

Principal Component Analysis (PCA)

Pablo E. Gutiérrez-Fonseca

8/26/2021

1. Primer paso: cargar las librerías que necesitas.

```
library(ggplot2)
library(dplyr)
library(missMDA) # Impute
library(ggfortify) # autoplot()
library(cluster) #pam
library(factoextra) #get_pca_var()
library(data.table) # data.table()
library(labdsv) #loadings.pca(pca)

library(devtools)

install_github("vqv/ggbiplot") #ggbiplot
library(ggbiplot)
```

2. Segundo paso: cargar los datos.

```
channel <- read.csv("data/channel_form.csv", header=TRUE)
#head(channel)
```

2.1 Vamos a examinar los datos

```
summary(channel)
```

```
##      Forma      NAN_Am      NADBO      NAtemp
## Length:138      Min.   :0.0200      Min.    : 1.310      Min.    :14.67
## Class :character 1st Qu.:0.0400      1st Qu.: 1.930      1st Qu.:24.30
## Mode  :character Median :0.2150      Median : 3.000      Median :26.05
##                      Mean  :0.3201      Mean   : 6.164      Mean   :25.84
##                      3rd Qu.:0.5000      3rd Qu.: 8.585      3rd Qu.:27.70
##                      Max.   :1.5000      Max.   :34.900      Max.   :32.18
##                      NA's    :35
##      nit      NASat02      Elevacion      Ancho
## Min.   : 0.00      Min.    : 23.43      Min.    : 3.00      Min.    : 1.000
## 1st Qu.: 0.40      1st Qu.: 86.24      1st Qu.: 25.25      1st Qu.: 2.000
## Median : 0.92      Median : 94.59      Median : 53.00      Median : 3.000
## Mean   :12.00      Mean   : 91.05      Mean   :230.89      Mean   : 3.822
## 3rd Qu.: 1.62      3rd Qu.:100.52      3rd Qu.:269.25      3rd Qu.: 3.000
## Max.   :324.11      Max.   :122.73      Max.   :2370.00      Max.   :16.000
```

```
## NA's :57 NA's :3
## Velocidad Rocas Canto grava
## Min. : 0.000 Min. : 0.00 Min. : 0.00 Min. : 0.0
## 1st Qu.: 3.000 1st Qu.: 0.00 1st Qu.: 0.00 1st Qu.: 2.5
## Median :11.000 Median :10.00 Median :25.00 Median :20.0
## Mean : 9.133 Mean :16.25 Mean :25.65 Mean :17.8
## 3rd Qu.:14.000 3rd Qu.:30.00 3rd Qu.:40.00 3rd Qu.:25.0
## Max. :16.000 Max. :90.00 Max. :80.00 Max. :80.0
## NA's :3 NA's :3 NA's :4 NA's :3
## arena Limo
## Min. : 0.00 Min. : 0.00
## 1st Qu.: 10.00 1st Qu.: 0.00
## Median : 15.00 Median : 10.00
## Mean : 19.79 Mean : 20.62
## 3rd Qu.: 25.00 3rd Qu.: 25.00
## Max. :100.00 Max. :100.00
## NA's :3 NA's :3
```

2.1 Remover la(s) variable(s) que tiene(n) mucho(s) NAs y las Etiquetas (a la funcion lo le gusta), luego las agregamos.

```
channel_1 <- select(channel, -Forma)
summary(channel_1)
```

```
## NAN_Am NADBO NAtemp nit
## Min. :0.0200 Min. : 1.310 Min. :14.67 Min. : 0.00
## 1st Qu.:0.0400 1st Qu.: 1.930 1st Qu.:24.30 1st Qu.: 0.40
## Median :0.2150 Median : 3.000 Median :26.05 Median : 0.92
## Mean :0.3201 Mean : 6.164 Mean :25.84 Mean :12.00
## 3rd Qu.:0.5000 3rd Qu.: 8.585 3rd Qu.:27.70 3rd Qu.: 1.62
## Max. :1.5000 Max. :34.900 Max. :32.18 Max. :324.11
## NA's :35 NA's :57
## NASat02 Elevacion Ancho Velocidad
## Min. : 23.43 Min. : 3.00 Min. : 1.000 Min. : 0.000
## 1st Qu.: 86.24 1st Qu.: 25.25 1st Qu.: 2.000 1st Qu.: 3.000
## Median : 94.59 Median : 53.00 Median : 3.000 Median :11.000
## Mean : 91.05 Mean : 230.89 Mean : 3.822 Mean : 9.133
## 3rd Qu.:100.52 3rd Qu.: 269.25 3rd Qu.: 3.000 3rd Qu.:14.000
## Max. :122.73 Max. :2370.00 Max. :16.000 Max. :16.000
## NA's :3 NA's :3
## Rocas Canto grava arena
## Min. : 0.00 Min. : 0.00 Min. : 0.0 Min. : 0.00
## 1st Qu.: 0.00 1st Qu.: 0.00 1st Qu.: 2.5 1st Qu.: 10.00
## Median :10.00 Median :25.00 Median :20.0 Median : 15.00
## Mean :16.25 Mean :25.65 Mean :17.8 Mean : 19.79
## 3rd Qu.:30.00 3rd Qu.:40.00 3rd Qu.:25.0 3rd Qu.: 25.00
## Max. :90.00 Max. :80.00 Max. :80.0 Max. :100.00
## NA's :3 NA's :4 NA's :3 NA's :3
## Limo
## Min. : 0.00
## 1st Qu.: 0.00
## Median : 10.00
## Mean : 20.62
```

```
## 3rd Qu.: 25.00
## Max.    :100.00
## NA's    :3
```

2.2 Vamos a imputar datos. Esto es comun para set de datos de campo, los cuales tienden a tener ceros (por mal funcionamiento de los equipos, condiciones climáticas adversas que no podemos ir al campo). Se realiza como un paso preliminar para para realizar un PCA en un set de datos completos.

