

# Análisis de Clústers, Componentes Principales y Mapas de Calor

Carlos E Martinez-Rodriguez

31 de mayo de 2023

## Índice

<b>1. Introducción</b>	<b>1</b>
1.1. Utilidad de los mapas de calor en la visualización de datos . . . . .	1
1.2. Construcción de mapas de calor . . . . .	2
1.3. Personalización de los mapas de calor . . . . .	2
1.4. Interpretación de los mapas de calor . . . . .	3
1.5. Preparación de datos . . . . .	3
1.6. Normalización de los datos . . . . .	4
1.7. Generación del mapa de calor . . . . .	4
1.7.1. Construcción de un mapa de calor básico con la librería heatmap . . . . .	4
1.7.2. Construcción de un mapa de calor con la librería ggplot2 . . . . .	4
<b>2. Análisis alternativo</b>	<b>5</b>
2.1. Teoría matemática de los mapas de calor . . . . .	5

## 1. Introducción

Un mapa de calor es una representación gráfica en la que los valores de una matriz se muestran mediante colores. Cada celda de la matriz se asigna a un color según su valor, lo que permite visualizar fácilmente las variaciones, distribución e intensidad de una variable en una matriz de datos. Los mapas de calor son especialmente útiles cuando se trabaja con matrices de datos grandes y multidimensionales y permiten visualizar patrones, tendencias y variaciones en los datos de manera efectiva.

La principal ventaja de los mapas de calor es su capacidad para resaltar las diferencias relativas entre los valores en una matriz. Al asignar colores a los valores, los mapas de calor permiten una fácil identificación de las regiones con valores altos o bajos en comparación con el resto de la matriz. Esto es especialmente útil cuando se trabaja con conjuntos de datos grandes y complejos, ya que los patrones pueden ser difíciles de detectar simplemente inspeccionando los valores numéricos.

Además, los mapas de calor permiten visualizar relaciones y patrones entre filas y columnas de la matriz. Al ordenar las filas y columnas según alguna métrica (por ejemplo, mediante agrupación jerárquica o clasificación), se pueden identificar agrupaciones y similitudes entre elementos. Esto facilita la identificación de clústeres, tendencias o relaciones entre diferentes variables o muestras.

### 1.1. Utilidad de los mapas de calor en la visualización de datos

Los mapas de calor son especialmente útiles cuando se trabaja con conjuntos de datos grandes y multidimensionales. Algunas de las aplicaciones comunes de los mapas de calor en la visualización de datos incluyen:

- Identificación de patrones y tendencias en datos climáticos y meteorológicos.
- Análisis de la distribución de la temperatura en imágenes térmicas.
- Visualización de datos de rendimiento en deportes para identificar fortalezas y debilidades.
- Análisis de expresión génica en biología molecular para identificar genes activos o inactivos.
- Representación de datos de mercado y finanzas para identificar áreas de crecimiento o declive.

## 1.2. Construcción de mapas de calor

Existen varias técnicas y herramientas para construir mapas de calor. En R, existen varias bibliotecas populares para generar mapas de calor. Algunas de las librerías más utilizadas son:

- **ggplot2**: Es una de las bibliotecas más populares y versátiles en R para visualización de datos. Si bien no es específicamente una biblioteca para mapas de calor, se puede utilizar para crear mapas de calor utilizando la función `geom_tile()` para representar los valores de los datos como colores en una cuadrícula.
- **heatmap**: Es una biblioteca base de R que proporciona funciones para generar mapas de calor. La función `heatmap()` crea un mapa de calor basado en una matriz de datos numéricos, donde los valores son representados por colores.
- **heatmaply**: Esta biblioteca utiliza la biblioteca `plotly` para generar mapas de calor interactivos en R. Proporciona funciones para crear mapas de calor 2D y 3D, y permite explorar los mapas de calor con características interactivas como zoom, rotación y selección de puntos.
- **pheatmap**: Es una biblioteca especializada en la generación de mapas de calor en R. Proporciona una amplia gama de opciones de personalización, como la configuración de colores, el ordenamiento de filas y columnas, y la adición de anotaciones a los mapas de calor.
- **ComplexHeatmap**: Es una biblioteca avanzada para la generación de mapas de calor complejos en R. Permite crear mapas de calor con múltiples capas de información, como anotaciones, dendrogramas y matrices de distancias.

