

---

# pre TP2: Agrupamiento de imágenes

Data Mining en Ciencia y Tecnología

---

Juan Gabriel Juara  
Universidad de Buenos Aires (estudiante)

jgjuara@gmail.com

## 1 4.1. Carga de datos y verificación

2 El dataset utilizado consiste en 210 imágenes en formato png. Se analizó el tamaño de cada  
3 imagen hallando que una de las imágenes (0208.png) no correspondía al tamaño estándar  
4 del dataset: dimensiones 208 x 208, en vez de 128 x 128. Dicha imagen fue achicada usando  
5 el metodo `cv2.resize` con el método de interpolación `INTER_AREA` que ajusta el valor del  
6 nuevo pixel como promedio del area original.

7 Luego de este procesamiento inicial se analizó la distribución de valores de cada canal para  
8 todas las imagenes. En la figura 1 se verifica que todos los valores para los tres canales de  
9 las imágenes se distribuyen entre 0 y 255.

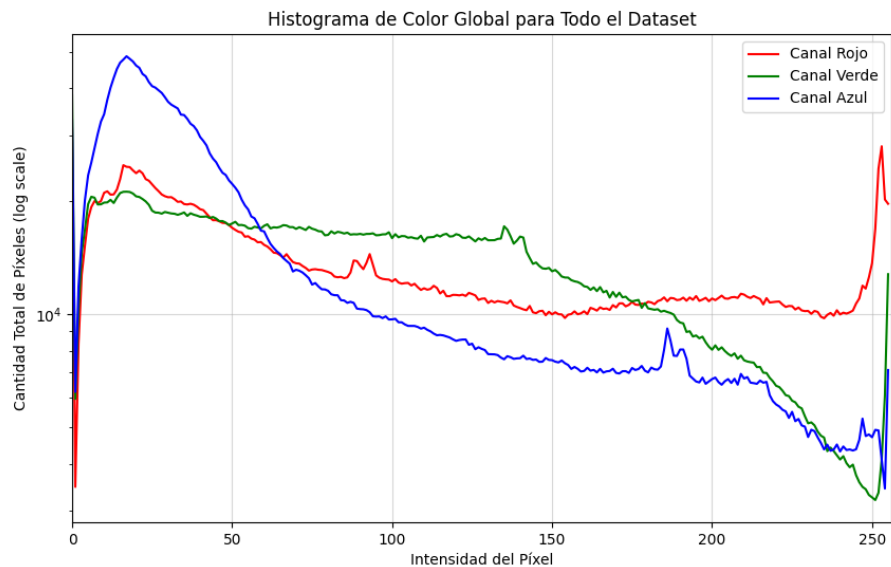


Figure 1: Figura 1. Histograma global de colores para el dataset completo

10 Figura 1. Histograma global de colores para el dataset completo

## 11 4.2. Exploración de subconjuntos por especie

12 El dataset de las imágenes contiene las etiquetas correspondientes a cada una. Se visualiza  
13 una imagen por especie a fines exploratorios de entender la representación visual de cada  
14 especie en la figura 2.

Muestra de una imagen por cada especie de flor



Figure 2: Figura 2. Grilla de especies

- 15 Figura 2. Grilla de imágenes de cada especie
- 16 5.1. Conversión a escala de grises y binarización
- 17 La figura 3 muestra el resultado de conversión a escala de grises de la imagen 0001.png.

Conversión a Escala de Grises

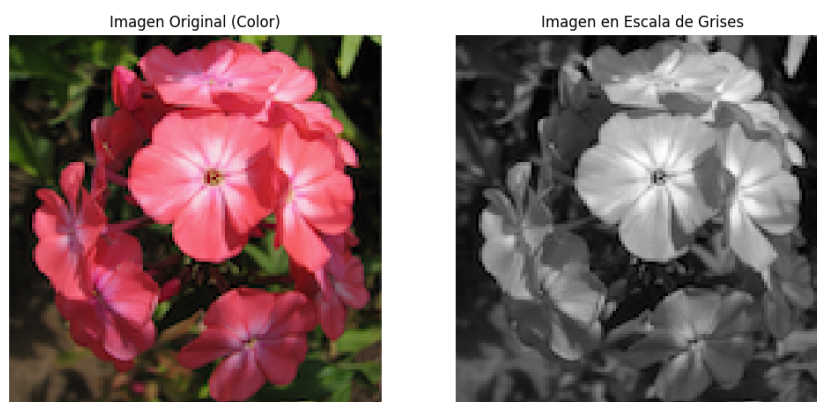


Figure 3: Figura 3. Conversión a grises

- 18 La conversión se ha realizado con el método `COLOR_BGR2GRAY` que utiliza la fórmula:

19 
$$Y = 0.299 \cdot R + 0.587 \cdot G + 0.114 \cdot B$$

20 Esta fórmula permite calcular el valor de luminancia a partir de los valores de los canales  
 21 de rojo, azul y verde de modo que sea adecuadamente percibido en escala de grises por el  
 22 ojo humano.

