

Projet Final R

Victor Simonin & Alexandre Lemonnier

Dataset : `decathlon.csv`

Import des bibliothèques

Pour commencer il est nécessaire d'importer les bibliothèques qui nous permettront d'utiliser leurs fonctions par la suite du projet

```
In [89]: library("WIM") # Pour la visualisation des valeurs manquantes
library("corplot") # Pour les matrices de corrélation.
library("factoextra") # Pour les graphes de l'ACP
library("FactoMineR") # Diagramme de dispersion
library("PerformanceAnalytics") # Diagramme de dispersion
library("pander") # Pour la régression logistique
```

Lecture du dataset

```
In [90]: df <- read.table('data/decathlon.csv', header = T, sep = ';', dec = '.', row.names = 1)
head(df)
```

Variable	Number of missings
X100m	0
Poids	0
Hauteur	0
X400m	0
Disque	0
Javelot	0
X1500m	0
Classement	0
Points	41
Competitions	0

Le dataset ci-dessus est composé de plusieurs données assez différentes. Celles-ci représentent les scores obtenus par des athlètes de l'épreuve du decathlon au **Jeux Olympiques** et au **Decastar**, donc les scores des 10 épreuves ainsi que leur classement et leurs points.

Visualisation des données manquantes

```
In [91]: missdata <- aggr(df, prop = F, number = T)
```



On peut voir ici que notre dataset ne semble pas contenir de données manquantes, ce qui facilite notre analyse et nous permet de continuer sereinement. Des choix étaient possibles si jamais le dataset contenait des données manquantes, soit la suppression des lignes qui ne sont pas complètes, ou alors il est possible d'imputer ces données manquantes par plusieurs algorithmes.

Il faudra néanmoins vérifier qu'il n'y a aucune valeurs aberrantes.

Description du dataset

Cette section va nous servir à comprendre un peu mieux notre dataset, notamment sur sa taille, ses attributs et quelques statistiques descriptives sur les colonnes.

En premier lieu, nous allons passer la colonne `Compétition` qui est une colonne catégorique en colonne numérique.

```
In [92]: table(df$Compétition)
compétition <- df$Compétition
```

```
Decastar
13
20
```

On peut voir ici que notre dataset contient les résultats de deux compétitions, les `JO` et le `Decastar`, on va donc simplement remplacer ces deux valeurs par des valeurs numériques en suivant : `JO` = 1 et `Decastar` = 2.

```
In [93]: df[1:28, 13] == 1
df[29:41, 13] == 2
df$Compétition == as.numeric(as.character(df$Compétition))
```

```
In [94]: df
```