## 1.3. Personalización de los mapas de calor

Los mapas de calor se pueden personalizar para adaptarse a las necesidades específicas de los datos y la visualización. Algunas de las opciones de personalización incluyen:

- Selección de colores: es posible elegir una paleta de colores adecuada para resaltar las variaciones en los datos. La elección de colores es esencial para la interpretación adecuada de los mapas de calor. R ofrece varias paletas de colores predefinidas, pero también es posible definir colores personalizados. A continuación se muestra un ejemplo de cómo seleccionar colores para un mapa de calor en R:
- Escalas de color: se pueden ajustar las escalas de color para resaltar áreas de alta o baja intensidad. La selección de escalas es un paso crítico en la construcción de mapas de calor. En general, existen dos tipos de escalas comúnmente utilizadas: escalas lineales y escalas logarítmicas. La elección de la escala depende de la naturaleza de los datos y del objetivo del análisis. Si los datos abarcan un rango amplio y se desea resaltar las diferencias en valores pequeños, es recomendable utilizar una escala logarítmica. Por otro lado, si los datos se distribuyen de manera más uniforme, una escala lineal puede ser más adecuada.

En este ejemplo, se utiliza la función `'colorRampPalette'` para definir una paleta de colores que va desde el azul hasta el rojo. Esta paleta personalizada se pasa como argumento a la función `'heatmap'` para construir el mapa de calor.

Además de los colores, es importante seleccionar la escala adecuada para los mapas de calor. R ofrece varias opciones de escalas, como escalas lineales y logarítmicas. A continuación se muestra un ejemplo de cómo seleccionar una escala para un mapa de calor en R:

```
# Cargar la librería pheatmap
library(pheatmap)

# Crear una matriz de ejemplo
matriz <- matrix(rnorm(100), nrow = 10)

# Construir el mapa de calor utilizando una escala logarítmica
pheatmap(matriz, scale = "log")
```

En este ejemplo, se utiliza la función `'pheatmap'` de la librería `pheatmap` para construir el mapa de calor. El argumento `'scale'` se establece en `"log"` para utilizar una escala logarítmica en el mapa de calor.

- Anotaciones y etiquetas: es posible agregar etiquetas y anotaciones a las celdas del mapa de calor para proporcionar información adicional. Las anotaciones pueden proporcionar información adicional en los mapas de calor, como etiquetas de fila y columna, valores numéricos o etiquetas personalizadas. A continuación se muestra un ejemplo de cómo agregar anotaciones a un mapa de calor en R:

```
# Cargar la librería ComplexHeatmap
library(ComplexHeatmap)

# Crear una matriz de ejemplo
matriz <- matrix(rnorm(100), nrow = 10)

# Construir el mapa de calor con anotaciones de fila y columna
Heatmap(matriz, name = "Valor", row_names_side = "left", column_names_side = "top")
```

En este ejemplo, se utiliza la función ‘Heatmap’ de la librería ComplexHeatmap para construir el mapa de calor. Los argumentos ‘row\_names\_side’ y ‘column\_names\_side’ se establecen en “left” y “top”, respectivamente, para agregar etiquetas de fila a la izquierda y etiquetas de columna en la parte superior del mapa de calor.

```
# Cargar la librería heatmap
library(heatmap)

# Crear una matriz de ejemplo
matriz <- matrix(rnorm(100), nrow = 10)

# Definir una paleta de colores personalizada
mi_paleta <- colorRampPalette(c("blue", "white", "red"))

# Construir el mapa de calor utilizando la paleta de colores personalizada
heatmap(matriz, col = mi_paleta)
```

## 1.4. Interpretación de los mapas de calor

La interpretación de los mapas de calor es crucial para comprender la distribución de los datos. La interpretación de los mapas de calor se basa en la capacidad para percibir y distinguir diferentes colores. Algunos aspectos a considerar al interpretar los mapas de calor son:

- Atención a la intensidad de los colores: los colores más intensos indican valores más altos o significativos, mientras que los colores más claros indican valores más bajos. Al utilizar una paleta de colores bien elegida, se pueden resaltar las diferencias y patrones de manera efectiva. Por ejemplo, se puede utilizar una paleta de colores divergente que resalte los valores extremos o una paleta de colores secuencial que muestre una progresión gradual de valores. La elección de la paleta de colores adecuada es crucial para garantizar una interpretación precisa y no sesgada de los datos.
- Identificación de patrones y tendencias: buscar áreas de concentración o dispersión en el mapa de calor para identificar patrones y tendencias en los datos.
- Comparación entre mapas de calor: comparar diferentes mapas de calor para identificar diferencias o similitudes en la distribución de la variable en diferentes conjuntos de datos.

