

Estadística interactiva con RNotebook.io

Tarea 1 de Innovación Docente e Iniciación a la Investigación Educativa

Antonio García Sevilla

Subtarea 1

Rnotebook.io

El Recurso Educativo en Abierto que vamos a analizar en este documento se llama “Rnotebook.io”¹.

¹ <https://rnotebook.io/>

Éste es un servicio online que nos ofrece la posibilidad de crear “Jupyter notebooks”² en la nube, visualizarlos con un navegador, o incluso distribuirlos entre nuestros alumnos simplemente compartiendo una dirección web.

² <https://jupyter.org/>

Jupyter Notebooks

Estos “notebooks” son documentos interactivos escritos en una mezcla de código y texto normal, que nos permiten usar un lenguaje de programación de manera dinámica con cualquier ordenador, sin tener que instalar nada.

En el caso que nos ocupa, el lenguaje que se utiliza es “R”, un conocido software estadístico.

R

R³ es a la vez un lenguaje de programación y una colección de recursos estadísticos. Es parecido a otros lenguajes para matemáticas, como Matlab u Octave, pero completamente orientado a la estadística y el análisis de datos.

³ <https://www.r-project.org/>

R nos permite escribir nuestras propias funciones matemáticas, pero su verdadera potencia consiste en que trae ya funciones para calcular operaciones comunes, como el modelado de distribuciones, test de hipótesis, etc. Incluso permite generar gráficas de manera dinámica, que se pueden visualizar inmediatamente en el mismo documento.

Valoración general

“R” es una herramienta muy utilizada en el análisis de datos. Puede resultar complicado de aprender, pues necesita de conocimientos de programación. Sin embargo, al tener muchas funciones ya codificadas en el propio lenguaje, se pueden usar notebooks que ya haya creado un experto sin tener que saber programar.

Valor didáctico

Desde el punto de vista de la docencia de la estadística, tiene muchos puntos fuertes. El análisis de datos es una disciplina muy visual, por lo que poder crear gráficas fácilmente en el ordenador es muy útil. La naturaleza dinámica de estas gráficas hace que podamos responder al alumno de manera inmediata y práctica, no sólo teórica.

Características prácticas

Desde el punto de vista pragmático, si se sabe utilizar, la herramienta es potentísima. Nos permite crear de antemano los documentos a usar en clase, e incluso imprimirlos en papel, pero luego modificarlos dinámicamente durante la lección y observar los resultados.

Además, se pueden crear documentos con preguntas, o con secciones sin rellenar, para que los alumnos los puedan utilizar interactivamente bien en clase bien como tarea para casa. Puede entregárseles a los alumnos una plantilla con las partes de código rellenas, y que ellos hagan su análisis exploratorio, visualicen las gráficas, y respondan a alguna pregunta.

Instalación

Una de las ventajas de la herramienta es que, al estar alojada en la nube, no requiere ningún tipo de instalación local, ahorrándonos posibles problemas.

Pero esto puede ser también una desventaja, pues es un servicio en fase de prueba y no se sabe cómo va a evolucionar. Sin embargo, las herramientas en las que se basa son de código abierto, por lo que con un poco de conocimiento técnico podría ser instalado en un aula o laboratorio de informática.

Subtarea 2

Demostración de análisis exploratorio y de componentes principales

En esta subtarea, utilizamos la herramienta “RNotebook” para crear un documento interactivo que nos ayude a explicar, de forma dinámica y visual, el análisis de datos exploratorio. El documento está disponible en la nube.⁴

Este documento se puede imprimir en papel o formato *PDF* para distribuir, pero su mayor punto fuerte es que también se puede utilizar de manera interactiva en la web. Así, durante la clase se puede ir desarrollando las instrucciones, utilizando el documento existente como guía, o incluso pedir a los alumnos que sean ellos mismos quienes van realizando el análisis sobre la marcha, siguiendo las instrucciones del “notebook” ya creado.

El notebook completo se puede visualizar en la URL 4, o también se puede encontrar en formato *PDF* a continuación de este documento.

⁴ <https://rnotebook.io/anon/ab58fdfb69f9ef90/notebooks/An%C3%A1lisis%20de%20componentes%20principales.ipynb>

Utilidad del recurso

El uso de “RNotebook” ha supuesto dos grandes ventajas.