A data frame: 41 × 13

	X100m	Longueur	Poids	Hauteur	X400m	X110m.H	Disque	Perche	Javelot	X1500m	Classement	Points	Compétition
<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<int>	<int>	<dbl>
Sebrle	10.85	7.84	16.36	2.12	48.36	14.05	48.72	5.00	70.52	280.01	1	8893	1
Clay	10.44	7.96	15.23	2.06	49.19	14.13	50.11	4.90	69.71	282.00	2	8820	1
Karpov	10.50	7.81	15.93	2.09	46.81	13.97	51.65	4.60	55.54	278.11	3	8725	1
Macey	10.89	7.47	15.73	2.15	48.97	14.56	48.34	4.40	58.46	265.42	4	8414	1
Warners	10.62	7.74	14.48	1.97	47.97	14.01	43.73	4.90	55.39	278.05	5	8343	1
Zisvoczky	10.91	7.14	15.31	2.12	49.40	14.95	45.62	4.70	63.45	269.54	6	8287	1
Hernu	10.97	7.19	14.65	2.03	48.73	14.25	44.72	4.80	57.76	264.35	7	8237	1
Noel	10.80	7.53	14.26	1.88	48.81	14.80	42.05	5.40	61.33	276.33	8	8235	1
Bernard	10.69	7.48	14.80	2.12	49.13	14.17	44.75	4.40	55.27	276.31	9	8225	1
Schwarzl	10.98	7.49	14.01	1.94	49.76	14.25	42.43	5.10	56.32	273.56	10	8102	1
Pogorelov	10.95	7.31	15.10	2.06	50.79	14.21	44.60	5.00	53.45	287.63	11	8064	1
Schoenbeck	10.90	7.30	14.77	1.88	50.30	14.34	44.41	5.00	60.89	278.82	12	8077	1
Barras	11.14	6.99	14.91	1.94	49.41	14.37	44.83	4.60	64.55	267.09	13	8067	1
Smith	10.85	6.81	15.24	1.91	49.27	14.01	49.02	4.20	61.52	272.74	14	8023	1
Averyanov	10.55	7.34	14.44	1.94	49.72	14.39	39.88	4.80	54.51	271.02	15	8021	1
Ojanieni	10.68	7.36	14.97	1.94	49.12	15.01	40.35	4.60	59.26	275.71	16	8006	1
Sminimov	10.89	7.07	13.98	1.94	49.11	14.77	42.47	4.70	60.88	263.31	17	7993	1
Qi	11.06	7.34	13.55	1.97	49.65	14.53	40.13	5.00	60.79	272.63	18	7993	1
Drews	10.87	7.38	13.07	1.88	48.51	14.01	45.11	5.00	51.53	274.21	19	7906	1
Parkhomenko	11.14	6.61	15.69	2.03	51.04	14.88	41.90	4.80	65.82	277.94	20	7918	1
Terek	11.02	6.94	15.15	1.94	49.56	15.12	45.62	5.30	60.62	290.36	21	7893	1
Gomez	10.98	7.26	14.57	1.85	48.61	14.21	40.95	4.40	57.17	269.70	22	7865	1
Turi	11.10	6.91	13.62	2.03	51.67	14.26	39.83	4.80	59.34	290.01	23	7708	1
Lorenzo	11.08	7.03	13.22	1.85	49.34	15.38	40.23	4.50	58.36	263.08	24	7592	1
Karlivans	11.31	7.26	13.30	1.97	50.54	14.98	43.34	4.50	52.92	278.67	25	7583	1
Korkizoglou	10.86	6.97	14.81	1.94	51.16	14.96	46.07	4.70	53.05	317.00	26	7573	1
Uzdal	11.36	6.96	13.53	1.85	50.95	15.09	43.01	4.50	60.00	281.70	27	7495	1
Casarsa	11.23	6.68	14.92	1.94	53.25	15.09	48.66	4.40	58.62	296.12	28	7404	1
SEBRLE	11.04	7.58	14.83	2.07	49.81	14.69	43.75	5.02	63.19	291.70	1	8217	2
CLAY	10.76	7.40	14.26	1.86	49.37	14.05	50.72	4.92	60.15	301.50	2	8122	2
KARPOV	11.02	7.30	14.77	2.04	48.37	14.09	48.95	4.92	59.31	300.20	3	8099	2
BERNARD	11.02	7.23	14.25	1.92	48.93	14.99	40.67	5.32	62.77	280.10	4	8067	2
YURKOV	11.34	7.09	15.19	2.10	49.65	15.31	46.26	4.72	63.44	276.40	5	8036	2
WARNERS	11.11	7.60	14.31	1.98	48.68	14.23	41.10	4.92	51.77	278.10	6	8030	2
ZISVOCZKY	11.13	7.30	13.48	2.01	48.62	14.17	45.67	4.42	56.37	268.00	7	8004	2
McMULLEN	10.83	7.31	13.76	2.13	49.91	14.38	44.41	4.92	56.37	285.10	8	7995	2
MARTINEAU	11.64	6.81	14.57	1.95	50.14	14.93	47.60	4.92	52.33	262.10	9	7802	2
HERNU	11.37	7.56	14.41	1.86	51.10	15.06	44.99	4.82	57.19	285.10	10	7733	2
BARRAS	11.33	6.97	14.09	1.95	49.48	14.48	42.10	4.72	55.40	282.00	11	7708	2
NOEL	11.33	7.27	12.68	1.98	49.20	15.29	37.92	4.62	57.44	266.60	12	7651	2
BOURGIGNON	11.36	6.80	13.46	1.86	51.16	15.69	40.49	5.02	54.68	291.70	13	7313	2

Nos compétitions ont bien été modifiées, ainsi que le type de la colonne, on a donc maintenant uniquement des colonnes numériques dans notre dataset.

```
In [95]: dim(df)
```

```
41 13
```

On comprend donc ici que notre dataset se compose de 41 individus qui ont chacun 13 données dans les colonnes qui servent à les qualifier et à les décrire.

```
In [96]: attributes(df)
```

\$names

"X100m" "Longueur" "Poids" "Hauteur" "X400m" "X110m.H" "Disque" "Perche" "Javelot" "X1500m" "Classement" "Points" "Compétition"

\$row.names

"Sebrle" "Clay" "Karpov" "Macey" "Warners" "Zisvoczky" "Hernu" "Noel" "Bernard" "Schwarzl" "Pogorelov" "Schoenbeck" "Barras" "Smith" "Averyanov" "Ojanieni" "Sminimov" "Qi" "Drews" "Parkhomenko" "Terek" "Gomez" "Turi" "Lorenzo" "Karlivans" "Korkizoglou" "Uzdal" "Casarsa" "SEBRLE" "CLAY" "KARPOV" "BERNARD" "YURKOV" "WARNERS" "ZISVOCZKY" "McMULLEN" "MARTINEAU" "HERNU" "BARRAS" "NOEL" "BOURGIGNON"

\$class

"data.frame"

On retrouve donc les différents attributs de notre dataset avec les noms des différentes colonnes qui représentent les épreuves des **Jeux Olympiques** et du **Decastar**. Les lignes représentent les individus et donc les scores des athlètes participant à ces compétitions.

```
In [97]: summary(df)
```