Mas informacion aca: <https://www.rdocumentation.org/packages/missMDA/versions/1.18/topics/imputePCA>

Primero separar e imputar los datos de sustrato y los fisicoquimicos por aparte.

```
# df0 <- channel_1[, 6:13]
# df0

df1 <- select(channel_1, Elevacion, Ancho, Velocidad, Rocas, Canto, grava, arena, Limo)

df1a <- imputePCA(df1, ncp=4, scale = TRUE, method = c("Regularized","EM"),
  row.w = NULL, ind.sup=NULL, quanti.sup=NULL, quali.sup=NULL,
  coeff.ridge = 1, threshold = 1e-06, seed = NULL, nb.init = 1,
  maxiter = 1000)

df2 <- select(channel_1, Elevacion, NAN_Am, NAtemp, NASatO2, nit, NADBO)
df2a <- imputePCA(df2, ncp=4, scale = TRUE, method = c("Regularized","EM"),
  row.w = NULL, ind.sup=NULL, quanti.sup=NULL, quali.sup=NULL,
  coeff.ridge = 1, threshold = 1e-06, seed = NULL, nb.init = 1,
  maxiter = 1000)
```

Unir las dos tablas y seleccionar las columnas para hacer el PCA.

```
df1b <- as.data.frame(df1a) # Sustrata
df2b <- as.data.frame(df2a) # Physicochemical

new_channel <- do.call("merge", c(lapply(list(df1b, df2b), data.frame, row.names=NULL),
  by = 0, all = TRUE, sort = FALSE))[-1]

new_channel2 <- select(new_channel,
  completeObs.Elevacion.x, completeObs.Ancho, completeObs.Velocidad,
  completeObs.Rocas, completeObs.Canto, completeObs.grava, completeObs.arena,
  completeObs.Limo, completeObs.NAtemp, completeObs.NASatO2,
  )
```

3. Vamos a correr el PCA

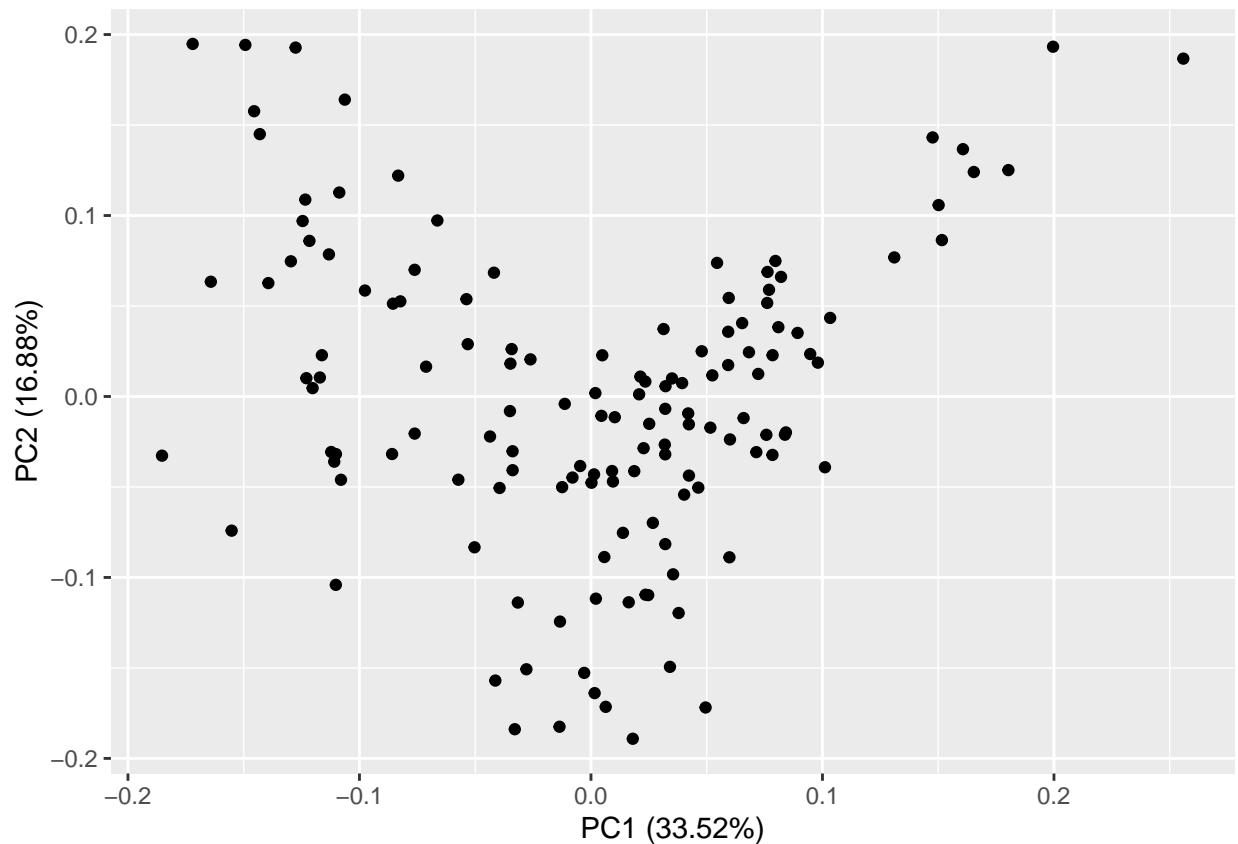
```
channel.pca <- prcomp(new_channel2, center = TRUE, scale =TRUE)
summary(channel.pca)
```

```
## Importance of components:
##              PC1      PC2      PC3      PC4      PC5      PC6      PC7
```