## 1.5. Preparación de datos

Antes de construir un mapa de calor en R, es importante preparar los datos adecuadamente. Es preciso asegurarse de tener una matriz o un dataframe con los datos que desees visualizar, además de que los datos estén en el formato correcto y que no falten valores. Preparación de los datos: Antes de construir un mapa de calor, es necesario realizar una preparación adecuada de los datos. Esto implica asegurarse de que los datos estén en el formato correcto y sean apropiados para su visualización en un mapa de calor. Si los datos están en forma de una matriz, es importante verificar que estén completos, es decir, no deben haber valores faltantes. En caso de que existan valores faltantes, es necesario tomar decisiones sobre cómo manejarlos, ya sea eliminando las filas o columnas correspondientes o imputando valores apropiados.

## 1.6. Normalización de los datos

Antes de construir un mapa de calor, es común aplicar técnicas de normalización a los datos. La normalización tiene como objetivo ajustar los valores de los datos a una escala específica y facilitar la comparación entre diferentes elementos. En el caso de los mapas de calor, se utilizan dos enfoques de normalización comunes: normalización por filas y normalización por columnas.

En la normalización por filas, los valores de cada fila se escalan para que tengan una suma de 1 o una media de 0. Esto permite comparar la contribución relativa de cada valor dentro de una fila y resalta los patrones de distribución entre filas.

En la normalización por columnas, los valores de cada columna se escalan para que tengan una suma de 1 o una media de 0. Esto permite comparar la contribución relativa de cada valor dentro de una columna y resalta los patrones de distribución entre columnas.

## 1.7. Generación del mapa de calor

Generación del mapa de calor: Una vez que los datos han sido preparados y normalizados, se procede a la construcción del mapa de calor. Esto implica asignar colores a los valores de la matriz normalizada y representarlos visualmente en una cuadrícula. La elección de la paleta de colores es importante para una interpretación adecuada del mapa de calor. Se deben seleccionar colores que sean perceptualmente distintos y que permitan una fácil identificación de los valores altos y bajos.

### 1.7.1. Construcción de un mapa de calor básico con la librería heatmap

A continuación se muestra un ejemplo de cómo construir un mapa de calor básico utilizando la librería heatmap en R:

```
# Cargar la librería heatmap
library(heatmap)

# Crear una matriz de ejemplo
matriz <- matrix(rnorm(100), nrow = 10)

# Construir el mapa de calor
heatmap(matriz)
```

Este código generará un mapa de calor básico utilizando la función `heatmap` de la librería heatmap. Es posible ajustar los parámetros de la función según las necesidades, como los colores utilizados, la escala, entre otros.

### 1.7.2. Construcción de un mapa de calor con la librería ggplot2

Para construir un mapa de calor más personalizable y estéticamente agradable, puedes utilizar la librería ggplot2. A continuación se muestra un ejemplo de cómo construir un mapa de calor utilizando la librería ggplot2:

```
# Cargar la librería ggplot2
library(ggplot2)

# Crear una matriz de ejemplo
matriz <- matrix(rnorm(100), nrow = 10)

# Convertir la matriz a un dataframe
df <- as.data.frame(matriz)

# Construir el mapa de calor utilizando ggplot2
ggplot(df, aes(x = factor(1), y = factor(1), fill = V1)) +
  geom_tile() +
  scale_fill_gradient(low = "blue", high = "red") +
  theme_void()
```

Este código generará un mapa de calor utilizando la función `ggplot` de la librería ggplot2. Puedes personalizar el mapa de calor ajustando los parámetros de la función, como los colores utilizados, la escala, entre otros.

## 2. Análisis alternativo

El análisis de clusters es una técnica fundamental en Bioinformática que permite agrupar objetos similares en conjuntos homogéneos. En el contexto de la investigación genómica y proteómica, el análisis de clusters se utiliza para identificar patrones, descubrir relaciones entre genes o proteínas, y clasificar muestras en función de su perfil molecular. Este enfoque juega un papel crucial en la comprensión de la estructura y función de los datos biológicos, así como en la identificación de biomarcadores y el descubrimiento de nuevas terapias.