X100m

Min.: 10.44 Min.: 6.61 Min.: 12.68 Min.: 11.850 Min.: 46.81

1st Qu.: 10.85 1st Qu.: 7.03 1st Qu.: 13.88 1st Qu.: 11.920 1st Qu.: 48.93

Median: 10.90 Median: 7.30 Median: 14.57 Median: 11.950 Median: 49.40

Mean: 11.10 Mean: 7.26 Mean: 14.48 Mean: 11.977 Mean: 49.62

3rd Qu.: 11.14 3rd Qu.: 7.48 3rd Qu.: 14.97 3rd Qu.: 12.040 3rd Qu.: 50.30

Max.: 13.64 Max.: 7.96 Max.: 16.36 Max.: 15.690 Max.: 53.20

X110m.H

Min.: 13.97 Min.: 37.92 Min.: 42.00 Min.: 45.800 Min.: 50.31

1st Qu.: 14.21 1st Qu.: 41.90 1st Qu.: 41.500 1st Qu.: 45.27

Median: 14.48 Median: 44.33 Median: 44.800 Median: 45.36

Mean: 14.61 Mean: 44.33 Mean: 44.762 Mean: 45.32

3rd Qu.: 14.68 3rd Qu.: 46.07 3rd Qu.: 45.308 3rd Qu.: 46.89

Max.: 15.67 Max.: 51.65 Max.: 54.600 Max.: 70.52

X1500m

Min.: 262.1 Min.: 1.00 Min.: 7313 Min.: 11.000

1st Qu.: 271.0 1st Qu.: 6.00 1st Qu.: 7802 1st Qu.: 1.000

Median: 278.1 Median: 11.00 Median: 8021 Median: 1.000

Mean: 279.0 Mean: 12.12 Mean: 8085 Mean: 1.317

3rd Qu.: 285.1 3rd Qu.: 18.00 3rd Qu.: 8122 3rd Qu.: 2.000

Max.: 317.0 Max.: 28.00 Max.: 8893 Max.: 2.000

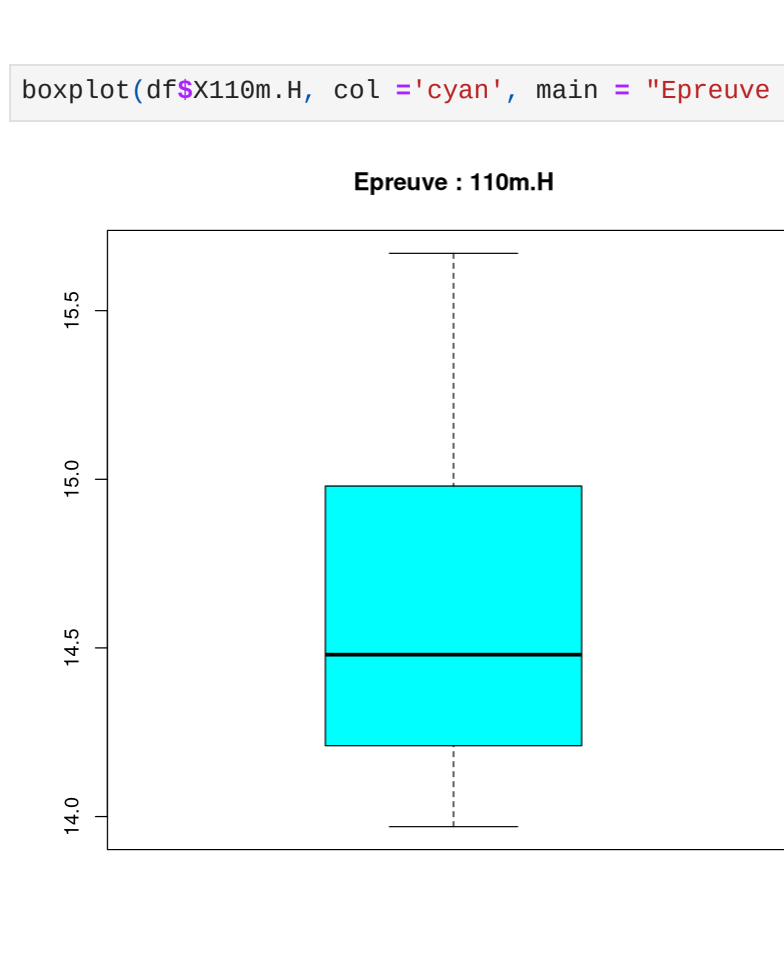
La commande `summary` nous permet d'afficher quelques informations statistiques intéressantes sur nos données.