23 En la figura 4 se muestra el resultado de 4 operaciones de binarización a diferentes valores  
 24 umbrales. El método de binarización utilizado mapea los valores de luminancia de la imagen  
 25 en escala de grises a 0 o el máximo de la escala utilizada (255 en nuestro caso), según si el  
 26 valor del pixel menor o mayor al valor umbral.

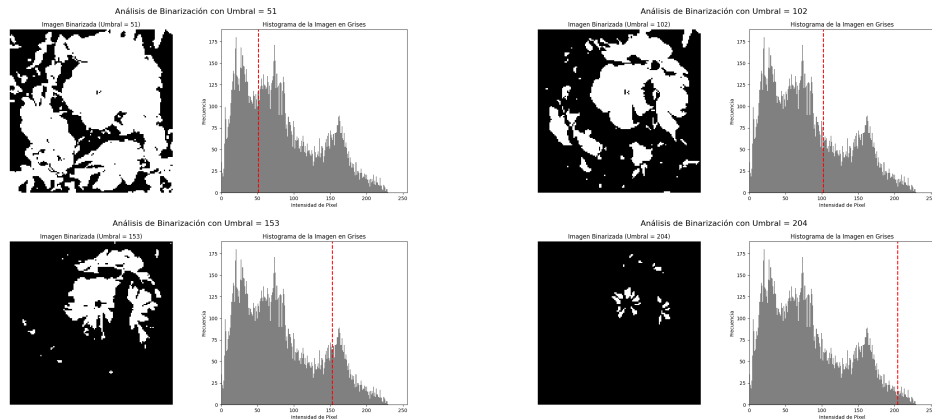


Figure 4: Figura 4. Binarización por 4 valores umbrales

27 En la figura 4 se acompaña el resultado de cada binarización con el histograma de la lu-  
 28 minancia de la imagen en escala de grises y se señala con una línea roja vertical el punto  
 29 de corte del valor umbral usado. Se puede apreciar como a mayores valores umbrales hay  
 30 mayor cantidad de píxeles mapeados a cero luminancia y viceversa a menor umbral, hay  
 31 mayor cantidad mapeada a máxima luminancia. Cabe destacar que entre los valores de  
 32 luminancia 100 y 150 existe un pozo en la distribución de píxeles que parece separar dos  
 33 máximos uno global y otro local de la distribución a izquierda y derecha respectivamente.

## 34 5.2. Generación de imágenes aleatorias

35 A partir de la imagen 0001.png se sintetizó una nueva imagen, para ello se tomaron los píxeles  
 36 en RGB, se aplanó la matriz de píxeles y se remuestrearon aleatoriamente sin reemplazo. El  
 37 resultado se visualiza en la figura 5.