Desde el punto de vista del docente, la realización en sí ha sido muy cómoda. Primero he ido realizando los experimentos, cambiando los datos y funciones, y eligiendo las gráficas a mostrar. Al poder hacerlo todo de manera interactiva en el mismo documento, la experiencia ha sido muy gratificante. Además, no ha sido necesario utilizar ningún programa adicional para los diagramas, y exportar el documento en *PDF* ha sido muy fácil.

Para el alumno, la posibilidad de interaccionar con el documento aporta un valor didáctico muy alto. Se pueden cambiar los datos y las funciones sobre la marcha, permitiendo responder a los “¿Y qué pasaría si...?” que son la base del aprendizaje conectivo. Las gráficas se generan con los parámetros elegidos por el usuario, cambian si cambiamos los datos, y nos permiten experimentar con las herramientas respondiendo a las preguntas según vayan surgiendo.

Conclusión

Antes de decidirme por “RNotebook”, una amiga que es profesora en un instituto me recomendó otra herramienta: **Geogebra**⁵. Yo también conocía previamente **Wolfram Alpha**⁶. Estas herramientas, así como “RNotebook”, sirven para hacer las matemáticas más interactivas. Las matemáticas muchas veces son muy visuales, y la exploración y experimentación ayudan a su entendimiento. Sin em-

⁵ <https://www.geogebra.org/>

⁶ <https://www.wolframalpha.com/>

bargo, los medios tradicionales son más limitados para esto, pues es más complicado dibujar una función o una gráfica sobre la marcha, o realizar cálculos complejos en los que el detalle no es tan relevante como la conclusión. Los medios digitales nos ayudan con esto, y nos permiten crear clases muy didácticas y dinámicas.

Entre “Geogebra”, “Wolfram Alpha” y “RNotebook”, la última es la menos enfocada a la educación. No está provista de conceptos como clase o lecciones, ni incluye recursos didácticos ya existentes. A pesar de ello, siempre he creído que uno de los déficits más grandes de los estudiantes actuales está en la estadística descriptiva. RNotebook me permitía realizar una pequeña lección sobre esto, y por eso la elegí.

Gracias a elaborar este documento he aprendido bastante sobre cómo estructurar una clase de estadística descriptiva, además de, como siempre, cimentar mis conocimientos en el área a explicar.

Notebook A continuación se encuentra anexo el documento creado, exportado en formato PDF.

Análisis exploratorio

Wikipedia:

El análisis exploratorio de datos es una forma de analizar datos definido por John W. Tukey (E.D.A.: Exploratory data analysis) es el tratamiento estadístico al que se someten las muestras recogidas durante un proceso de investigación en cualquier campo científico. Para mayor rapidez y precisión, todo el proceso suele realizarse por medios informáticos, con aplicaciones específicas para el tratamiento estadístico. Los E.D.A., no necesariamente, se llevan a cabo con una base de datos al uso, ni con una hoja de cálculo convencional; no obstante el programa SPSS y R (lenguaje de programación) son las aplicaciones más utilizadas, aunque no las únicas.

Tras recoger una muestra de datos, antes de procesarlos para la tarea correspondiente, es necesario hacer primero un análisis exploratorio. En este análisis lo que hacemos es *explorar* los datos, mirándolos "a ojo" para descubrir su estructura y relación.

Pero los números por sí solos no dicen nada, o peor, pueden llevarnos a conclusiones equivocadas. Por ello, aunque el análisis sea "a ojo", no quiere decir que no se pueda hacer de manera formal. Existen herramientas estadísticas que nos permiten hacer este análisis de manera rigurosa, para que las conclusiones (**¡eso sí, preliminares!**) que extraigamos de los datos sean sólidas.

Si los datos que tenemos son muchos, este análisis es arduo de hacer a mano. Afortunadamente, la mayor parte de las herramientas estadísticas que queramos usar están implementadas computacionalmente. Veamos un ejemplo:

Dataset

El conjunto de datos "*Iris*" es un conjunto de datos multivariante usado muy a menudo como un ejemplo de análisis discriminante. Nosotros lo vamos a utilizar para el análisis exploratorio, debido a su pequeño tamaño y facilidad de uso, y a que está contenido en *R* por defecto.