On retrouve les calculs suivant :

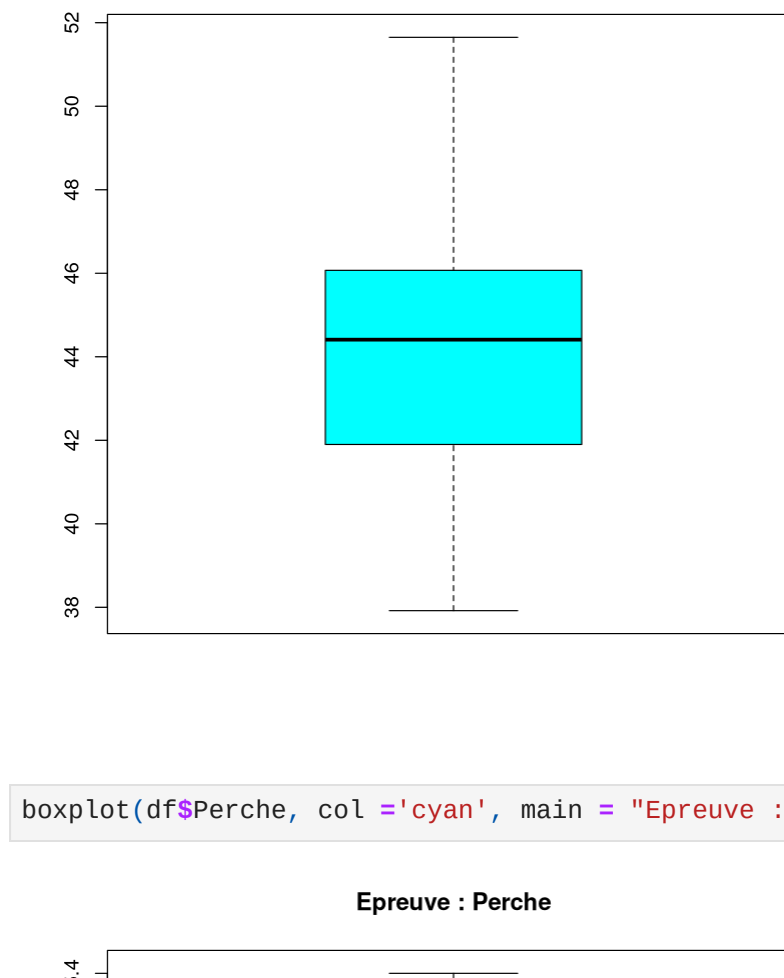
- Minimum
- 1er quartile
- Médiane
- Moyenne
- 3ème quartile
- Maximum

Cela nous permet de comprendre un peu mieux notre dataset, et aussi de commencer à repérer certains outliers et valeurs aberrantes. Ici les colonnes avec les valeurs les plus dispersées sont le **javelot** et le **1500m**. Les statistiques sur le classement et la compétition ne sont pas intéressantes.

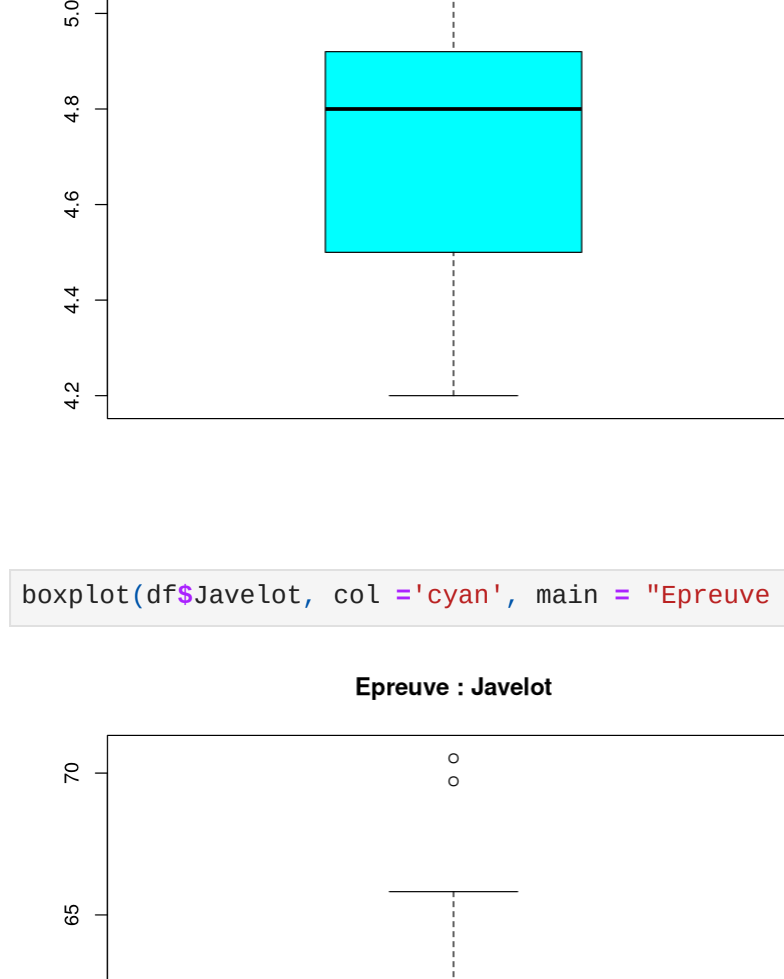
```
In [98]: boxplot(df$X1500m, col = "cyan", main = "Epreuve : 1500m")
```