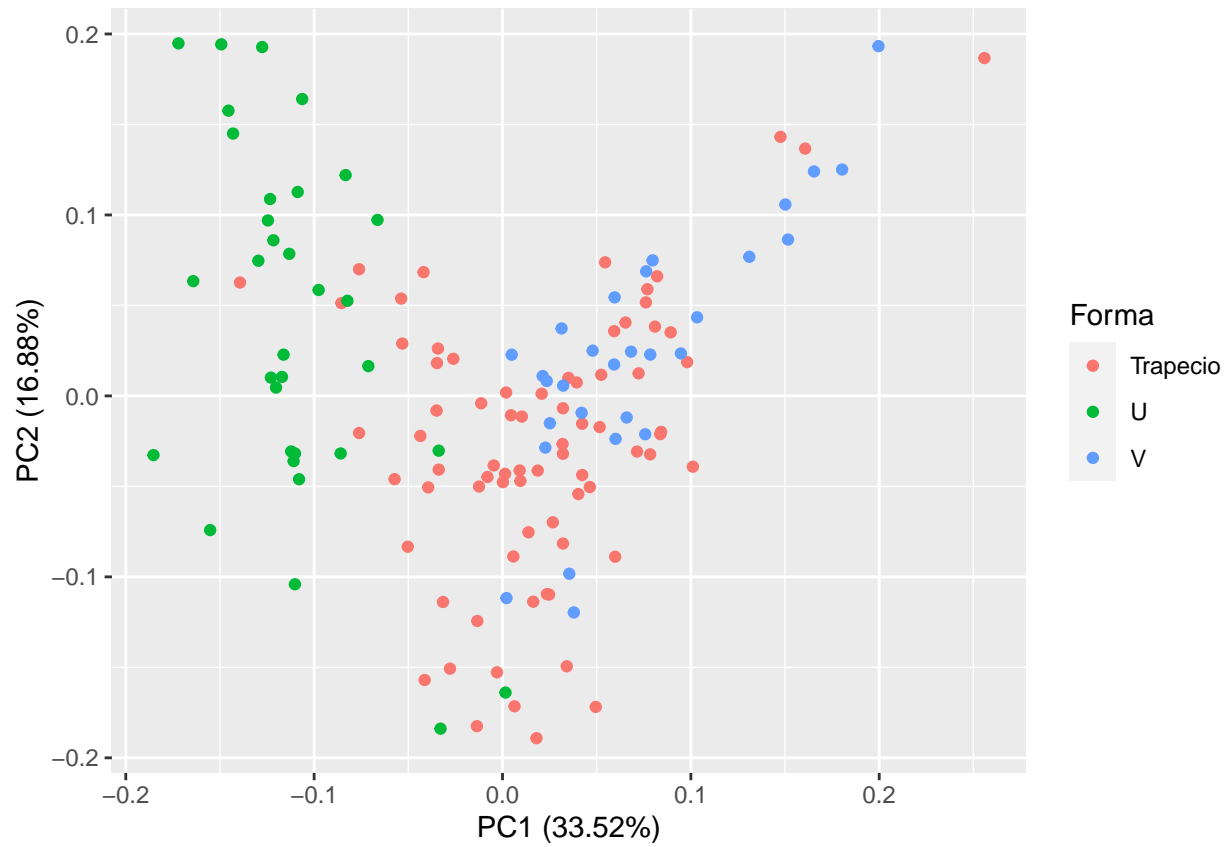
```
## Standard deviation      1.8308 1.2991 1.1906 1.0728 0.90108 0.75065 0.70276
## Proportion of Variance 0.3352 0.1688 0.1417 0.1151 0.08119 0.05635 0.04939
## Cumulative Proportion 0.3352 0.5040 0.6457 0.7608 0.84198 0.89833 0.94772
##                          PC8    PC9    PC10
## Standard deviation      0.61343 0.3756 0.07403
## Proportion of Variance 0.03763 0.0141 0.00055
## Cumulative Proportion 0.98535 0.9994 1.00000
```

