category
28 # 298    0   10    5          C    cat1
29 # 299   NA    1    5          E    cat1
30 # 300    2    3    5          A    cat1
31 # 301    2   NA    8          B    cat2
```

Missing value imputation using 'knn'

```
1 # knn imputation using caret and 5 'neighbours'
2 set.seed(2023)
3 dataset.mis.model = preProcess(dataset.mis %>%
4                               dplyr::select(names(dataset.mis)),
5                               "knnImpute", k = 5, knnSummary = mean)
6 dataset.mis.model
7 dataset.mis.pred = predict(dataset.mis.model, dataset.mis) # variables are
  normalized
8 dataset.mis.pred[298:301, ]
9 #           var1          var2          var3 subcategory category
10 # 298 -1.37892580  3.0235674 -0.4137490          C        cat1
11 # 299  0.70767513 -1.4093204 -0.4137490          E        cat1
12 # 300  0.01214149 -0.4242342 -0.4137490          A        cat1
13 # 301  0.01214149 -0.1287083  0.8754844          B        cat2
14
15 # values in original scale
16 complete.dataset <- data.frame(col = names(dataset.mis[,1:3]),
17                                mean = dataset.mis.model$mean,
18                                sd = dataset.mis.model$std)
19 for(i in complete.dataset$col){
20   dataset.mis.pred[i] <- dataset.mis.pred[i]*dataset.mis.model$std[i] +
21     dataset.mis.model$mean[i] }
22
23 # now the dataset is complete
24 complete.data..dataset <- dataset.mis.pred
25 complete.data.dataset[298:301, ]
26 #           var1 var2 var3 subcategory category
27 # 298         0 10.0   5          C        cat1
28 # 299         3  1.0   5          E        cat1
29 # 300         2  3.0   5          A        cat1
30 # 301         2  3.6   8          B        cat2
```

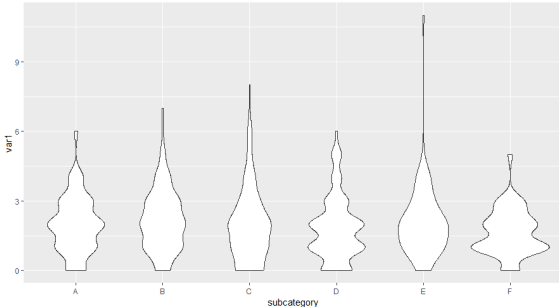
Rearranging and computing means using tidyverse functions

```
1 dataset <- complete.data.dataset
2
3 # compute the mean of the variable 'var1' for each 'subcategory' group
4 dataset2 <- dataset %>%
5   group_by(subcategory) %>%
6   mutate(Mean_var1 = mean(var1))
7
8 # compute the mean of the variable 'var2' for each 'subcategory' group
9 dataset3 <- dataset %>%
10  group_by(subcategory) %>%
11  mutate(Mean_var2 = mean(var2))
12
13 # compute the mean of the variable 'var3' for each 'subcategory' group
14 dataset4 <- dataset %>%
15   group_by(subcategory) %>%
16   mutate(Mean_var3 = mean(var3))
17
18 dataset4[298:301, ]
19 # A tibble: 4      6
20 # Groups:   subcategory [2]
21 #   var1 var2 var3 subcategory category Mean_var3
22 #   <dbl> <dbl> <dbl> <chr>      <chr>      <dbl>
23 #   1     4     3     3 C        cat1         5.60
24 #   2     3     3     6 D        cat1         6.08
25 #   3     4     1     3 C        cat1         5.60
26 #   4     0     3     5 D        cat2         6.08
```

Creating a violin plot with minimal code

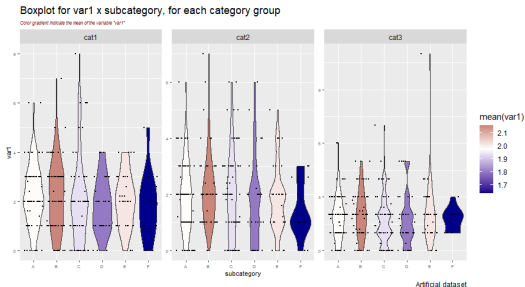
Two lines of code and with help of tidyverse including the powerful visualization library ggplot2 yield a 'quick and dirty' violin plot on which we will improve. We can have an idea of the distribution of the variable 1 for each subcategory.