## Generación de Imagen con Píxeles Mezclados Aleatoriamente

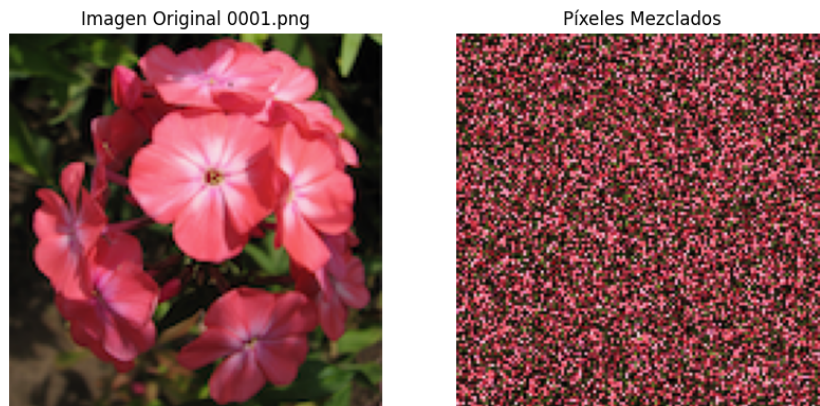


Figure 5: Figura 5. Aleatorización de píxeles de una imagen

38 A partir de las imágenes '0207.png', '0005.png', '0014.png' y '0167.png' se sintetizó una nueva  
39 imagen tomando de cada una la esquina izquierda superior, la esquina derecha superior, la  
40 esquina izquierda inferior y la esquina derecha inferior respectivamente. Para ello se toman  
41 las dimensiones de las imágenes, se toma el punto medio del ancho y del alto, y para cada  
42 imagen se toma la sección correspondiente de la matriz.

### 43 5.3. Aplicación de filtros

44 En la figura 6 se muestra el resultado de aplicar el filtro del operador Sobel sobre la dirección  
45 X y la dirección Y. Este operador transforma el valor de cada pixel en función de la magnitud  
46 de cambio de los valores de los píxeles adyacentes en una de las direcciones X o Y. En  
47 la implementación elegida se usa `filters.sobel_v` y `filters.sobel_h` que equivalen a  
48 realizar una convolución sobre la matriz aplicando los siguientes kernels:

```
sobel_v_kernel = np.array([[ -1,  0,  1],  
                           [ -2,  0,  2],  
                           [ -1,  0,  1]], dtype=np.float32)  
  
sobel_h_kernel = np.array([[ -1, -2, -1],  
                           [  0,  0,  0],  
                           [  1,  2,  1]], dtype=np.float32)
```

#### Detección de Bordes con Sobel (usando scikit-image)

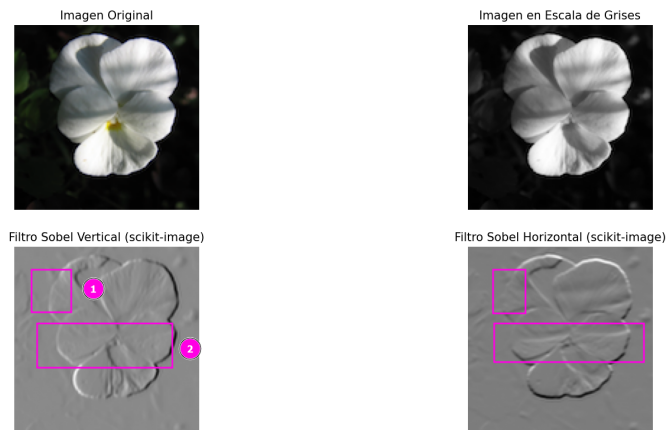


Figure 6: Figura 6. Aplicación de filtros

49 Se destaca que aplicar el filtro vertical ha permitido resaltar el contorno de la figura dentro  
50 del recuadro 1, que en la imagen original mostraba un contraste tenue. Mientras que aplicar  
51 el filtro horizontal ha resaltado por ejemplo el contraste de la proyección de sombras sobre el  
52 propio pétalo de las flores en el recuadro 2, además de los bordes de la flor. Por otra parte,  
53 cabe señalar que el gradiente resultante de la aplicación de los filtros puede tomar valores  
54 positivos o negativos, de modo que los bordes se verán oscuros o claros en la representación  
55 gráfica según sea la distribución de brillo. Por ejemplo se observa que la aplicación del filtro  
56 horizontal resulta en un gradiente positivo sobre el borde superior de la flor, mientras que  
57 da un gradiente negativo en el borde inferior.

#### 58 5.4. Imagen promedio

59 En las figuras 7 y 8 se muestra la representación gráfica de promediar todas las imágenes  
60 de cada especie tanto a color como en blanco y negro. Para promediar las imágenes en  
61 blanco y negro se realizó en este caso la binarización de las mismas mediante el método de  
62 Otsu. Se utilizó este método ya que las diferencias de iluminación y escala de brillo entre  
63 las diferentes imágenes y diferentes especies podía implicar introducir un sesgo inesperado  
64 al elegir un umbral fijo para todo el conjunto de imágenes.

65 Se observa en la figura 7 que la principal diferencia que se percibe a simple vista entre las  
66 especies radica en el color predominante de cada una.