3.1 Vamos a ver el grafico.

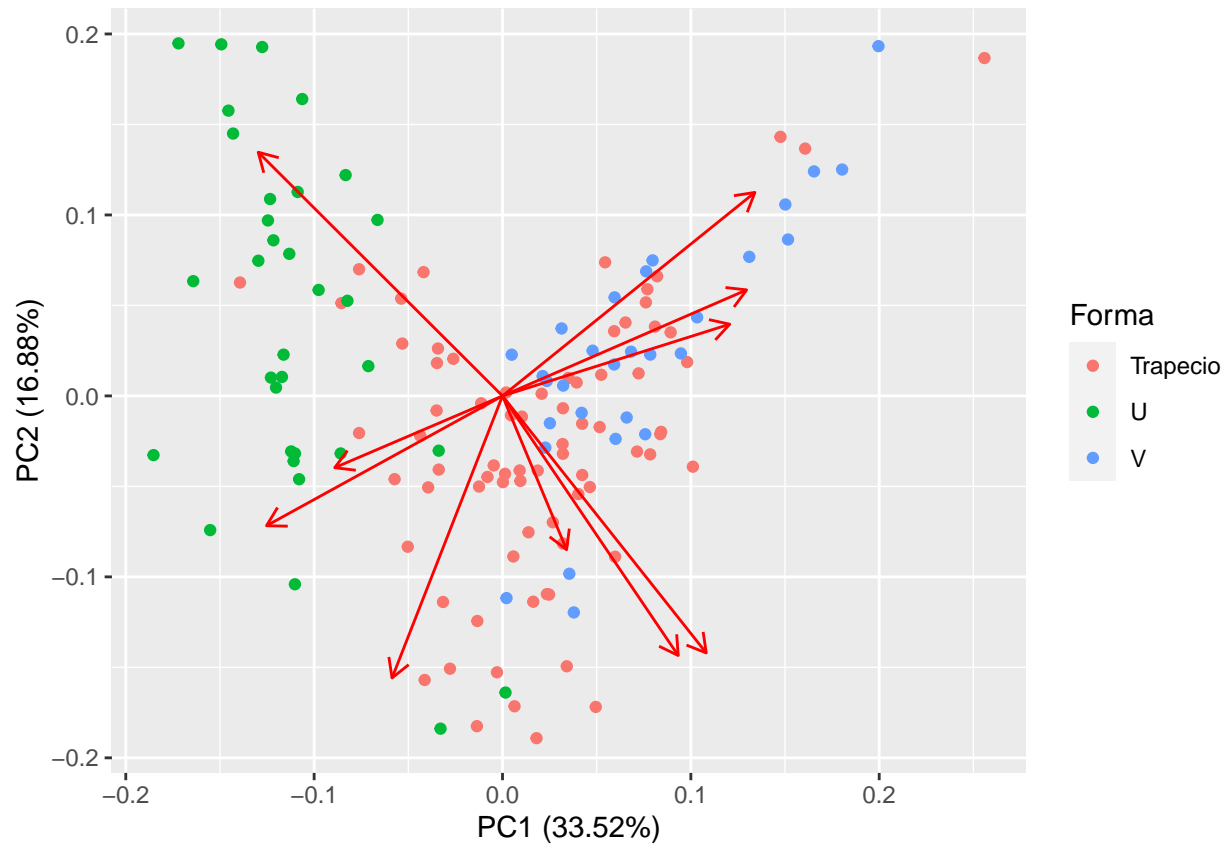
```
autoplot(channel.pca)
```



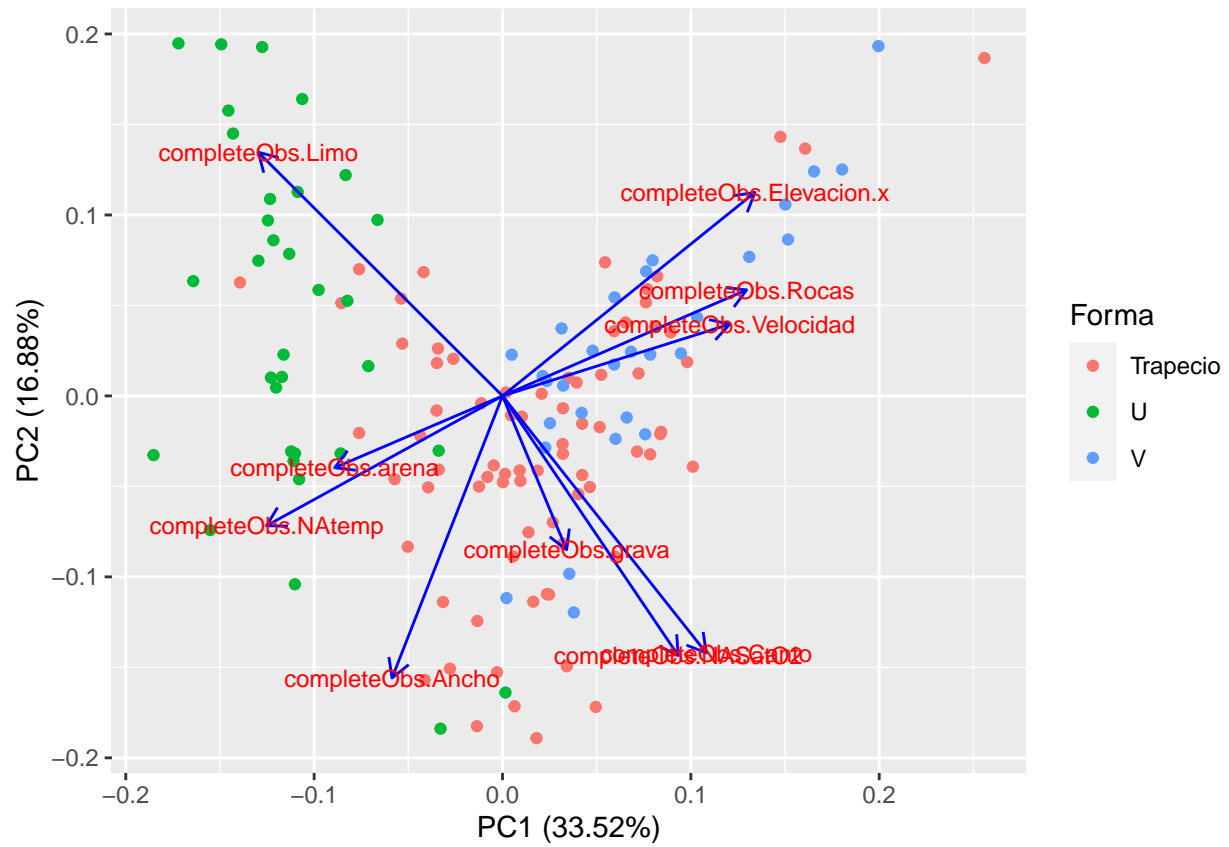
```
autoplot(channel.pca, data = channel, colour = 'Forma')
```



```
autoplot(channel.pca, data = channel, colour = 'Forma', loadings = TRUE)
```

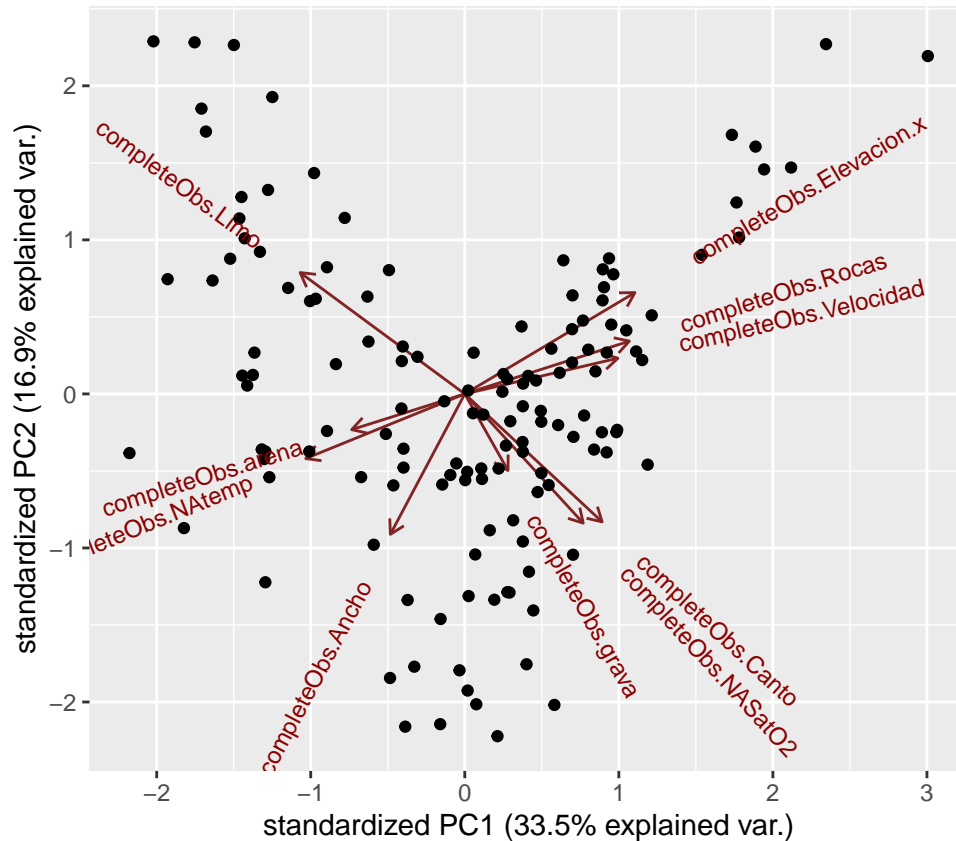


```
autoplot(channel.pca, data = channel, colour = 'Forma', loadings = TRUE,
         loadings.colour = 'blue',
         loadings.label = TRUE, loadings.label.size = 3)
```



Otra manera de ver el grafico

```
ggbiplot(channel.pca, labels=rownames(channel$Forma))
```