Este conjunto de datos se refiere a una población de especímenes de flor de tres especies distintas. Sobre estos *individuos*, se realizaron cuatro medidas: la longitud y anchura de los sépalos, y las de los pétalos. Además, cada espécimen viene identificado con el nombre de la especie a la que pertenece.

```
In [14]: iris[0:10,]
```

Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
5.4	3.9	1.7	0.4	setosa
4.6	3.4	1.4	0.3	setosa
5.0	3.4	1.5	0.2	setosa
4.4	2.9	1.4	0.2	setosa
4.9	3.1	1.5	0.1	setosa

El primer análisis que podemos hacer es observar la distribución observada de las variables, usando el commando `summary`.

```
In [2]: summary(iris)
```

```
 Sepal.Length Sepal.Width Petal.Length Petal.Width
Min.   :4.300  Min.   :2.000  Min.   :1.000  Min.   :0.100
1st Qu.:5.100  1st Qu.:2.800  1st Qu.:1.600  1st Qu.:0.300
Median :5.800  Median :3.000  Median :4.350  Median :1.300
Mean   :5.843  Mean   :3.057  Mean   :3.758  Mean   :1.199
3rd Qu.:6.400  3rd Qu.:3.300  3rd Qu.:5.100  3rd Qu.:1.800
Max.   :7.900  Max.   :4.400  Max.   :6.900  Max.   :2.500

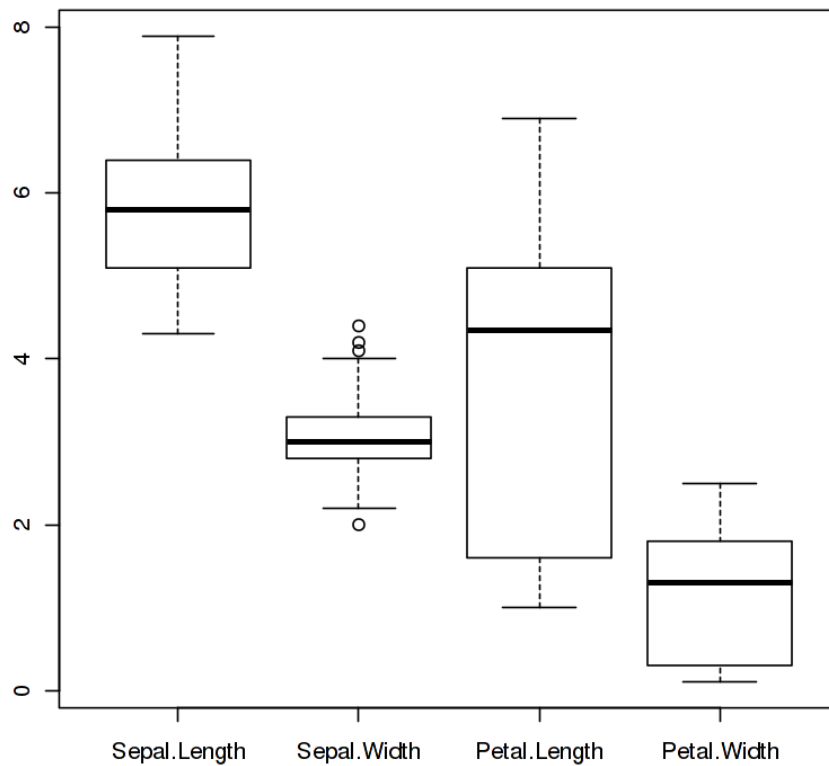
Species
setosa   :50
versicolor:50
virginica :50
```

Separamos las medidas continuas de la etiqueta de especie para facilitar el análisis posterior.

```
In [3]: medidas <- iris[1:4]
        especies <- iris$Species
```