Imágenes Promedio por Especie (Color)

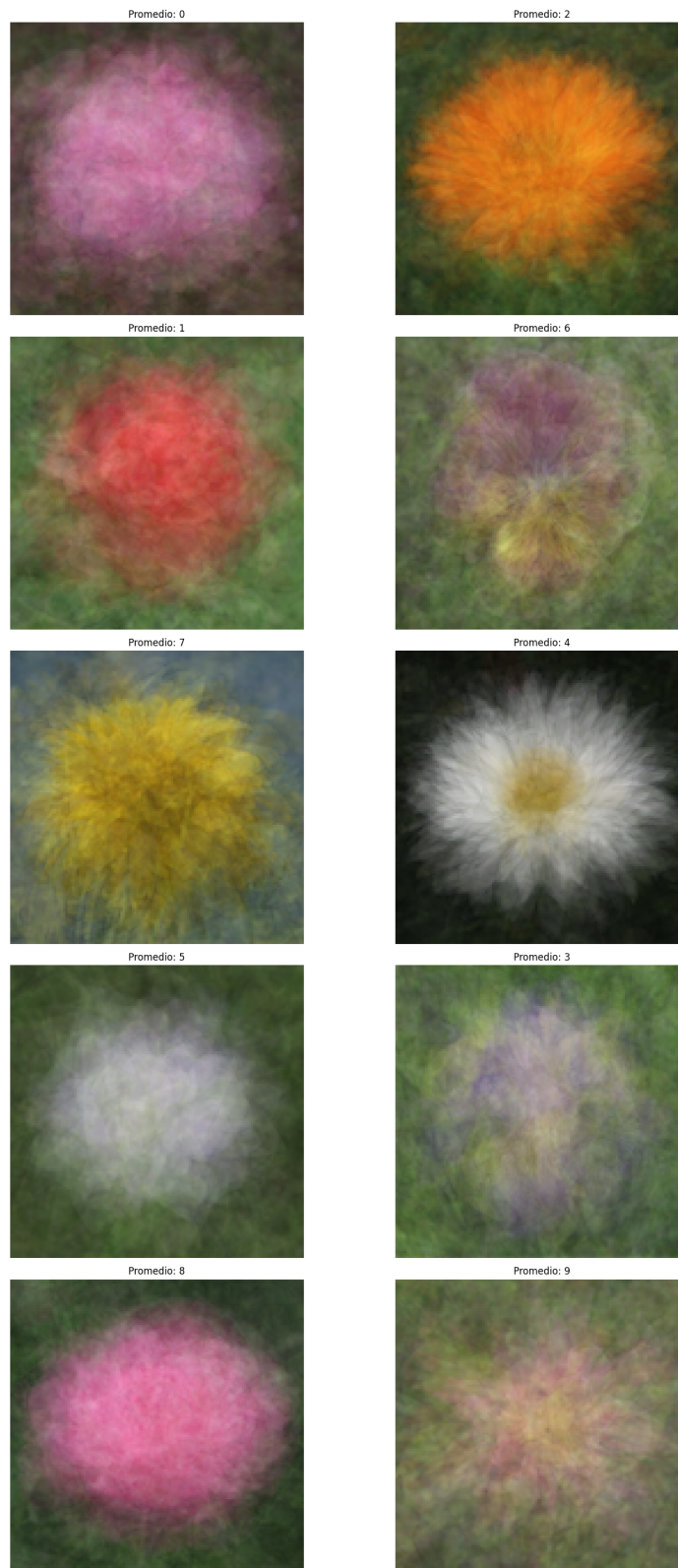


Figure 7: Figura 7. Promedio de imágenes a color

67 Como muestra la figura 8 al binarizar las mismas los atributos resultantes que se perciben a  
68 simple vista son la forma del contorno de la flor y el brillo promedio la misma, sin embargo en  
69 un análisis visual no parece haber diferencias particulares a nivel especie aunque es posible  
70 agruparlas por similitud. a) especie 0, 5 y 8, b) especies 2 y 4, y c) las especies restantes.