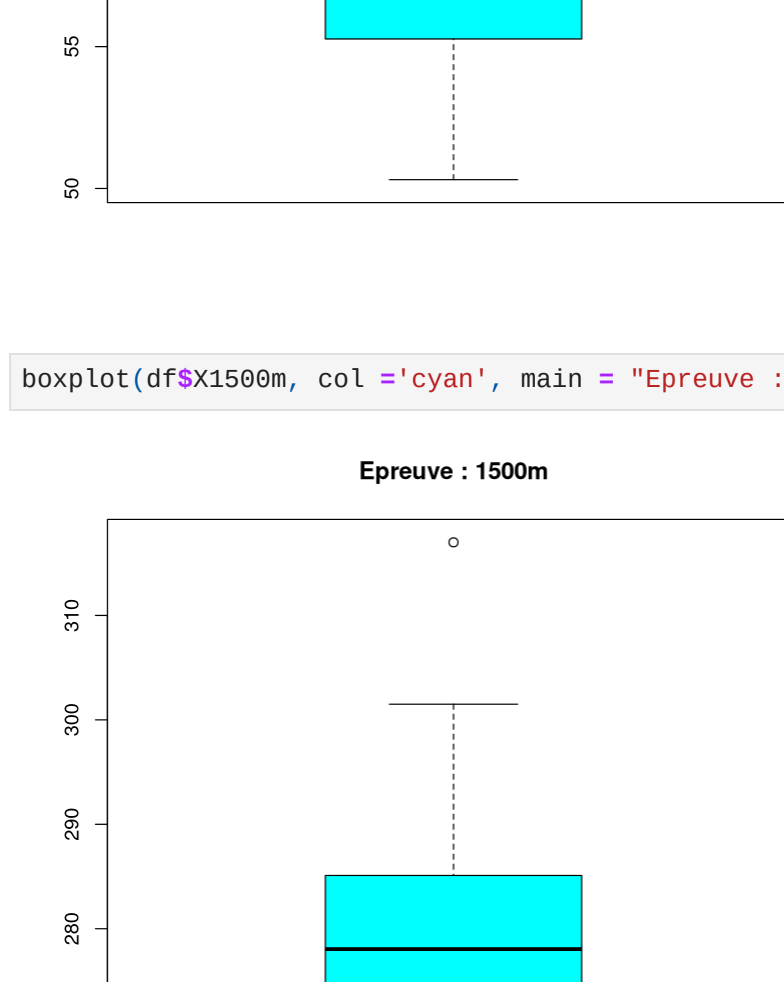
```
In [99]: boxplot(df$Longueur, col = "cyan", main = "Epreuve : Longueur")
```



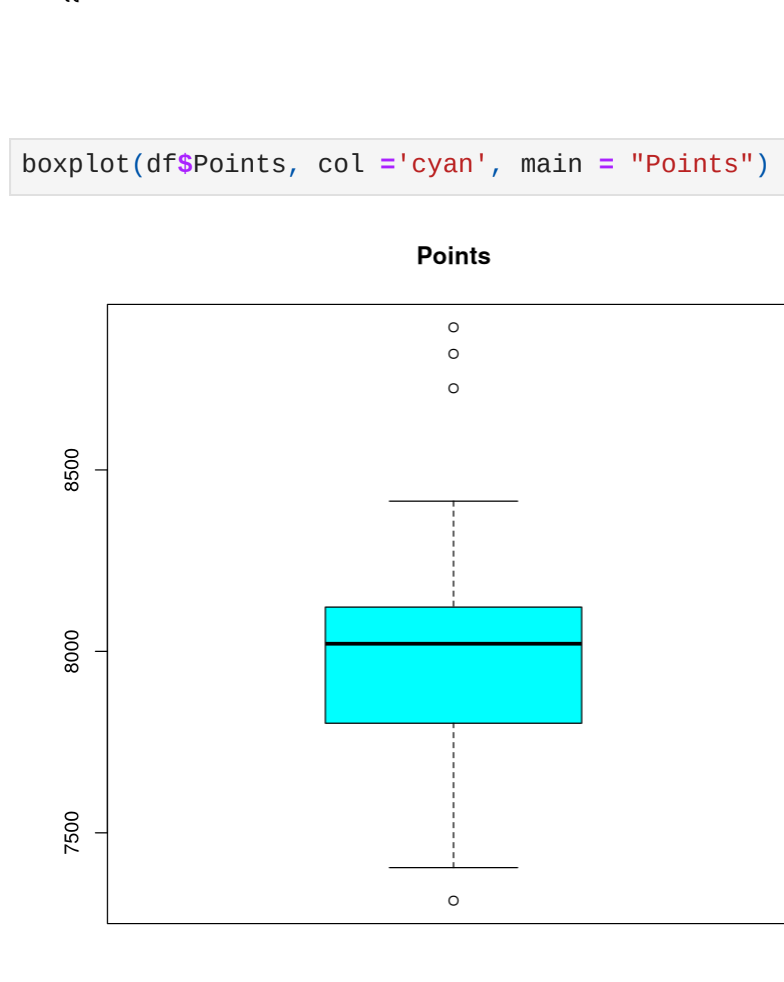
```
In [100]: boxplot(df$Poids, col = "cyan", main = "Epreuve : Poids")
```