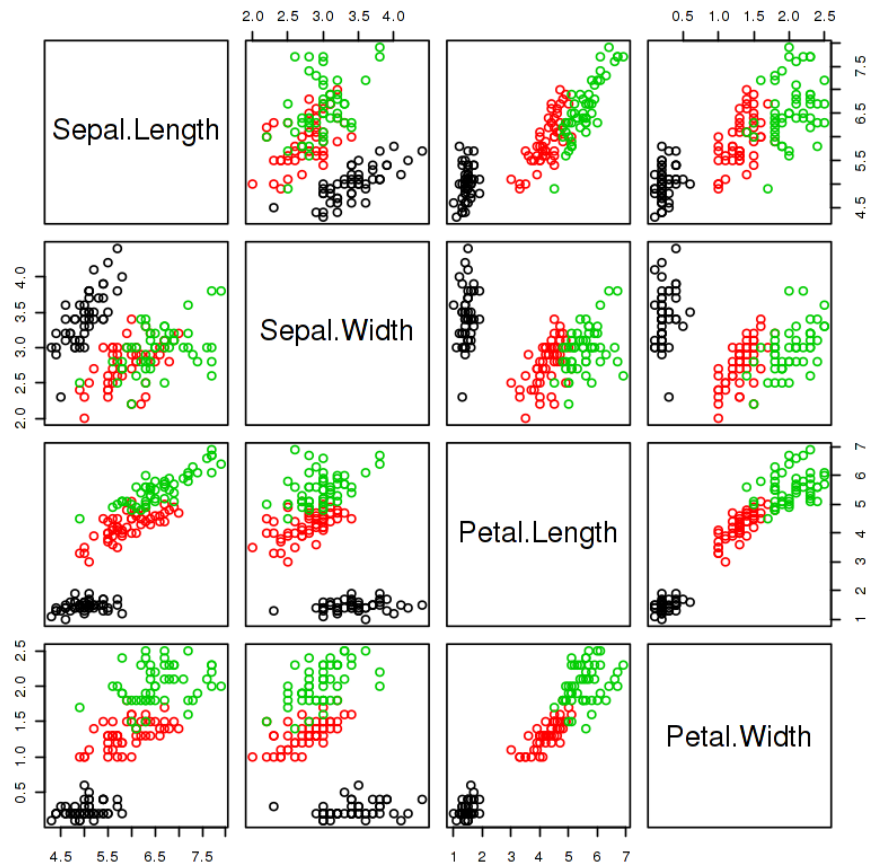
Una buena gráfica siempre resulta más fácil de visualizar que una tabla de números. En este caso, para observar la distribución de las variables, podemos usar un **diagrama de cajas y bigotes**. Este diagrama muestra los cuartiles y varianzas, con lo que nos da una muy buena idea de la distribución.

```
In [4]: boxplot(medidas)
```



Para visualizar la relación de las variables, podemos también usar **diagramas de puntos**, poniendo en los ejes las variables a pares. En este caso, podemos pedirle a R también que nos coloree los puntos según la especie. Aquí empezamos a ver los **clusters**, agrupaciones de individuos que en este caso se corresponden claramente con las especies.

```
In [5]: plot(medidas, col=especies)
```