3.2 Vamos a ver la contribucion de cada una de las variables.

```
variance <- (channel.pca$sdev)^2

# Cargar los loadings
loadings <- channel.pca$rotation
round(loadings, 2)[ , 1:3]
```

```
##           PC1  PC2  PC3
## completeObs.Elevacion.x  0.40  0.33 -0.20
## completeObs.Ancho       -0.17 -0.46 -0.45
## completeObs.Velocidad   0.36  0.12  0.42
## completeObs.Rocas       0.38  0.17 -0.19
## completeObs.Canto       0.32 -0.42 -0.08
## completeObs.grava       0.10 -0.25  0.61
## completeObs.arena      -0.26 -0.12 -0.12
## completeObs.Limo       -0.38  0.40 -0.07
## completeObs.NAtemp     -0.37 -0.21  0.36
## completeObs.NASatO2     0.28 -0.42 -0.13
```

```
print(channel.pca)
```

```
## Standard deviations (1, ..., p=10):
## [1] 1.83078436 1.29914774 1.19059462 1.07275990 0.90108032 0.75064777
## [7] 0.70275954 0.61343402 0.37555814 0.07402975
```



```
##
## Rotation (n x k) = (10 x 10):
##
```

	PC1	PC2	PC3	PC4
## completeObs.Elevacion.x	0.3956074	0.3319346	-0.19970962	0.262046675
## completeObs.Ancho	-0.1730797	-0.4600987	-0.44918134	0.088453085
## completeObs.Velocidad	0.3561030	0.1169566	0.41659923	0.004781697
## completeObs.Rocas	0.3826021	0.1728855	-0.18751252	-0.096969115
## completeObs.Canto	0.3193348	-0.4192288	-0.07809951	-0.331598120
## completeObs.grava	0.1004574	-0.2508480	0.61179737	0.324655198
## completeObs.arena	-0.2629454	-0.1171209	-0.11791669	0.692074307
## completeObs.Limo	-0.3828069	0.3975173	-0.07389293	-0.324593592
## completeObs.NAtemp	-0.3698073	-0.2117725	0.36336704	-0.310643419
## completeObs.NASat02	0.2753892	-0.4234359	-0.13356891	-0.133776783

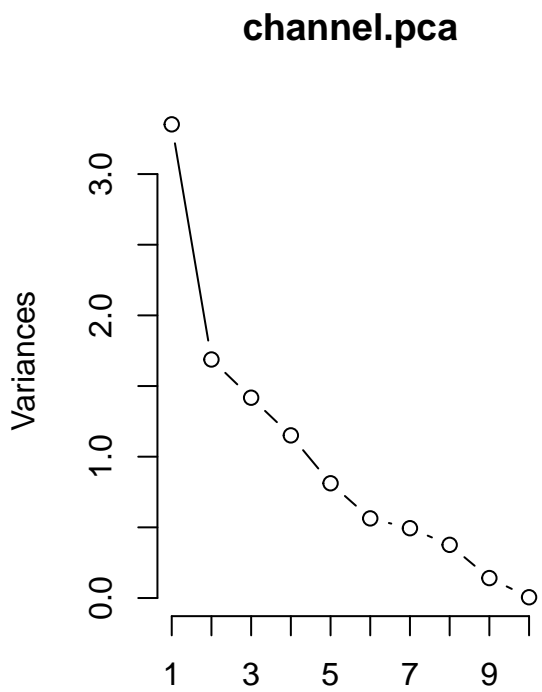
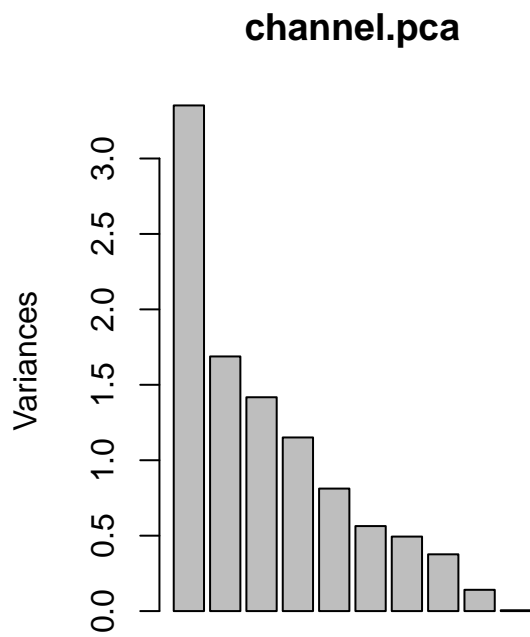
	PC5	PC6	PC7	PC8
## completeObs.Elevacion.x	-0.26927221	-0.05245244	-0.18409079	0.13170674
## completeObs.Ancho	-0.28352198	-0.22183848	0.21212512	-0.58059003
## completeObs.Velocidad	0.06851242	0.48059224	0.19734141	-0.63431592
## completeObs.Rocas	0.51040363	-0.47888909	0.33062616	-0.07874652
## completeObs.Canto	-0.31634667	0.34773472	0.22985746	0.36337021
## completeObs.grava	-0.28727042	-0.46271562	-0.19568813	-0.03041943
## completeObs.arena	0.33320571	0.36772659	0.05763329	0.11070131
## completeObs.Limo	-0.16997384	0.06742557	-0.35248106	-0.25014094
## completeObs.NAtemp	0.32096625	-0.08770394	0.16443429	0.10651154
## completeObs.NASat02	0.39543136	0.07243312	-0.72603361	-0.13370722

	PC9	PC10
## completeObs.Elevacion.x	-0.70558764	0.020323763
## completeObs.Ancho	-0.19257159	-0.001435631
## completeObs.Velocidad	-0.09363279	-0.014273068
## completeObs.Rocas	0.04081289	-0.414891885
## completeObs.Canto	-0.02935453	-0.446639620
## completeObs.grava	0.06570440	-0.327292380
## completeObs.arena	-0.07241297	-0.396487492
## completeObs.Limo	-0.03948325	-0.602335535
## completeObs.NAtemp	-0.66357759	0.020739455
## completeObs.NASat02	-0.04809287	0.014074170

```
rownames(loadings) <- colnames(new_channel2)
scores <- channel.pca$x
```