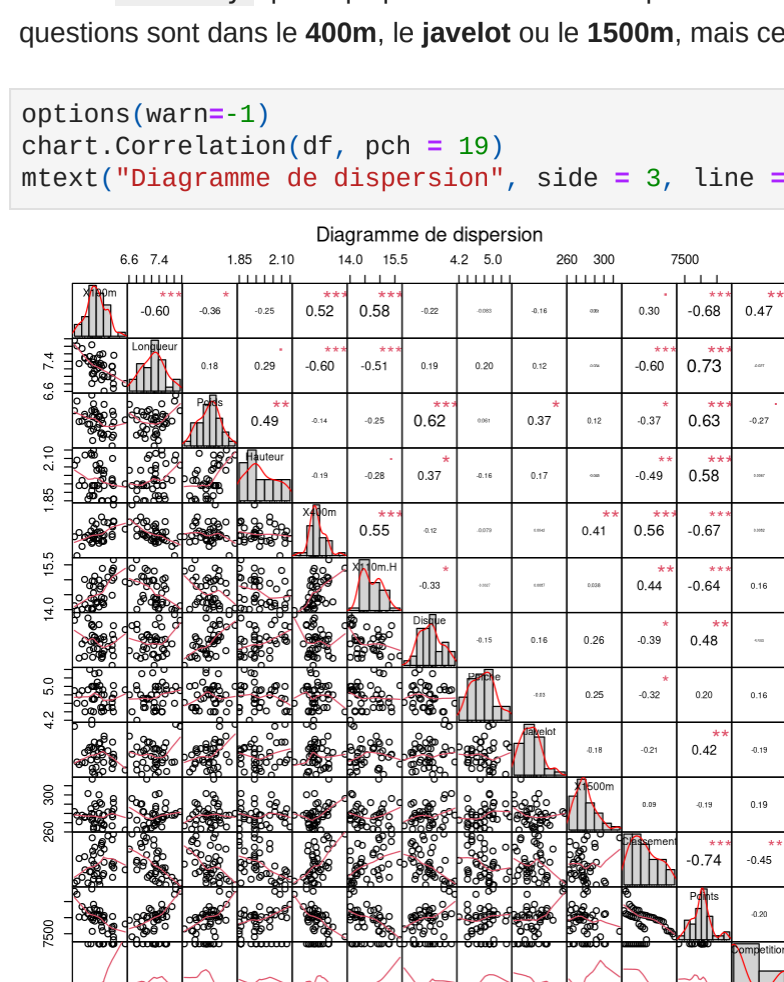
```
In [101]: boxplot(df$X400m, col = "cyan", main = "Epreuve : 400m")
```



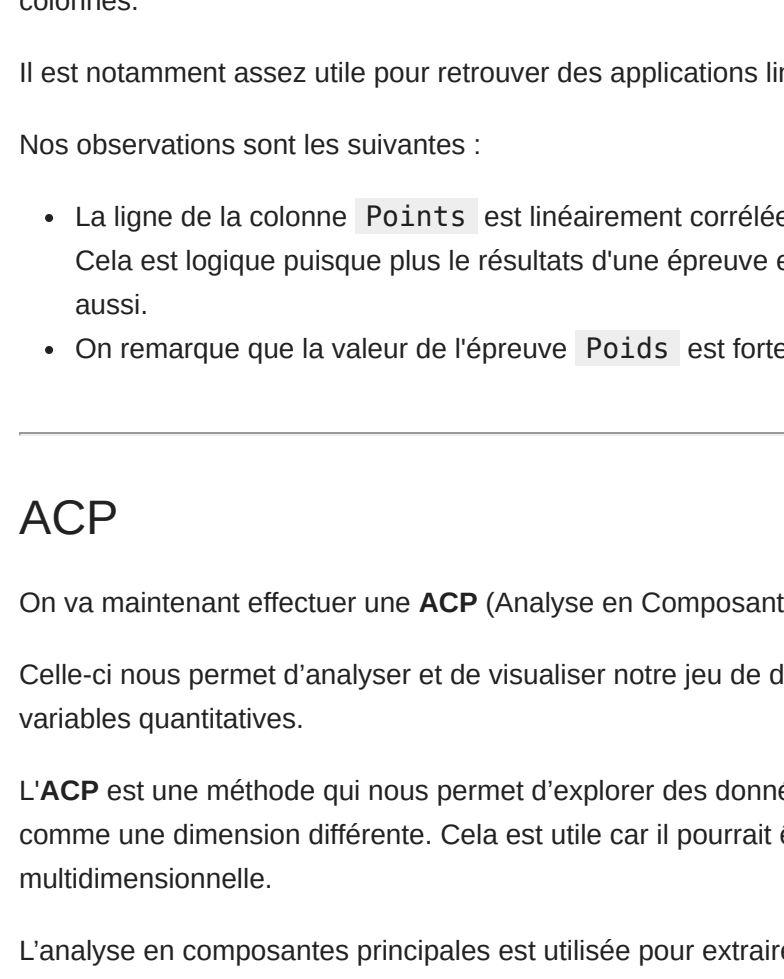
```
In [102]: boxplot(df$X110m.H, col = "cyan", main = "Epreuve : 110m.H")
```