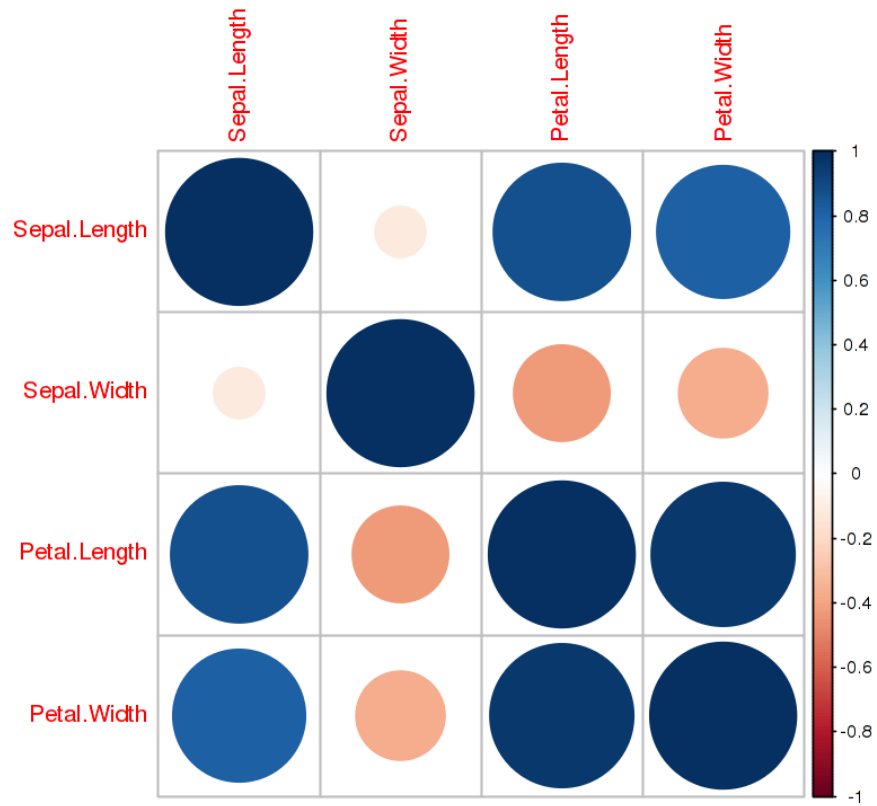


Otra cuestión que podemos ver en los diagramas anteriores, es que hay variables que presentan una fuerte relación lineal. Por ejemplo, observad la longitud y anchura de los pétalos.

Para analizar esto formalmente, podemos obtener la correlación de las medidas, e incluso visualizarla en una gráfica.


```
In [6]: cor(medidas)
library(corrplot)
corrplot(cor(medidas))
```

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width
Sepal.Length	1.0000000	-0.1175698	0.8717538	0.8179411
Sepal.Width	-0.1175698	1.0000000	-0.4284401	-0.3661259
Petal.Length	0.8717538	-0.4284401	1.0000000	0.9628654
Petal.Width	0.8179411	-0.3661259	0.9628654	1.0000000



Como hemos visto antes, hay variables que muestran una gran correlación. Por ejemplo, el ancho y largo de los pétalos. Intuitivamente, esto podríamos verlo como "el tamaño de los pétalos". Pensamos que los pétalos tienden a tener una proporción definida, y simplemente pueden ser más grandes o pequeños. Esta intuición se puede trasladar formalmente de varias maneras. En nuestro análisis exploratorio, una herramienta muy habitual es el Análisis de Componentes Principales.

ACP

De la wikipedia:

[El A]nálisis de componentes principales (en español ACP, en inglés, PCA) es una técnica utilizada para describir un set de datos en términos de nuevas variables ("componentes") no correlacionadas. Los componentes se ordenan por la cantidad de varianza original que describen, por lo que la técnica es útil para reducir la dimensionalidad de un conjunto de datos.

De manera resumida, el ACP busca maneras de "reorientar el espacio" en el que las nuevas dimensiones (componentes) sean independientes. Es decir, cambia las medidas originales por otras, combinación de las anteriores, que no dependan entre ellas. En nuestro caso, cambiaría longitud y anchura de los pétalos por tamaño y proporción, por ejemplo.

Pero, ¿qué ganamos con esto? La idea más importante del ACP es que podemos ordenar las nuevas componentes por "importancia" (calculada por ejemplo a partir de la desviación estándar o la varianza). Así, vemos que la componente "tamaño del pétalo" presenta una gran varianza, y "proporción" poca. Por tanto, podemos eliminar la componente "proporción", y quedarnos sólo con "tamaño". Así, reducimos el número de variables a tratar.