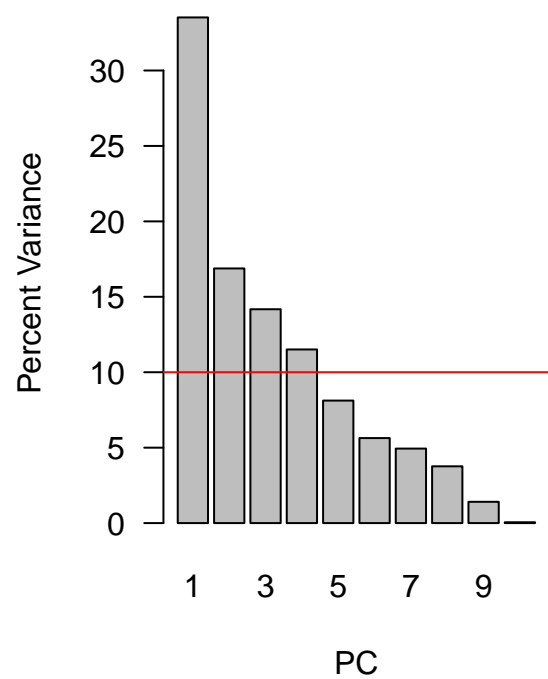
3.3 Ver graficamente lo que explica cada axis.

```
layout(matrix(1:2, ncol=2))
screeplot(channel.pca)
screeplot(channel.pca, type="lines")
```



```
varPercent <- variance/sum(variance) * 100
barplot(varPercent, xlab='PC', ylab='Percent Variance',
names.arg=1:length(varPercent), las=1, col='gray') +
abline(h=1/ncol(new_channel2)*100, col="red")
```

