



T106: TAREA DE INVESTIGACIÓN: FUNCIONES IMPORTANTES DE R

Ing. Maximiliano Carsi Castrejón – Extracción y
Conocimiento en Bases de Datos

DESCRIPCIÓN BREVE

Este documento trata sobre la definición de R y
RStudio, así como su proceso de instalación

Luis Eduardo Bahena Castillo

9°C IDyGS



INTRODUCCION

Tarea de investigación: Funciones importantes de R

Objetivo:

- Identificar y describir las funciones más importantes de R para el análisis de datos.
- Explorar ejemplos prácticos de cómo usar estas funciones.
- Evaluar la utilidad de estas funciones para diferentes tareas.

Instrucciones:

1. Investiga:

- Funciones para la manipulación de datos:
 - `read.csv()`: Leer datos desde un archivo CSV.
 - `str()`: Visualizar la estructura de un objeto de datos.
 - `summary()`: Obtener un resumen estadístico de un objeto de datos.
 - `subset()`: Filtrar un conjunto de datos.
 - `mutate()`: Agregar nuevas variables a un conjunto de datos.
- Funciones para la visualización de datos:
 - `ggplot2`: Crear gráficos de alta calidad.
 - `plot()`: Visualizar variables de un conjunto de datos.
 - `barplot()`: Crear gráficos de barras.
 - `piechart()`: Crear gráficos circulares.
 - `boxplot()`: Crear diagramas de caja.
- Funciones para el análisis estadístico:
 - `t.test()`: Realizar una prueba t para dos muestras.
 - `aov()`: Realizar un análisis de varianza.
 - `lm()`: Ajustar un modelo de regresión lineal.
 - `cor.test()`: Calcular el coeficiente de correlación.
 - `chisq.test()`: Realizar una prueba de chi-cuadrado.

2. Para cada función:

- Define la función: Explica qué hace la función y qué argumentos recibe.
- Proporciona un ejemplo práctico: Muestra cómo usar la función con un conjunto de datos real.
- Explica la utilidad de la función: Describe en qué casos es útil la función y qué tipo de problemas permite solucionar.

3. Elabora un informe:

Contenido del informe:

- **Introducción:**
 - Describe brevemente el lenguaje R y su importancia en el análisis de datos.
 - Menciona las ventajas de usar R para el análisis de datos.
- **Desarrollo:**
 - Explica en detalle las funciones investigadas.
 - Proporciona ejemplos prácticos de cada función.
 - Explica la utilidad de cada función.
- **Conclusiones:**

- Resume los puntos clave sobre las funciones importantes de R.
- Explica por qué estas funciones son herramientas valiosas para el análisis de datos.

Recursos adicionales:

Sitio web oficial de R: <https://cran.r-project.org/>

Sitio web oficial de RStudio: <https://www.rstudio.com/products/rstudio/>

Tutoriales de R:

<https://www.r-bloggers.com/>

Libros sobre R:

- [The R Programming Language by Hadley Wickham](#)
- [R for Data Science by Hadley Wickham and Garrett Grolemund](#)

DESARROLLO

Lenguaje R:

R es un lenguaje de programación y un entorno de software especializado en estadísticas y análisis de datos. Se ha convertido en una herramienta fundamental para profesionales en campos como la ciencia de datos, la estadística, la investigación académica y la ingeniería. Su popularidad radica en su versatilidad, potencia y comunidad activa de usuarios y desarrolladores.

Ventajas de R para el Análisis de Datos:

- **Gran cantidad de paquetes:** R cuenta con una amplia variedad de paquetes que cubren prácticamente todas las áreas del análisis de datos, desde la importación de datos hasta la visualización y modelado estadístico.
- **Flexibilidad:** Permite realizar análisis estadísticos complejos y personalizados gracias a su capacidad para manipular datos de diversas formas y aplicar una amplia gama de técnicas analíticas.
- **Calidad gráfica:** Sus capacidades gráficas, especialmente a través de paquetes como ggplot2, permiten crear visualizaciones de alta calidad que son fácilmente personalizables y estéticamente agradables.
- **Comunidad activa:** La comunidad de usuarios de R es muy activa, lo que significa que siempre hay recursos disponibles, como tutoriales, foros de discusión y paquetes actualizados.
- **Open-source y gratuito:** R es un software de código abierto y gratuito, lo que lo hace accesible para cualquier persona interesada en el análisis de datos, sin importar su ubicación geográfica o sus recursos económicos.