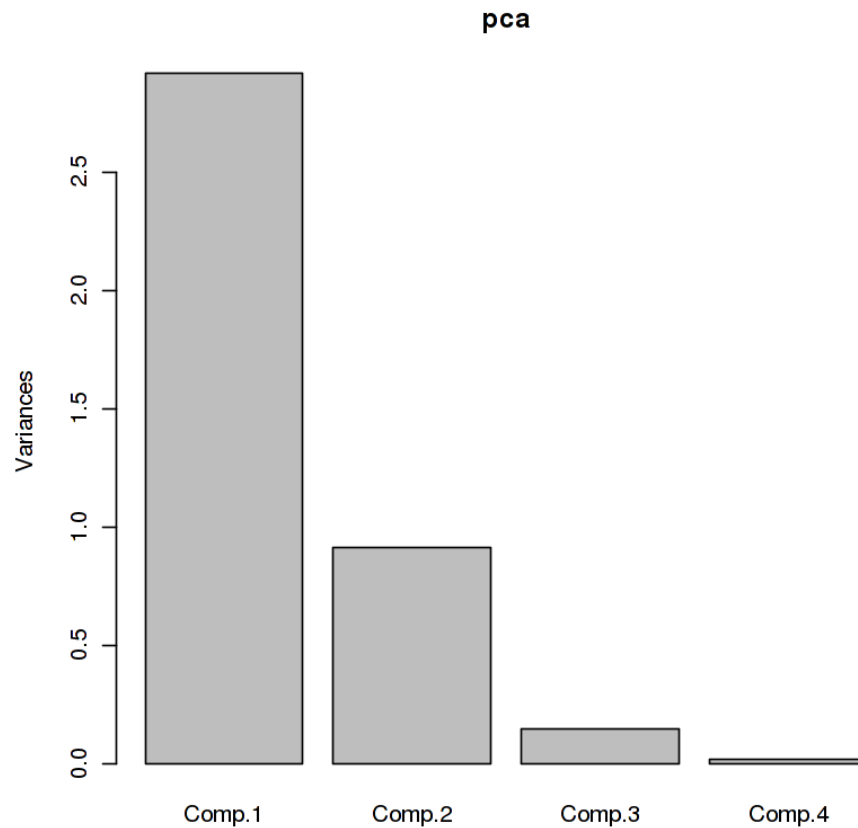
Sin embargo, hay que tener cuidado cuando damos nombres a las componentes. El ACP sólo nos da componente 1, 2, 3, etc. La interpretación que le demos nosotros luego es eso, una interpretación, y por tanto hay que ser cuidadoso con ella. Es más, en ocasiones tendremos componentes muy explicativas (es decir, que justifican gran parte de la varianza de los individuos) pero que no podremos interpretar claramente con un concepto intuitivo.

```
In [7]: pca <- princomp(medidas, cor = T)
summary(pca)
```

```
Importance of components:
                    Comp.1   Comp.2   Comp.3
Comp.4
Standard deviation  1.7083611 0.9560494 0.38308860 0.143
926497
Proportion of Variance 0.7296245 0.2285076 0.03668922 0.005
178709
Cumulative Proportion 0.7296245 0.9581321 0.99482129 1.000
000000
```

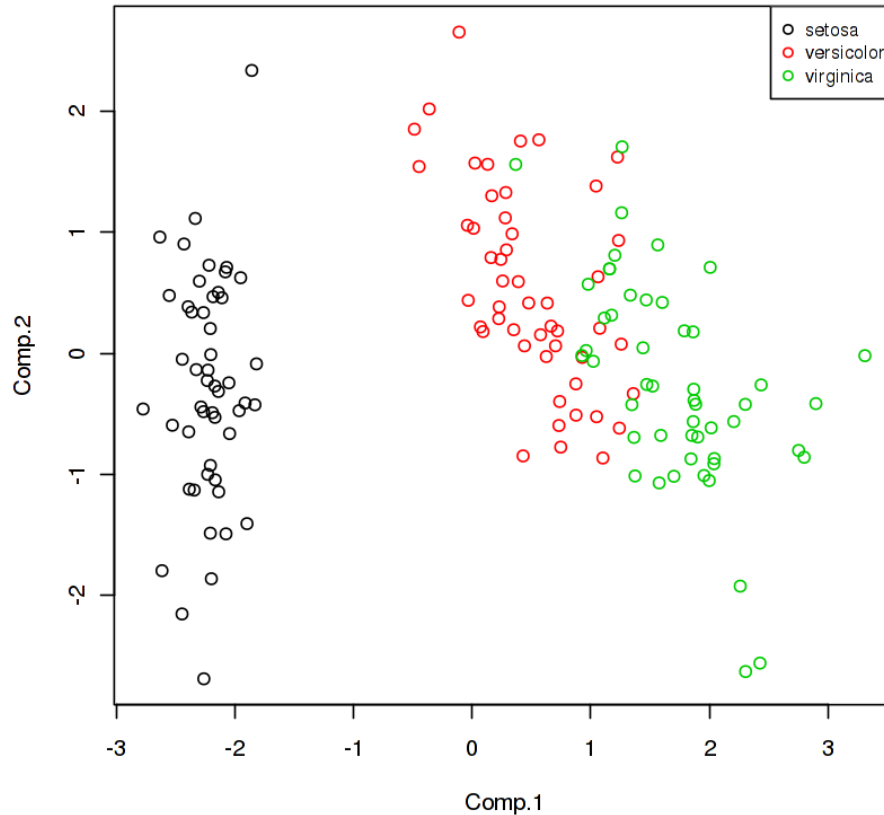
Podemos dibujar las varianzas, para visualizar la "importancia" de las componentes. Tras este paso, lo habitual es sólo quedarnos con las más explicativas (principales).

```
In [8]: plot(pca)
```