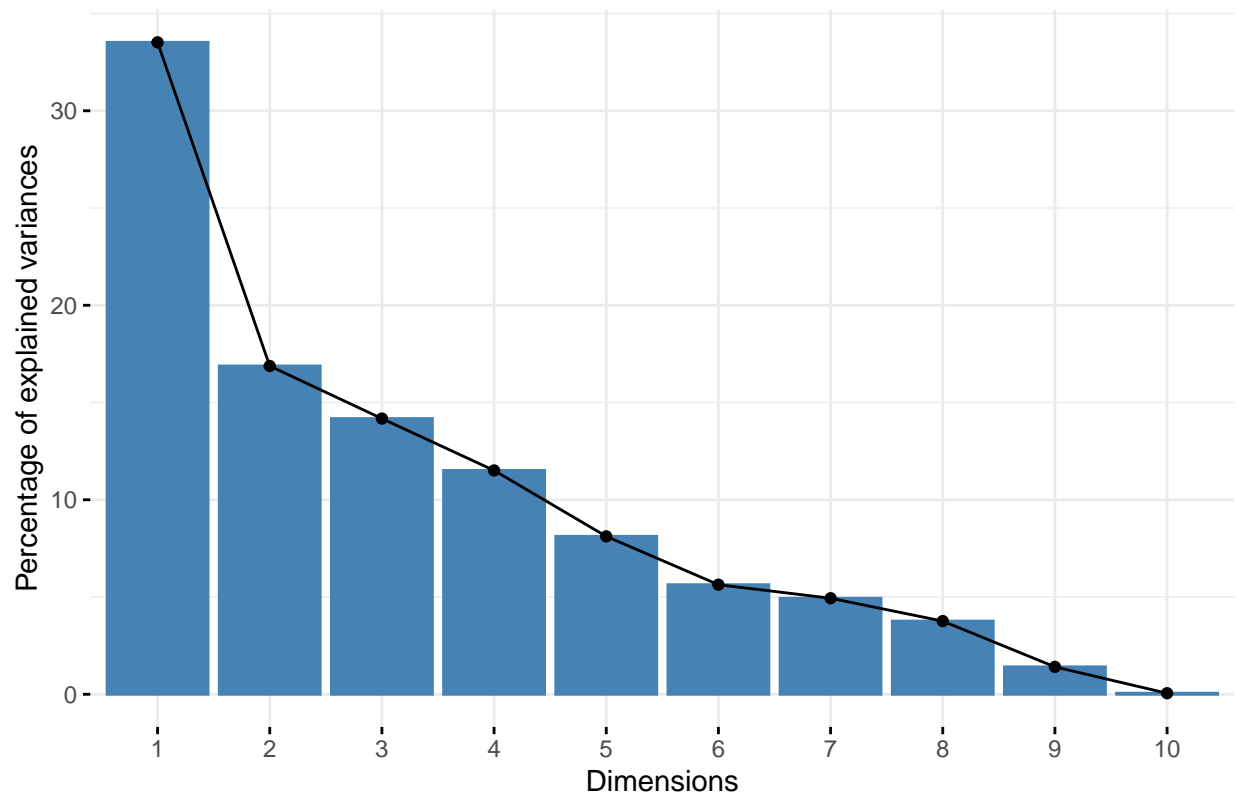
```
## numeric(0)
```



4 Otras formas de visualizar los datos.

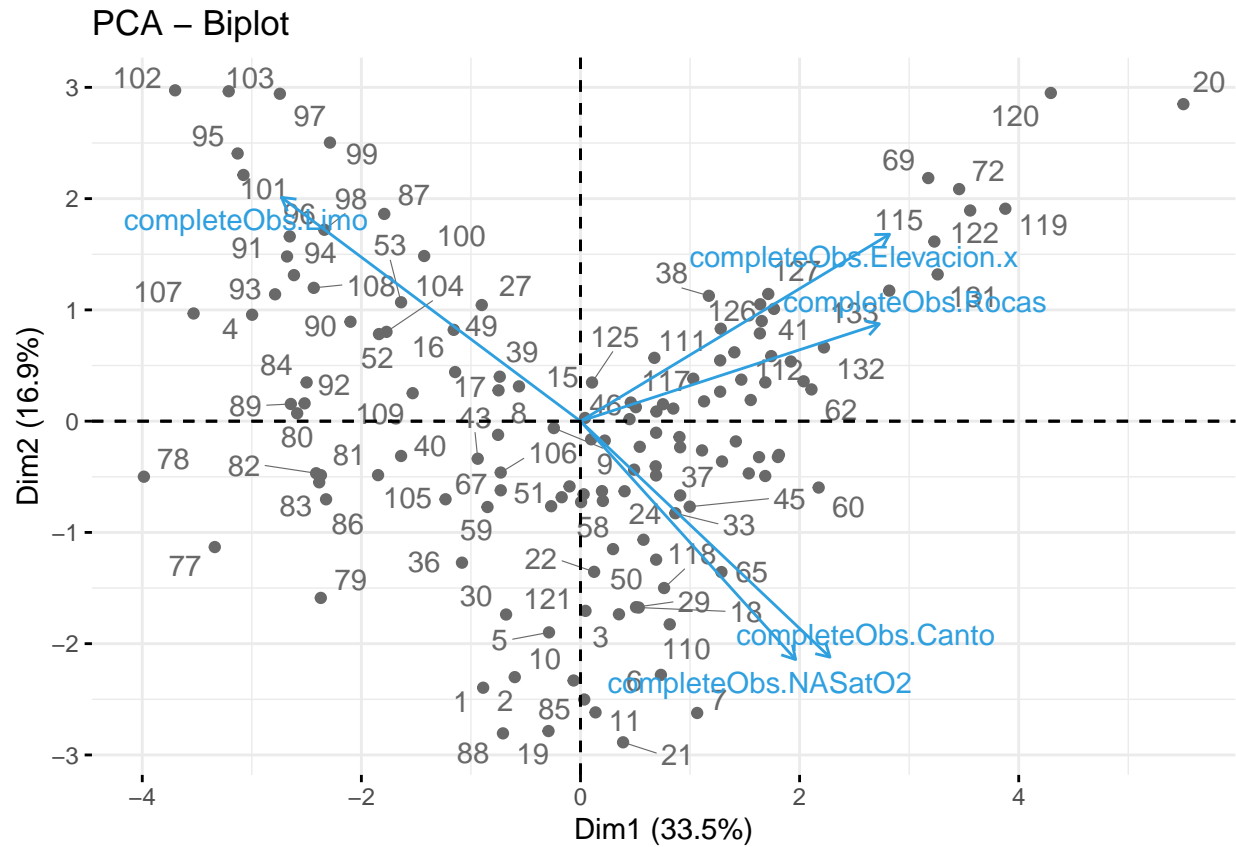
```
fviz_eig(channel.pca)
```

Scree plot



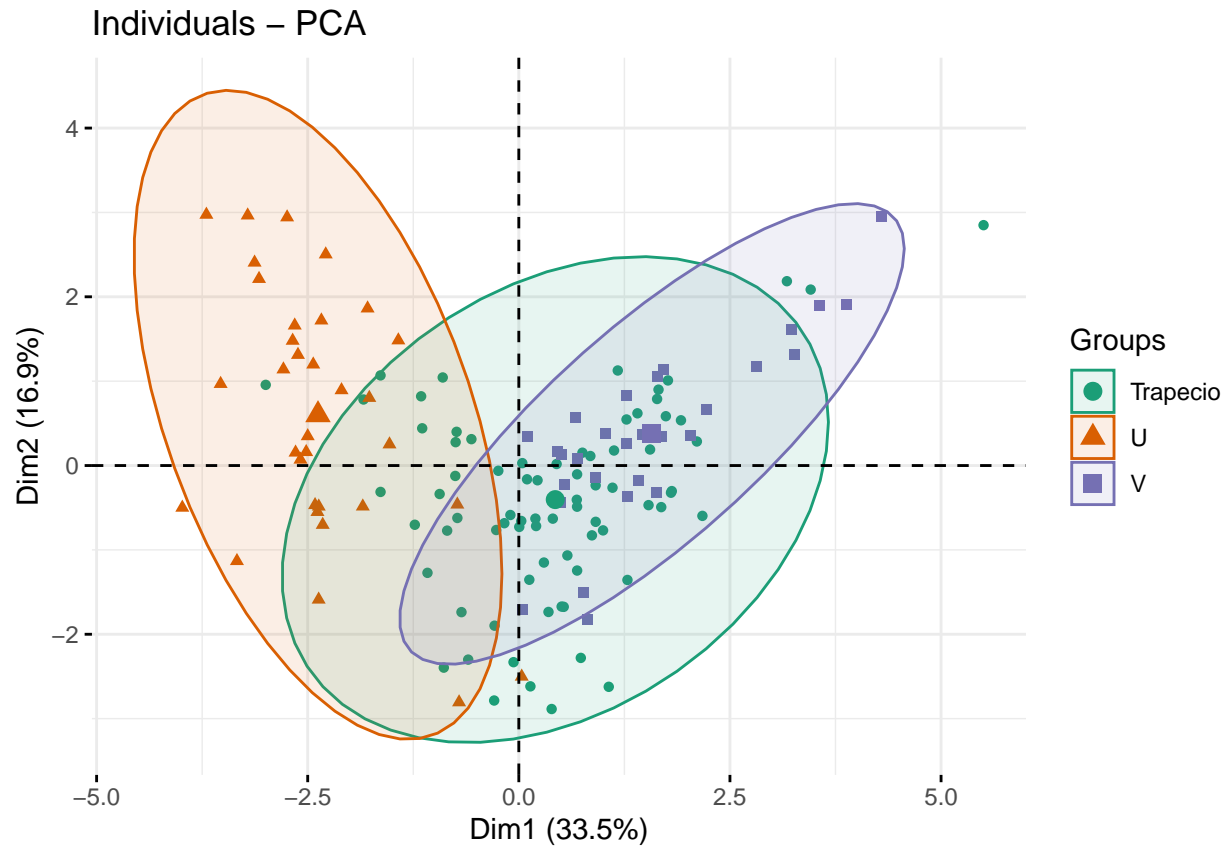
```
fviz_pca_biplot(channel.pca, repel = TRUE,  
  col.var = "#2E9FDF", # Variables color  
  col.ind = "#696969", # Individuals color  
  select.var = list(contrib = 5))
```

```
## Warning: ggrepel: 43 unlabeled data points (too many overlaps). Consider  
## increasing max.overlaps
```