```
1 ggplot(dataset, aes(x = subcategory, y = var1)) +  
2   geom_violin()
```



Multiple plots with mean color gradient

```
1 ggplot(dataset2, aes(x = subcategory, y = var1)) +  
2   geom_violin(aes(fill=Mean_var1)) +  
3   geom_point(aes(x = subcategory, y = var1), position = 'jitter', size = 0.4) +  
4   scale_fill_gradient2('mean(var1)', low = "blue4",  
5                         mid = "white", high = "firebrick4",  
6                         midpoint = mean(dataset2$Mean_var1)) +  
7   facet_wrap(~category, scales="free") +  
8   labs(title = 'Boxplot for var1 x subcategory, for each category group',  
9        subtitle = "Color gradient indicate the mean of the variable 'var1'",  
10       caption = "Artificial dataset") +  
11   theme(axis.text=element_text(size=5),  
12         axis.title=element_text(size=8),  
13         plot.subtitle=element_text(size=6, face="italic", color="darkred"))
```



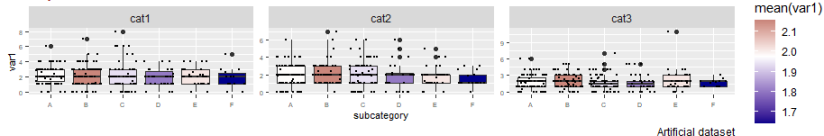
R code to reproduce the figure of slide 1

```
1 # Violin plot with facet_wrap
2 p1 <- ggplot(dataset2, aes(x = subcategory, y = var1)) +
3   geom_violin(aes(fill=Mean_var1)) +
4   geom_point(aes(x = subcategory, y = var1), position = 'jitter', size = 0.4) +
5   scale_fill_gradient2('mean(var1)', low = "blue4",
6     mid = "white", high = "firebrick4",
7     midpoint = mean(dataset2$Mean_var1)) +
8   facet_wrap(~category, scales="free") +
9   labs(title = 'Violin plot for var1 x subcategory, for each category group',
10     subtitle = "Color gradient indicate the mean of the variable 'var1'",
11     caption = "Artificial dataset") +
12   theme(axis.text=element_text(size=5),
13     axis.title=element_text(size=8),
14     plot.subtitle=element_text(size=6, face="italic", color="darkred"))
15
16 # Violin plot with facet_wrap
17 p2 <- ggplot(dataset3, aes(x = subcategory, y = var2)) +
18   geom_violin(aes(fill= Mean_var2)) +
19   ...
20
21 # Violin plot with facet_wrap
22 p3 <- ggplot(dataset4, aes(x = subcategory, y = var3)) +
23   geom_violin(aes(fill= Mean_var3)) +
24   ...
25
26 final.plot <- grid.arrange(p1, p2, p3)
```

Objective: multiple boxplots

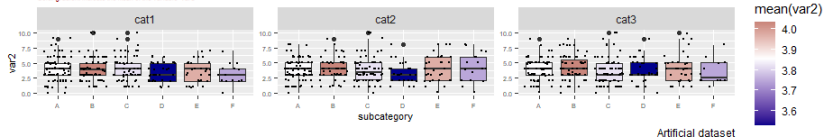
Boxplot for var1 x subcategory, for each category group

Color gradient indicate the mean of the variable "var1"



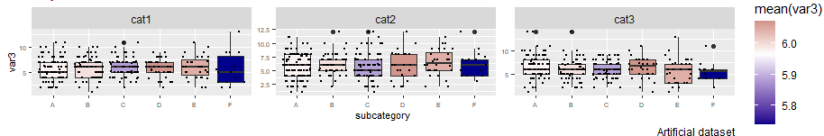
Boxplot for var2 x subcategory, for each category group

Color gradient indicate the mean of the variable "var2"



Boxplot for var3 x subcategory, for each category group

Color gradient indicate the mean of the variable "var3"



What is a boxplot?

A Boxplot is a non-parametric descriptive statistical way to summarize and visualize the distribution of grouped data. It is used when we have a quantitative variable and a qualitative variable (nominal or ordinal). It gives an indication of the Median (measure of central tendency), the Interquartile Range (measure of dispersion), the 95% range for the observations, as well as outlying observations (observations outside the 95% range). We give the definition of those statistics, F^{-1} being the inverse CDF or Quantile function.

$$\text{First quartile (Q1)} = F^{-1}(0.25)$$

$$\text{Median (Q2)} = F^{-1}(0.5)$$

$$\text{Third quartile (Q3)} = F^{-1}(0.75)$$

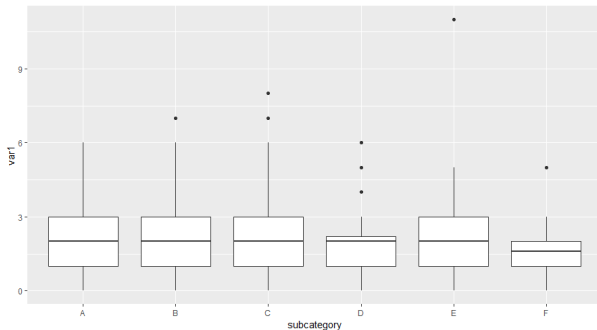
$$\text{Interquartile Range (IQR)} = Q3 - Q1$$

$$\text{Outlying observations : } x_i > F^{-1}(0.95) \text{ or } x_i < F^{-1}(0.05)$$

Create a boxplot with minimal code

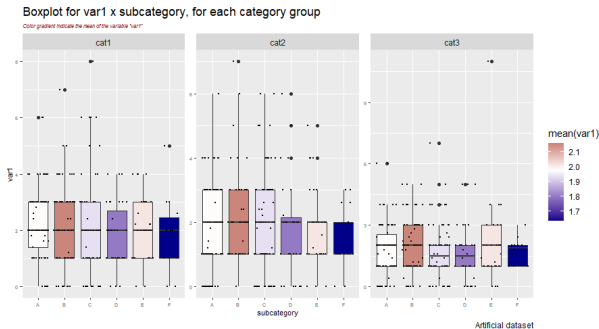
Two lines of code and with help of tidyverse including the powerful visualization library ggplot2 yield a 'quick and dirty' boxplot on which we will improve.