El análisis de clusters consiste en agrupar objetos similares en conjuntos homogéneos, donde la similitud se basa en medidas específicas según el tipo de datos (por ejemplo, distancias genéticas, correlaciones de expresión génica o similitud de secuencias). El análisis de clusters es esencial en Bioinformática por varias razones:

1. Descubrimiento de patrones y relaciones: Permite identificar grupos de genes o proteínas con perfiles de expresión similares, lo que sugiere su función común en procesos biológicos específicos. Además, puede ayudar a descubrir relaciones entre diferentes muestras, como en el caso de la clasificación de pacientes en subtipos de enfermedades o en la identificación de especies relacionadas.
2. Identificación de biomarcadores: Los biomarcadores son características moleculares que se asocian con enfermedades, estados patológicos o respuestas a tratamientos. El análisis de clusters puede revelar patrones moleculares distintivos en diferentes grupos de pacientes, lo que ayuda a identificar biomarcadores potenciales para el diagnóstico, pronóstico y selección de terapias personalizadas.
3. Visualización de datos: Los resultados del análisis de clusters se pueden representar gráficamente mediante dendrogramas, mapas de calor u otras visualizaciones, lo que facilita la interpretación y comunicación de los patrones encontrados. Estas representaciones visuales permiten una comprensión más intuitiva de los datos y pueden guiar investigaciones posteriores.
4. Agrupamiento de secuencias biológicas: En el campo de la genómica y la proteómica, el análisis de clusters es esencial para agrupar secuencias de ADN, ARN o proteínas similares. Esto facilita la identificación de familias de genes o proteínas relacionadas, así como la predicción de funciones basadas en la similitud de secuencias.
5. Exploración de datos de alto rendimiento: Con la creciente disponibilidad de conjuntos de datos de alto rendimiento, como los datos de secuenciación masiva o los perfiles de expresión a gran escala, el análisis de clusters se ha convertido en una herramienta imprescindible. Permite analizar grandes volúmenes de datos de manera eficiente y extraer información valiosa de ellos.

### 2.1. Teoría matemática de los mapas de calor

Los mapas de calor son una representación visual utilizada para mostrar la distribución y la intensidad de los valores en una matriz de datos. Esta representación se basa en la asignación de colores a los valores numéricos, donde los colores más claros indican valores altos y los colores más oscuros indican valores bajos. La teoría matemática detrás de los mapas de calor involucra los siguientes conceptos:

- **Normalización:** Antes de generar un mapa de calor, es común realizar una normalización de los datos. La normalización tiene como objetivo escalar los valores de la matriz de datos a un rango común, lo que permite una comparación más precisa entre los valores. Una técnica común de normalización es la normalización por filas o por columnas, donde se calcula el valor relativo de cada dato con respecto a los demás en la misma fila o columna.
- **Normalización por filas:** Dada una matriz de datos  $X$  de tamaño  $m \times n$ , la normalización por filas se realiza calculando el valor relativo de cada dato con respecto a los demás en la misma fila. Para cada elemento  $x_{ij}$  de la matriz, la normalización por filas se calcula de la siguiente manera:

$$x'_{ij} = \frac{x_{ij}}{\sqrt{\sum_{k=1}^n x_{ik}^2}}$$

La normalización por renglones se realiza calculando el valor relativo de cada dato con respecto a los demás valores en la misma fila. Este enfoque permite resaltar los patrones de variación dentro de cada renglón. El proceso de normalización por renglones se puede describir de la siguiente manera:

$$\text{Dato normalizado por renglones}_{ij} = \frac{\text{Dato original}_{ij}}{\sum_{j=1}^n \text{Dato original}_{ij}} \quad (1)$$

donde  $\text{Dato original}_{ij}$  representa el elemento en la fila  $i$  y columna  $j$  de la matriz de datos original, y  $n$  es el número total de columnas.

- **Normalización por columnas:** Dada una matriz de datos  $X$  de tamaño  $m \times n$ , la normalización por columnas se realiza calculando el valor relativo de cada dato con respecto a los demás en la misma columna. Para cada elemento  $x_{ij}$  de la matriz, la normalización por columnas se calcula de la siguiente manera:

$$x'_{ij} = \frac{x_{ij}}{\sqrt{\sum_{k=1}^m x_{kj}^2}}$$

En ambas fórmulas,  $x'_{ij}$  representa el valor normalizado del elemento  $x_{ij}$ ,  $\sum_{k=1}^n x_{ik}^2$  es la suma de los cuadrados de los elementos en la misma fila y  $\sum_{k=1}^m x_{kj}^2$  es la suma de los cuadrados de los elementos en la misma columna. La normalización por filas y por columnas asegura que los valores en cada fila o columna tengan una magnitud comparable y permite una comparación más precisa entre los valores de la matriz de datos.