Imágenes Promedio por Especie (Grises)

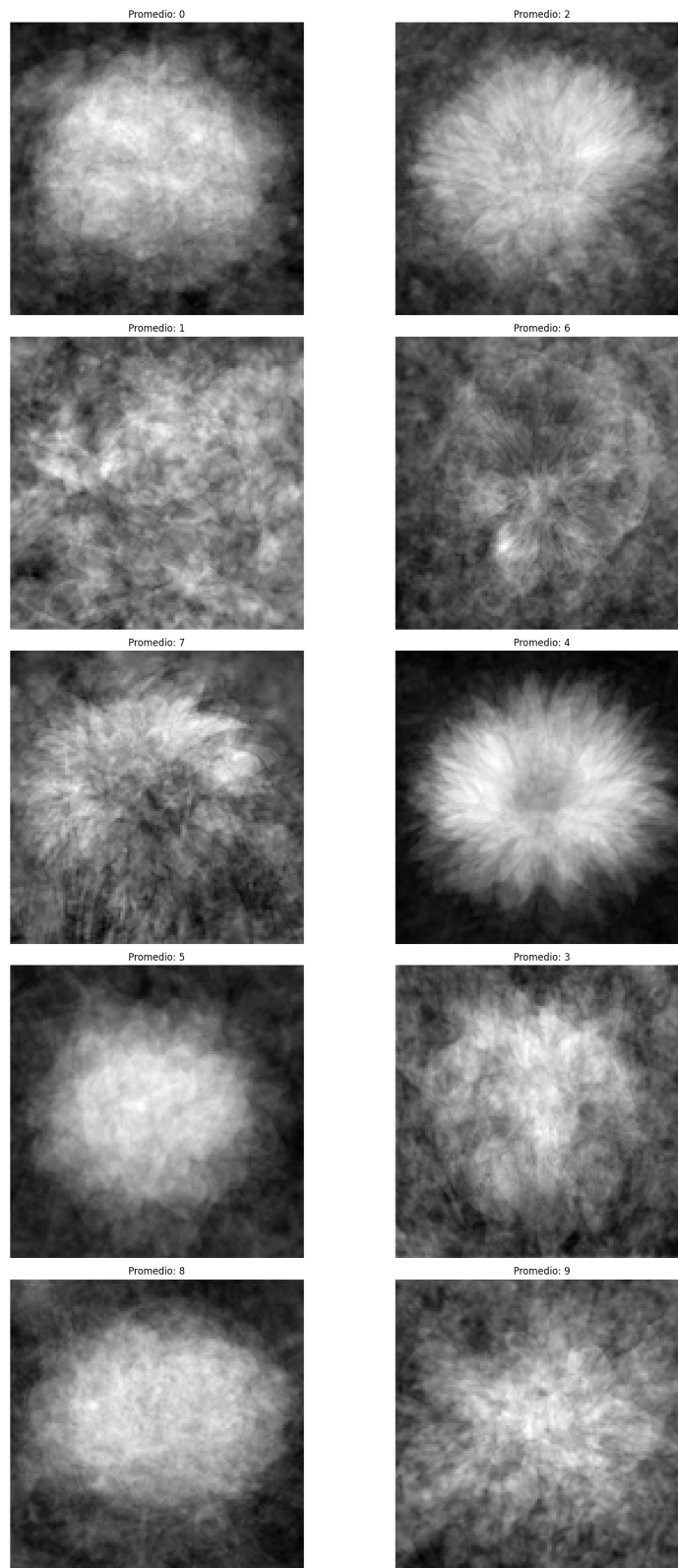


Figure 8: Figura 8. Promedio de imágenes en blanco y negro



72 A partir del conjunto de imagenes de cada especie se grafico en la figura 9 la distribución  
73 de densidad de pixeles por valor para cada uno de los canal RGB.