Si nos quedamos sólo con las componentes 1 y 2, podemos dibujar una nube de puntos con los especímenes (coloreados por especie) y observar que con estas nuevas medidas las especies se distinguen bien.

```
In [9]: plot(pca$scores, col=especies)
legend('topright', legend = levels(especies), col = 1:3, ce
x = 0.8, pch = 1)
```

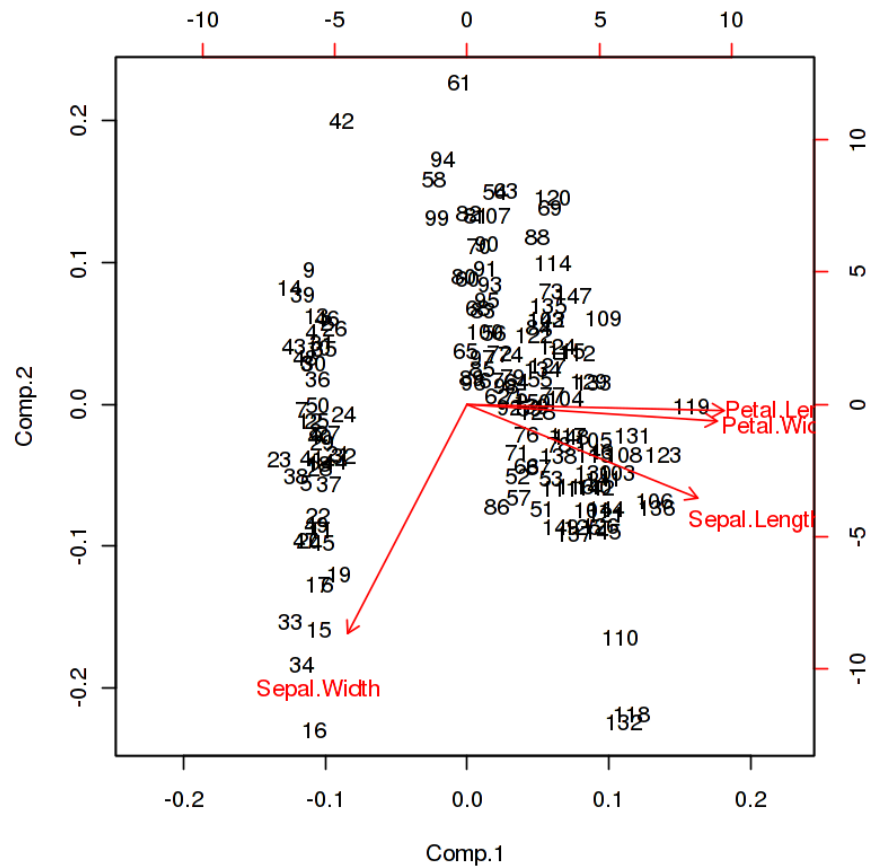


También podemos dibujar, en el espacio de las nuevas componentes, cómo varían las medidas originales.

En este caso, podemos ver justificada nuestra intuición de que el tamaño del pétalo es una componente importante, y que el ancho y la longitud dependen de ella. Esto se ve en que estas medidas son casi horizontales (es decir, paralelas a la componente 1), y coinciden en gran medida.

Sin embargo, ¿cuál sería la interpretación para la componente 2? Esto no es tan inmediato, y es posible que no haya una interpretación natural con mejor nombre que "segunda componente principal".

```
In [10]: biplot(pca)
```



Preguntas

- ¿Cómo describirías tú la componente 2?
- ¿Qué quiere decir que los especímenes de "setosa" formen una línea vertical?
- A la vista del diagrama de componentes, y de las nubes de puntos en las medidas dos a dos (segundo diagrama), ¿qué herramientas estadísticas usarías para decidir, dada una nueva flor de la que no conocemos la especie, a qué especie pertenece?