Funciones para la manipulación de datos:

1. read.csv():

- **Definición:** La función read.csv() se utiliza para leer datos desde un archivo CSV (Comma Separated Values) y cargarlos en un objeto de datos en R.
- **Ejemplo:** `datos <- read.csv("datos.csv")`
- **Utilidad:** Es útil cuando se necesitan cargar datos almacenados en un archivo CSV para su análisis y manipulación en R.

2. str():

- **Definición:** La función str() es utilizada para visualizar la estructura de un objeto de datos en R, mostrando la estructura de los datos y el tipo de cada variable.
- **Ejemplo:** `str(datos)`
- **Utilidad:** Proporciona una descripción detallada de la estructura de los datos, lo que ayuda a comprender su composición y facilita la manipulación y análisis subsiguientes.

3. summary():

- **Definición:** La función summary() se utiliza para obtener un resumen estadístico de un objeto de datos en R, mostrando estadísticas descriptivas para cada variable.
- **Ejemplo:** `summary(datos)`

- **Utilidad:** Proporciona una visión general rápida de los datos, incluyendo medidas como la media, la mediana, los valores mínimo y máximo, entre otros, para cada variable.
4. **subset():**
- **Definición:** La función `subset()` se utiliza para filtrar un conjunto de datos en base a ciertos criterios, seleccionando solo las filas que cumplen con ciertas condiciones.
 - **Ejemplo:** `subset(datos <- subset(datos, edad > 30))`
 - **Utilidad:** Permite realizar selecciones específicas de datos según condiciones definidas, lo que facilita el análisis de subconjuntos de interés.
5. **mutate():**
- **Definición:** La función `mutate()` se utiliza para agregar nuevas variables calculadas a un conjunto de datos existente en R, utilizando operaciones basadas en las variables existentes.
 - **Ejemplo:** `datos <- mutate(datos, IMC = peso / (altura^2))`
 - **Utilidad:** Permite crear nuevas variables derivadas de las variables existentes, lo que puede ser útil para realizar análisis más avanzados o para preparar los datos para modelos estadísticos o de aprendizaje automático.

Funciones para la visualización de datos:

6. **ggplot2:**
- **Definición:** ggplot2 es un paquete en R utilizado para crear gráficos de alta calidad y personalizables, basados en la gramática de gráficos.
 - **Ejemplo:** `ggplot(data = iris, aes(x = Sepal.Length, y = Sepal.Width)) + geom_point()`
 - **Utilidad:** Permite crear una amplia variedad de gráficos, incluyendo gráficos de dispersión, líneas, barras, entre otros, con la capacidad de personalizar aspectos como colores, etiquetas y temas.
7. **plot():**
- **Definición:** La función `plot()` se utiliza para visualizar variables de un conjunto de datos en forma de gráficos básicos, como gráficos de dispersión, de líneas o de puntos.
 - **Ejemplo:** `plot(iris$Sepal.Length, iris$Sepal.Width)`
 - **Utilidad:** Proporciona una forma rápida y sencilla de visualizar relaciones entre variables o distribuciones de datos.
8. **barplot():**
- **Definición:** La función `barplot()` se utiliza para crear gráficos de barras que muestran la distribución de una variable categórica.
 - **Ejemplo:** `barplot(table(iris$Species))`
 - **Utilidad:** Es útil para comparar la frecuencia o proporción de diferentes categorías en un conjunto de datos.
9. **piechart():**
- **Definición:** Aunque no existe una función `piechart()` en R base, se pueden crear gráficos circulares utilizando funciones como `pie()` o ggplot2.
 - **Ejemplo:** `pie(table(iris$Species))`
 - **Utilidad:** Los gráficos circulares son útiles para mostrar la distribución de una variable categórica como partes de un todo.