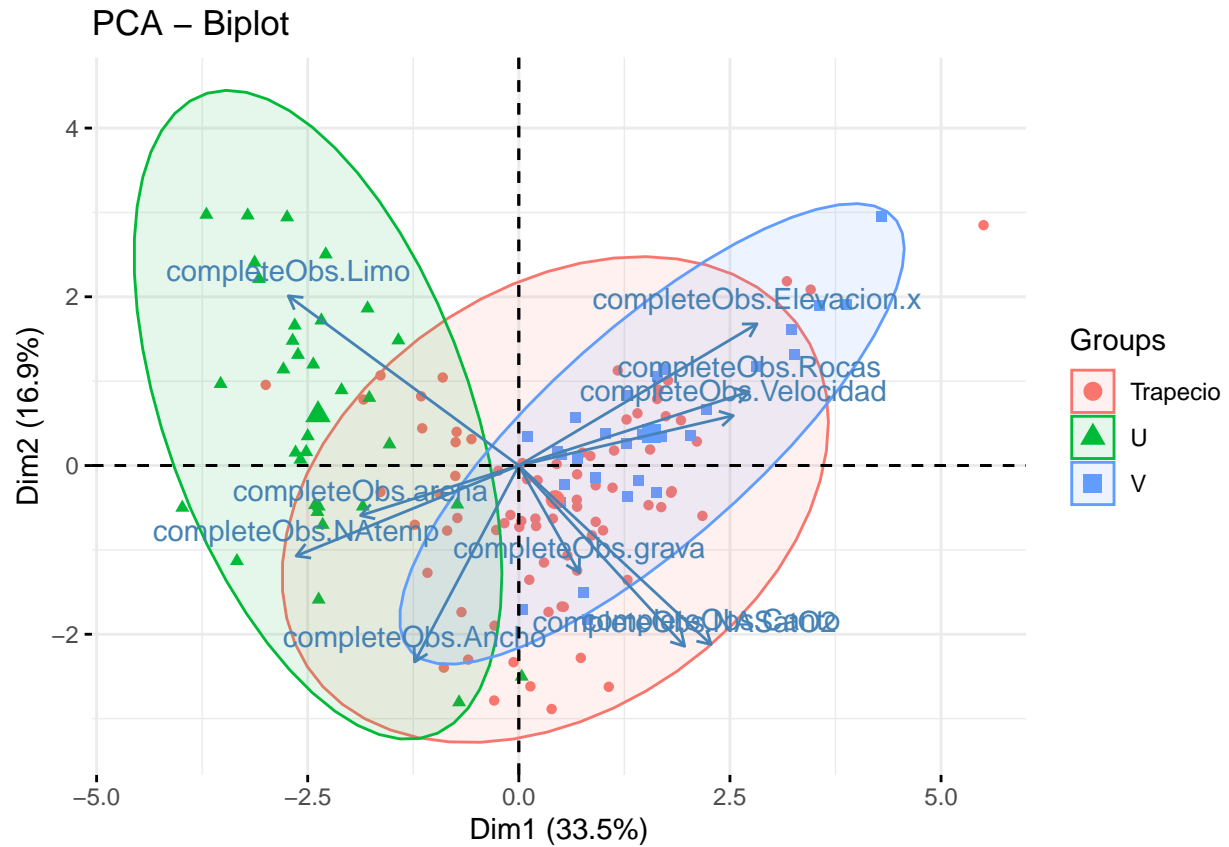
4.1 Con las elipses.

```
fviz_pca_ind(channel.pca, label="none", habillage=channel$Forma,
  addEllipses=TRUE, ellipse.level=0.95, palette = "Dark2")
```



4.2

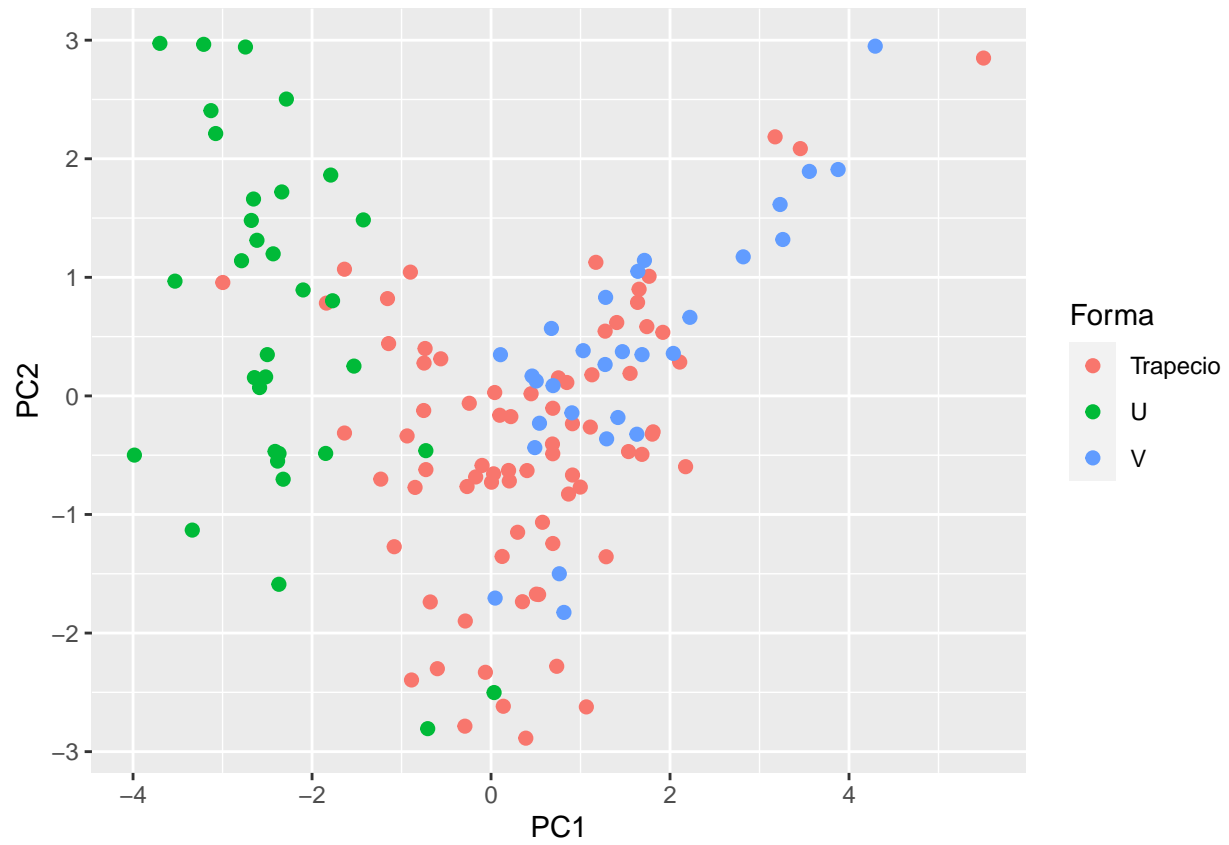
```
fviz_pca_biplot(channel.pca, label = "var", habillage=channel$Forma,  
  addEllipses=TRUE, ellipse.level=0.95,  
  ggtheme = theme_minimal())
```



5. Convertirlo en una data.frame para trabajarlo en ggplot2

```
data <- data.table(PC1=channel.pca$x[,1], PC2=channel.pca$x[,2], Forma= channel[,1])
data <- data[order(channel$Forma),]

ggplot(data, aes(x=PC1,y=PC2)) +
  geom_point(size = 2, aes(color=Forma))
```



6. Otras enlaces de interes.

<http://www.sthda.com/english/articles/31-principal-component-methods-in-r-practical-guide/118-princip>