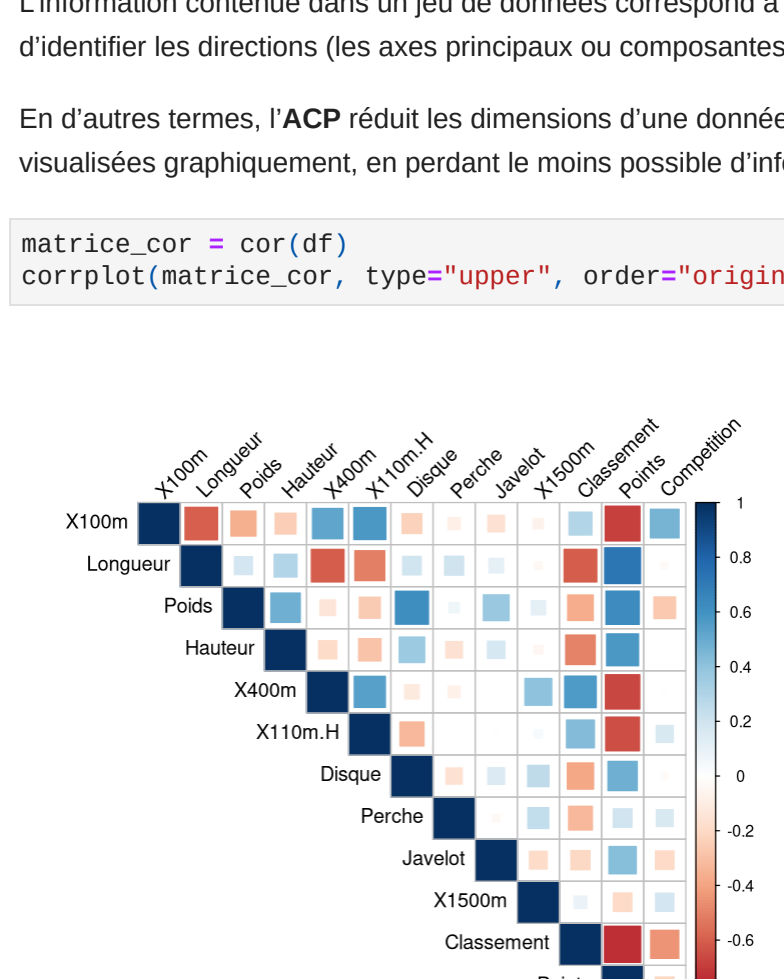
```
In [103]: boxplot(df$Disque, col = "cyan", main = "Epreuve : Disque")
```



```
In [104]: boxplot(df$Perche, col = "cyan", main = "Epreuve : Perche")
```



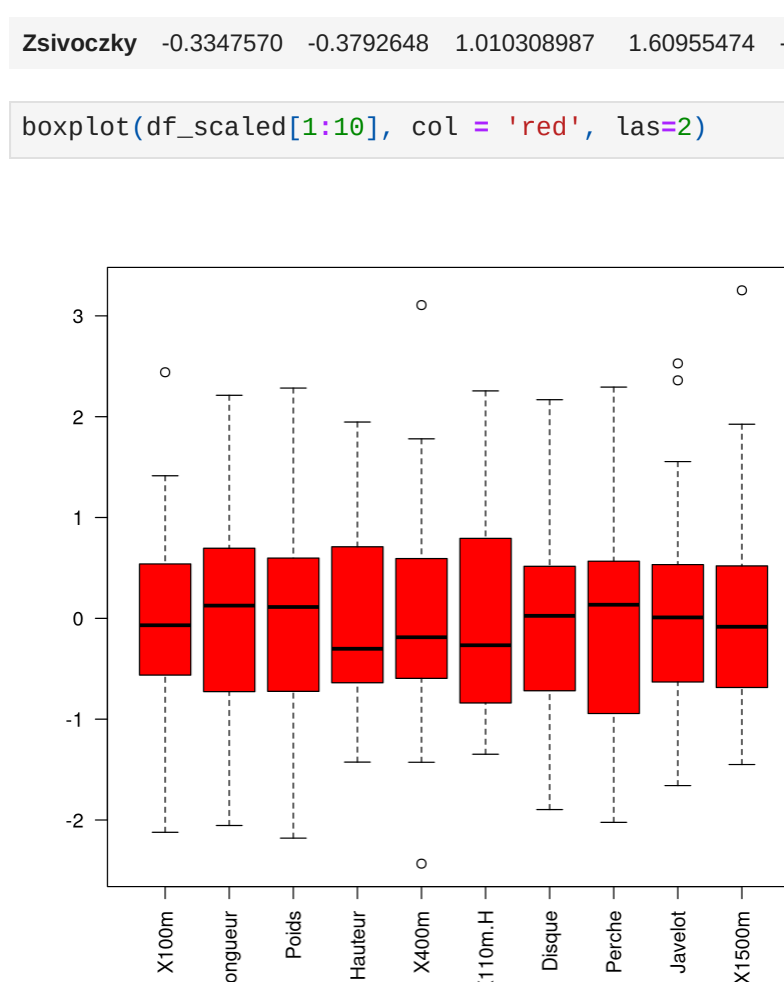
```
In [105]: boxplot(df$Javelot, col = "cyan", main = "Epreuve : Javelot")
```



```
In [106]: boxplot(df$X1500m, col = "cyan", main = "Epreuve : 1500m")
```