La normalización por columnas se realiza calculando el valor relativo de cada dato con respecto a los demás valores en la misma columna. Este enfoque permite resaltar los patrones de variación dentro de cada columna. El proceso de normalización por columnas se puede describir de la siguiente manera:

$$\text{Dato normalizado por columnas}_{ij} = \frac{\text{Dato original}_{ij}}{\sum_{i=1}^m \text{Dato original}_{ij}} \quad (2)$$

donde  $\text{Dato original}_{ij}$  representa el elemento en la fila  $i$  y columna  $j$  de la matriz de datos original, y  $m$  es el número total de filas.

La normalización de datos es esencial para asegurar que los mapas de calor reflejen de manera precisa los patrones y relaciones presentes en los datos originales. Los mapas de calor normalizados permiten una visualización más efectiva y facilitan la identificación de las regiones con mayores y menores valores en la matriz de datos.

Adicionalmente, es común aplicar una transformación logarítmica a los datos antes de la normalización para resaltar mejor las diferencias en valores pequeños y grandes. Una transformación común es el  $\log_2$  debido a su interpretación en términos de cambios relativos y la frecuente aplicación en análisis genómicos. La transformación de los datos por  $\log_2$  se puede describir de la siguiente manera:

$$\text{Dato transformado}_{ij} = \log_2(\text{Dato original}_{ij}) \quad (3)$$

donde  $\text{Dato transformado}_{ij}$  representa el elemento transformado en la fila  $i$  y columna  $j$  de la matriz de datos original.

- **Colormap:** La colormap es una función que asigna valores numéricos a colores. Define la correspondencia entre los valores de los datos y los colores que se mostrarán en el mapa de calor. Existen diversas colormaps disponibles, como el mapa de colores *jet*, *hot*, *cool*, entre otros. Cada colormap tiene una gama de colores que van desde un color inicial (por ejemplo, azul) hasta un color final (por ejemplo, rojo) que representan los valores mínimos y máximos de la escala.

La función de asignación en una colormap define la correspondencia entre los valores numéricos de los datos y los colores que se mostrarán en el mapa de calor. Dada una colormap, se utiliza una función que mapea un valor numérico a un color específico en la escala de colores. La función de asignación toma como entrada un valor numérico y devuelve el color correspondiente según la escala establecida en la colormap.

En una colormap, los valores numéricos de los datos se asignan a colores específicos siguiendo una escala de colores predefinida. Por ejemplo, en la colormap *jet*, los valores mínimos se representan con colores más fríos como el azul, mientras que los valores máximos se representan con colores más cálidos como el rojo. La función de asignación mapea los valores numéricos en la escala de los valores mínimos y máximos definidos en la colormap, asignando colores correspondientes a cada valor.

Cabe mencionar que existen diversas colormaps disponibles, cada una con una gama de colores específica y una asignación única de valores numéricos a colores. Algunas colormaps comunes incluyen *jet*, *hot*, *cool*, entre otras. La elección de la colormap adecuada depende del tipo de datos y la representación visual deseada.

$$\text{color} = f(\text{valor}) \quad (4)$$

Donde  $\text{color}$  representa el valor de color asignado y "valor" es el valor numérico al que se le asigna el color. La función  $f$  define la correspondencia específica entre el valor numérico y el color correspondiente en la escala de colores. Esta función puede variar dependiendo de la colormap utilizada y cómo se defina la escala de colores.

- **Asignación lineal:**

$$\text{color} = m \cdot \text{valor} + b \quad (5)$$

Esta función asigna colores de manera lineal, donde  $m$  y  $b$  son constantes que determinan la pendiente y el desplazamiento de la línea.

- **Asignación logarítmica:**

$$\text{color} = a \cdot \log(\text{valor}) + b \quad (6)$$

Esta función asigna colores utilizando una escala logarítmica, donde  $a$  y  $b$  son constantes que afectan la amplitud y el desplazamiento de la función logarítmica.