```
1 ggplot(dataset, aes(x = subcategory, y = var1)) +  
2   geom_boxplot()
```



Multiple plots with mean color gradient

```
1 ggplot(dataset2, aes(x = subcategory, y = var1)) +  
2   geom_boxplot(aes(fill=Mean_var1)) +  
3   geom_point(aes(x = subcategory, y = var1), position = 'jitter', size = 0.4) +  
4   scale_fill_gradient2('mean(var1)', low = "blue4",  
5                         mid = "white", high = "firebrick4",  
6                         midpoint = mean(dataset2$Mean_var1)) +  
7   facet_wrap(~category, scales="free") +  
8   labs(title = 'Boxplot for var1 x subcategory, for each category group',  
9        subtitle = "Color gradient indicate the mean of the variable 'var1'",  
10       caption = "Artificial dataset") +  
11   theme(axis.text=element_text(size=5),  
12         axis.title=element_text(size=8),  
13         plot.subtitle=element_text(size=6, face="italic", color="darkred"))
```



R code to reproduce the figure of slide 8

```
1 # to access the function grid.arrange() for multiple plotting
2 library(gridExtra)
3
4 # Boxplot with facet_wrap
5 p1 <- ggplot(dataset2, aes(x = subcategory, y = var1)) +
6   geom_boxplot(aes(fill=Mean_var1)) +
7   geom_point(aes(x = subcategory, y = var1), position = 'jitter', size = 0.4) +
8   scale_fill_gradient2('mean(var1)', low = "blue4",
9                        mid = "white", high = "firebrick4",
10                       midpoint = mean(dataset2$Mean_var1)) +
11   facet_wrap(~category, scales="free") +
12   labs(title = 'Boxplot for var1 x subcategory, for each category group',
13        subtitle = "Color gradient indicate the mean of the variable 'var1'",
14        caption = "Artificial dataset") +
15   theme(axis.text=element_text(size=5),
16         axis.title=element_text(size=8),
17         plot.subtitle=element_text(size=6, face="italic", color="darkred"))
18
19 # Boxplot with facet_wrap
20 p2 <- ggplot(dataset3, aes(x = subcategory, y = var2)) +
21   geom_boxplot(aes(fill= Mean_var2)) +
22   ...
23
24 # Boxplot with facet_wrap
25 p3 <- ggplot(dataset4, aes(x = subcategory, y = var3)) +
26   geom_boxplot(aes(fill= Mean_var3)) +
27   ...
28
29 final.plot <- grid.arrange(p1, p2, p3)
```