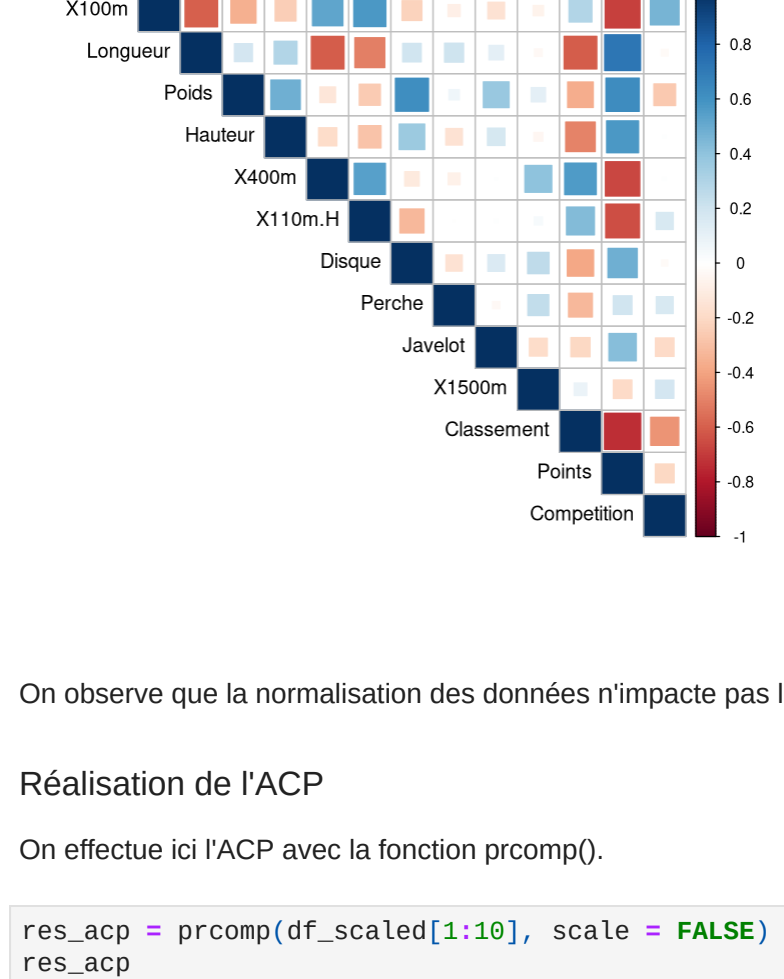


```
In [107]: boxplot(df$Points, col = "cyan", main = "Points")
```



Grâce aux boxplots, il est possible de mieux comprendre la dispersion de nos données dans les différentes colonnes. On voit ici comme avec le `summary` que la plupart des données ne présentent pas d'anomalies. Les seules données qui pourraient soulever des questions sont dans le **400m**, le **javelot** ou le **1500m**, mais celles-ci restent dans des normes tout à fait acceptable.

```
In [108]: options(warn=1)
chart.Correlation(df, pch = 10)
htext("Diagramme de dispersion", side = 3, line = 3)
```



On retrouve ici le **Diagramme de dispersion** appliqué sur nos données ainsi que les valeurs de corrélations entre les différentes colonnes.

Il est notamment assez utile, pour retrouver des applications linéaires entre les colonnes de nos données.

Nous allons passer à nos observations précédentes dans ce graphe de corrélation.

On peut ajouter que :

- les résultats des épreuves de courses courtes, donc **X100m**, **X400m** et **X110m.H**, sont fortement corrélés entre elles. Assez étonnamment elles sont