- **Asignación exponencial:**

$$\text{color} = a \cdot \exp(b \cdot \text{valor}) \quad (7)$$

Esta función asigna colores utilizando una escala exponencial, donde  $a$  y  $b$  son constantes que controlan la amplitud y la tasa de crecimiento exponencial.

- **Interpolación:** La interpolación se utiliza para asignar colores a los valores que no están representados explícitamente en la matriz de datos. La interpolación se realiza mediante algoritmos que estiman los valores entre los puntos de datos conocidos. Por ejemplo, si hay una celda vacía en la matriz de datos, se puede utilizar la interpolación para estimar el valor correspondiente y asignarle un color en el mapa de calor.

La interpolación es una técnica utilizada para estimar valores intermedios entre puntos de datos conocidos. En el contexto de asignar colores a los valores en un mapa de calor, la interpolación se utiliza cuando hay celdas vacías o valores faltantes en la matriz de datos. La interpolación nos permite estimar esos valores faltantes y asignarles un color correspondiente.

Una forma común de realizar la interpolación en este caso es utilizar interpolación lineal. Supongamos que tenemos un conjunto de puntos de datos conocidos  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , donde  $x_i$  representa la posición en el eje  $x$  y  $y_i$  es el valor asociado en ese punto. Queremos estimar el valor  $y$  en un punto intermedio  $x$  que no está explícitamente representado en la matriz de datos.

La interpolación lineal se basa en la idea de que el cambio en los valores  $y$  es proporcional al cambio en los valores  $x$  entre los puntos conocidos. La fórmula de interpolación lineal se puede expresar matemáticamente como:

$$y = y_i + \frac{(x - x_i)}{(x_{i+1} - x_i)} \cdot (y_{i+1} - y_i) \quad (8)$$

donde  $x_i$  y  $x_{i+1}$  son los puntos de datos conocidos más cercanos al punto intermedio  $x$ , y  $y_i$  y  $y_{i+1}$  son los valores correspondientes asociados a esos puntos.

Utilizando esta fórmula de interpolación lineal, podemos estimar el valor faltante en la matriz de datos y asignarle un color en el mapa de calor. La interpolación nos permite obtener una representación visual más completa de los datos, rellenando las celdas vacías con valores estimados basados en la información disponible en los puntos conocidos.

- **Interpolación polinómica:** La interpolación polinómica es un método que utiliza polinomios para ajustar los puntos de datos conocidos y estimar los valores intermedios. Uno de los métodos más comunes es la interpolación de Lagrange. Dado un conjunto de puntos  $(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)$ , donde  $x_i$  son las abscisas y  $y_i$  son las ordenadas de los puntos, el polinomio de Lagrange se define como:

$$P(x) = \sum_{i=0}^n y_i \cdot \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j} \quad (9)$$

Este polinomio pasa exactamente por todos los puntos conocidos y se utiliza para estimar los valores intermedios.

- Interpolación spline:

La interpolación spline utiliza funciones suaves llamadas splines para estimar los valores faltantes. Un spline es una curva definida en segmentos que se ajusta a los puntos de datos conocidos. En la interpolación spline cúbica, se utilizan polinomios de tercer grado para cada segmento. El spline cúbico se puede representar mediante ecuaciones polinómicas de la forma:

$$S_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3 \quad \text{para } x_i \leq x \leq x_{i+1} \quad (10)$$

donde  $a_i$ ,  $b_i$ ,  $c_i$  y  $d_i$  son coeficientes que se determinan utilizando condiciones de continuidad y suavidad en los puntos conocidos.

- Interpolación vecino más cercano:

La interpolación vecino más cercano asigna el valor del punto de datos conocido más cercano al punto intermedio sin realizar cálculos adicionales. En esta técnica, el color asignado al punto intermedio será el mismo color que el punto conocido más cercano. Matemáticamente, se puede representar como:

$$f(x) = f(x_{\text{vecino más cercano}}) \quad (11)$$

donde  $f(x)$  es el valor estimado para el punto intermedio y  $x_{\text{vecino más cercano}}$  es el punto de datos conocido más cercano a  $x$ .