10. boxplot():

- **Definición:** La función boxplot() se utiliza para crear diagramas de caja que muestran la distribución de una variable numérica o categórica a través de sus cuartiles.
- **Ejemplo:** boxplot(iris\$Sepal.Length)
- **Utilidad:** Es útil para visualizar la distribución, la dispersión y los valores atípicos de una variable numérica.

Funciones para el análisis estadístico:

11. t.test():

- **Definición:** La función t.test() se utiliza para realizar una prueba t para comparar las medias de dos muestras independientes.
- **Ejemplo:** t.test(x = datos1, y = datos2)
- **Utilidad:** Es útil para determinar si hay diferencias significativas entre las medias de dos grupos o poblaciones.

12. aov():

- **Definición:** La función aov() se utiliza para realizar un análisis de varianza (ANOVA) para comparar las medias de más de dos grupos.
- **Ejemplo:** aov(formula = y ~ factor(grupo), data = datos) summary(modelo)
- **Utilidad:** Es útil para determinar si hay diferencias significativas entre las medias de varios grupos y para identificar qué grupos difieren entre sí.

13. lm():

- **Definición:** La función lm() se utiliza para ajustar un modelo de regresión lineal a los datos.
- **Ejemplo:** modelo <- lm(formula = y ~ x, data = datos) summary(modelo)
- **Utilidad:** Es útil para modelar la relación entre una variable dependiente y una o más variables independientes, y para hacer predicciones basadas en esta relación.

14. cor.test():

- **Definición:** La función cor.test() se utiliza para calcular el coeficiente de correlación entre dos variables y realizar una prueba de hipótesis sobre la fuerza y dirección de la relación.
- **Ejemplo:** cor.test(x = datos\$variable1, y = datos\$variable2)
- **Utilidad:** Es útil para determinar si existe una relación lineal significativa entre dos variables y cuantificar la fuerza y dirección de esta relación.

15. chisq.test():

- **Definición:** La función chisq.test() se utiliza para realizar una prueba de chi-cuadrado de independencia entre dos variables categóricas.
- **Ejemplo:** chisq.test(tabla_de_contingencia)
- **Utilidad:** Es útil para determinar si existe una asociación significativa entre dos variables categóricas en una población, es decir, si son independientes o no.

CONCLUSIÓN

Las funciones proporcionadas para la manipulación de datos, la visualización de datos y el análisis estadístico en R, se puede concluir que R es un lenguaje de programación poderoso y versátil para la ciencia de datos y el análisis estadístico.

- Manipulación de datos: R ofrece una variedad de funciones como `read.csv()`, `subset()`, y `mutate()` que permiten cargar, filtrar y transformar conjuntos de datos de manera eficiente, facilitando la preparación de datos para análisis posteriores.
- Visualización de datos: Con herramientas como `ggplot2`, `plot()`, y `barplot()`, R proporciona capacidades avanzadas para crear visualizaciones claras y atractivas que ayudan a comprender y comunicar patrones y tendencias en los datos.
- Análisis estadístico: R ofrece una amplia gama de funciones estadísticas, como `t.test()`, `aov()`, `lm()`, `cor.test()`, y `chisq.test()`, que permiten realizar desde pruebas de hipótesis básicas hasta análisis de regresión complejos y pruebas de asociación entre variables categóricas.

En resumen, R es una herramienta integral que proporciona todas las capacidades necesarias para llevar a cabo un análisis de datos completo, desde la importación y manipulación de datos hasta la visualización y el análisis estadístico avanzado. Su flexibilidad, potencia y amplia comunidad de usuarios y desarrolladores lo convierten en una opción popular para profesionales en campos como la estadística, la ciencia de datos, la investigación y la academia.