## Distribución de Intensidad de Píxeles por Especie y Canal de Color

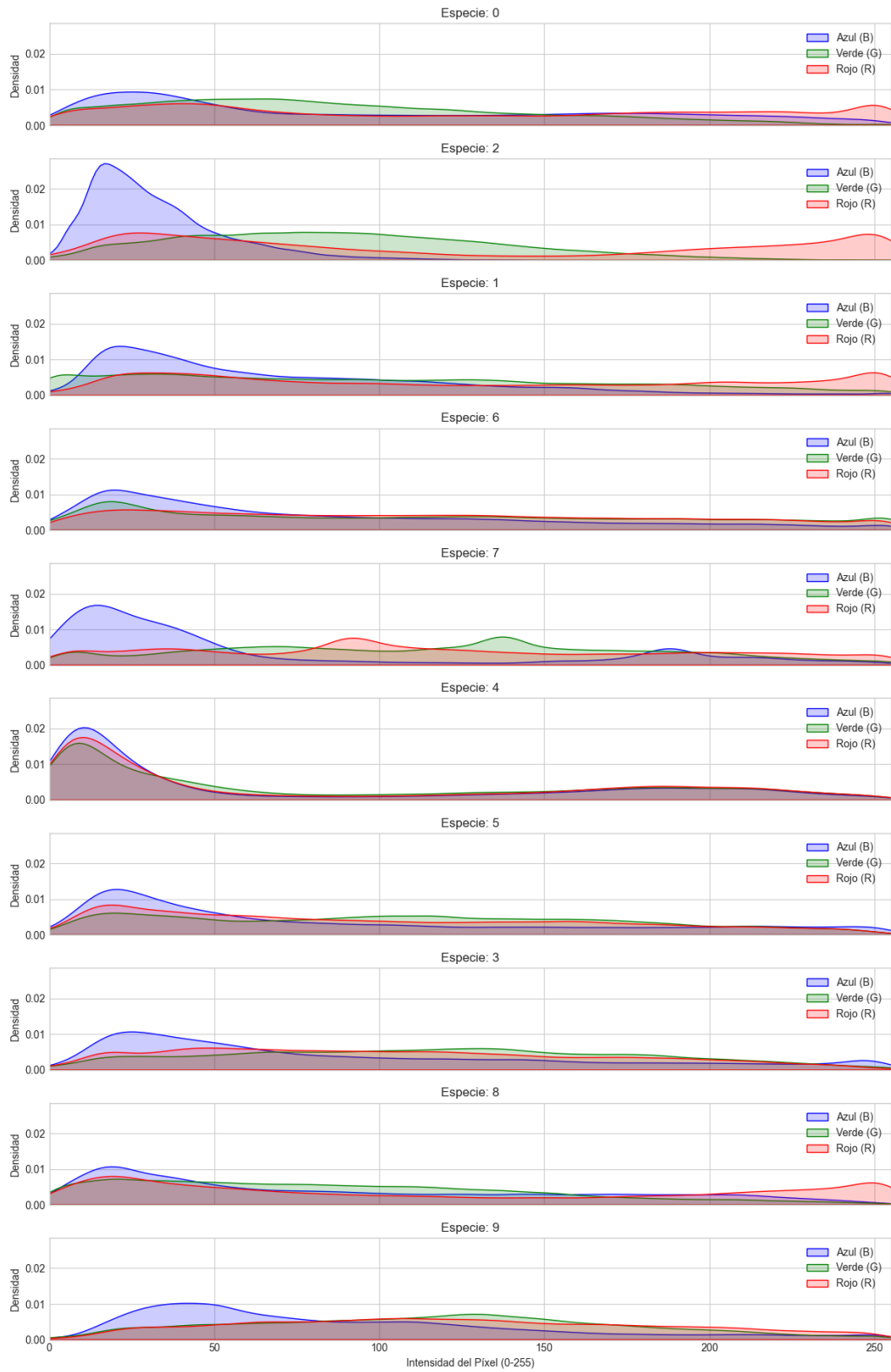


Figure 9: Figura 9. Densidad de píxeles por intensidad para cada uno de los canales RGB por especie

Se puede observar en la figura 9 que algunas especies son distinguibles a simple vista según su distribución. Las especies 2, 7, y 4 son en una lectura inicial las que presentan distribuciones más características. La especie 2 tiene una alta densidad de pixeles de baja intensidad en el canal azul junto una densidad comparativamente alta (respecto al resto de las especies) de pixeles en el canal rojo de alta intensidad. La especie 7 también muestra una alta densidad de pixeles en el canal azul a baja intensidad pero en ese caso acompañados por una concentración de pixeles en el canal verde a media intensidad. Por último la especie 4 se caracteriza por tener una alta densidad de pixeles de baja intensidad en los tres canales en simultáneo.

## 6.2. Análisis de Componentes Principales (PCA)

Se realizó un análisis de componentes principales (PCA) del dataset. La dimensión total del dataset consiste en 210 registros de 49125 variables. Primero se estandarizaron los valores respecto a la media de cada variable y luego se aplicó PCA. La figura 10 muestra el % de varianza explicada en función de la cantidad de componentes tomadas.