- Interpolación Kriging:

El Kriging es un método avanzado que considera la estructura espacial y la correlación de los datos para realizar la interpolación. Se basa en modelos geoestadísticos y se utiliza principalmente en análisis espacial y geología. La estimación Kriging se calcula como una combinación lineal ponderada de los valores conocidos en función de su distancia y correlación espacial. La ecuación general del Kriging es:

$$Z(x) = \sum_{i=1}^n \lambda_i Z(x_i) \quad (12)$$

donde  $Z(x)$  es el valor estimado para el punto intermedio,  $Z(x_i)$  son los valores conocidos,  $\lambda_i$  son los pesos asignados a cada valor conocido y  $n$  es el número de puntos conocidos.

Estas son algunas de las funciones de interpolación utilizadas en la generación de mapas de calor.

- **Escalado:** El escalado se utiliza para ajustar la intensidad de los colores en el mapa de calor. Permite resaltar o atenuar la diferencia de intensidad entre los valores para una mejor visualización. Se pueden utilizar diferentes métodos de escalado, como el escalado lineal, el escalado logarítmico o el escalado basado en percentiles, dependiendo de las características de los datos y el propósito de la representación.

- Escalado lineal:

El escalado lineal es el método más simple y directo de escalado. Consiste en ajustar los valores de los datos a un rango específico, generalmente entre 0 y 1 o entre un valor mínimo y máximo predefinidos. La fórmula para el escalado lineal es:

$$\text{valor escalado} = \frac{\text{valor original} - \text{mínimo}}{\text{máximo} - \text{mínimo}} \quad (13)$$

El escalado lineal preserva la proporción de las diferencias entre los valores originales, pero puede verse afectado por valores atípicos y puede no resaltar adecuadamente las diferencias entre valores cercanos.



- Escalado logarítmico:

El escalado logarítmico se utiliza cuando los datos abarcan un rango muy amplio y se desea resaltar las diferencias en los valores más pequeños. El escalado logarítmico aplica una transformación logarítmica a los datos antes de realizar el escalado lineal. La fórmula para el escalado logarítmico es:

$$\text{valor escalado} = \frac{\log(\text{valor original} - \text{mínimo} + 1)}{\log(\text{máximo} - \text{mínimo} + 1)} \quad (14)$$

El escalado logarítmico comprime los valores más grandes y expande los valores más pequeños, lo que puede revelar mejor las diferencias en la parte inferior de la escala. Sin embargo, es importante tener en cuenta que no se puede aplicar a valores negativos o cero.

- Escalado basado en percentiles:

El escalado basado en percentiles se utiliza cuando se desea resaltar las diferencias en ciertos rangos de valores. En este método, los valores se escalan en función de sus posiciones relativas dentro de la distribución. Por ejemplo, se puede asignar el percentil 5 al valor mínimo y el percentil 95 al valor máximo. La fórmula para el escalado basado en percentiles es:

$$\text{valor escalado} = \frac{\text{rango percentil del valor original} - \text{percentil mínimo}}{\text{percentil máximo} - \text{percentil mínimo}} \quad (15)$$

El escalado basado en percentiles resalta las diferencias en ciertos rangos de valores, lo que puede ser útil para destacar valores atípicos o concentraciones de datos en áreas específicas.

- **Eliminación de datos faltantes:** En el análisis de datos, es común enfrentarse a valores faltantes o datos ausentes en el conjunto de datos. Estos valores pueden ser el resultado de diversas circunstancias, como errores de medición, pérdida de datos durante la recolección o problemas técnicos.

En el caso específico de los mapas de calor, la presencia de datos faltantes puede afectar la precisión y la interpretación de los resultados. A continuación, se presenta una justificación para eliminar los datos faltantes antes de generar mapas de calor en R:

- Preservar la coherencia visual: Los mapas de calor se utilizan para visualizar patrones y relaciones en los datos, y se basan en la comparación de los valores numéricos. Cuando se eliminan los datos faltantes, se asegura que los patrones y las relaciones sean consistentes y no se distorsionen debido a los valores ausentes. Esto permite una interpretación más precisa de los resultados.
- Facilitar el análisis y la interpretación: Los datos faltantes pueden introducir incertidumbre y complejidad en el análisis de los datos. Al eliminar los datos faltantes, se simplifica el conjunto de datos y se facilita su análisis y la interpretación de los resultados. Los patrones emergentes se vuelven más evidentes y se reduce la posibilidad de conclusiones erróneas debido a los valores ausentes.
- Evitar sesgos en los resultados: La presencia de datos faltantes puede introducir sesgos en el análisis de datos y en los resultados obtenidos. Dependiendo de la naturaleza de los datos faltantes, pueden existir sesgos sistemáticos en ciertas variables o subconjuntos de datos. Al eliminar los datos faltantes, se minimiza la posibilidad de sesgos y se promueve la imparcialidad en el análisis.

Sin embargo, es importante destacar que la eliminación de datos faltantes debe realizarse con cuidado y considerando el contexto específico del análisis. Algunas consideraciones adicionales son:

- Evaluar el patrón de los datos faltantes: Es importante examinar si los datos faltantes siguen un patrón específico, como datos faltantes aleatorios o datos faltantes sistemáticos. Esto puede influir en la decisión de eliminar o imputar los valores faltantes.
- Evaluar el impacto en el tamaño y la representatividad de la muestra: Eliminar los datos faltantes puede reducir el tamaño de la muestra y afectar la representatividad de los datos. Es fundamental evaluar si la eliminación de los datos faltantes introduce algún sesgo en la muestra y si la muestra resultante sigue siendo adecuada para el análisis.
- Considerar técnicas de imputación de datos: En algunos casos, puede ser apropiado utilizar técnicas de imputación de datos para estimar los valores faltantes en lugar de eliminarlos por completo. Estas técnicas permiten preservar el tamaño de la muestra y proporcionar una aproximación de los valores ausentes. Sin embargo, es importante evaluar la idoneidad de las técnicas de imputación y los supuestos asociados.

- **Generación:** Una vez que se ha realizado la normalización de los datos, se ha definido la colormap y se han aplicado las técnicas de interpolación y escalado, se procede a generar el mapa de calor. Esto implica asignar colores a cada valor en la matriz de datos según la colormap y visualizar la matriz resultante en forma de una imagen o gráfico.

- Normalización de los datos:

La normalización de los datos es un paso importante antes de generar el mapa de calor. El objetivo es escalar los valores de la matriz de datos a un rango común para permitir una comparación más precisa entre los valores. Una técnica común de normalización es la normalización por filas o por columnas. La fórmula para la normalización por filas es:

$$\text{valor normalizado} = \frac{\text{valor original} - \text{valor mínimo en la fila/columna}}{\text{valor máximo en la fila/columna} - \text{valor mínimo en la fila/columna}} \quad (16)$$

La normalización por columnas se realiza de manera similar, pero utilizando los valores mínimo y máximo en la columna correspondiente.

- Definición de la colormap:

La colormap es una función que asigna valores numéricos a colores. Define la correspondencia entre los valores de los datos y los colores que se mostrarán en el mapa de calor. Existen diversas colormaps disponibles, como el mapa de colores jet, hot, cool, entre otros. Cada colormap tiene una gama de colores que representa los valores mínimos y máximos de la escala. La asignación de valores a colores puede realizarse mediante una función específica que mapea los valores al espacio de color.

- Aplicación de técnicas de interpolación y escalado:

Después de la normalización y la definición de la colormap, se aplican técnicas de interpolación y escalado para asignar colores a cada valor en la matriz de datos. La interpolación se utiliza para asignar colores a los valores que no están representados explícitamente en la matriz de datos. Se utilizan algoritmos que estiman los valores entre los puntos de datos conocidos. Por ejemplo, si hay una celda vacía en la matriz de datos, se puede utilizar la interpolación para estimar el valor correspondiente y asignarle un color en el mapa de calor.

El escalado se utiliza para ajustar la intensidad de los colores en el mapa de calor, resaltando o atenuando la diferencia de intensidad entre los valores para una mejor visualización. Se pueden utilizar diferentes métodos de escalado, como el escalado lineal, el escalado logarítmico o el escalado basado en percentiles, dependiendo de las características de los datos y el propósito de la representación.

- Generación del mapa de calor:

Una vez que se ha realizado la normalización, se ha definido la colormap y se han aplicado las técnicas de interpolación y escalado, se procede a generar el mapa de calor. Esto implica asignar colores a cada valor en la matriz de datos según la colormap y visualizar la matriz resultante en forma de una imagen o gráfico.

Estos son los pasos principales involucrados en la generación de un mapa de calor.

Es importante destacar que los mapas de calor son una representación visual que facilita la interpretación y comprensión de los patrones y tendencias en los datos. Sin embargo, es necesario tener en cuenta que la elección de la colormap, la normalización y el escalado pueden influir en la percepción y la interpretación de los datos representados en el mapa de calor.