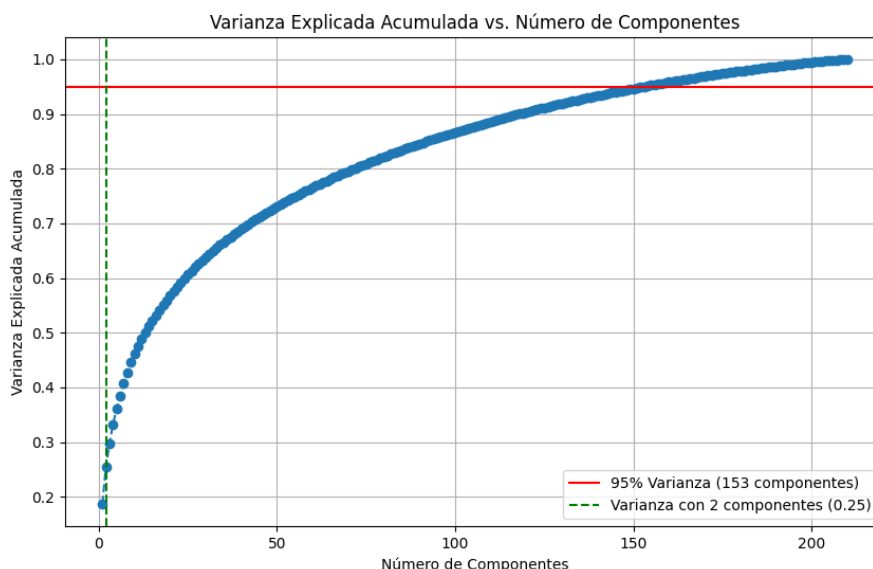


Figure 10: Figura 10. Porcentaje de varianza explicada según cantidad de componentes

Se observa que crecimiento varianza explicada alcanza el 25% con 2 componentes y el 95% con 153 componentes. Esto implica que la dimension del dataset puede ser reducida desde las 49125 hasta las 153 componentes con solo un 5% de pérdida de la varianza.

En la figura 11 se representan todos los registros en las dimensiones de los dos primeros componentes principales que como se expresó anteriormente explican el 25% de la varianza.

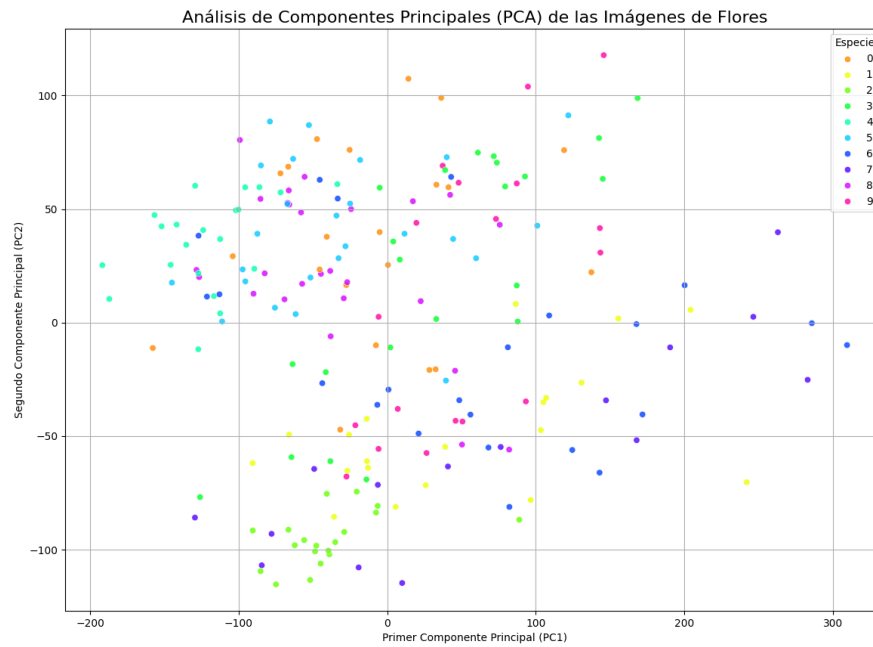


Figure 11: Figura 11. PCA 1 vs PCA 2

92 Se observa que si bien algunas especies parecen agruparse en ciertas regiones del plano: la  
 93 especie 2 en la esquina inferior izquierda, la especie 4 en la region superior izquierda; existe  
 94 un solapamiento entre los puntos que hace difícil una separación espacial clara con solo 2  
 95 componentes. Es recomendable utilizar más componentes para poder separar las especies  
 96 con precisión.