

Theorem 1: If  $K$  is a kernel,  $rK$  is a kernel for  $r > 0$ .

Proof: If  $\phi$  is feature map for  $K$ ,  $r\phi$  is feature map for  $rK$ .

Theorems from lectures:

•② 2: If  $K_1, K_2$  are kernels,  $K_1 + K_2$  is a kernel

3: If  $K_1, K_2$  are kernels,  $K_1 K_2$  is a kernel

4:  $x^T y$  is a kernel and  $(x^T y)^n$  is also a kernel

Theorem 5: If  $K$  is a kernel, then  $e^K$  is also a kernel

Proof:  $e^K = I + \frac{K}{1!} + \frac{K^2}{2!} + \dots$

It can be seen as a kernel using ③, ① and ④

From ④, ③ and ①, each term is a kernel because,  $K^n$  is obtained by multiplying  $K$   $n$  times, which by ③ (or by ⑦) makes  $K^n$  a kernel and by making  $r$  as  $\frac{1}{n!}$  in ①,  $\frac{K^n}{n!}$  is a kernel and by ④, summation of all such kernels is a kernel.

$\therefore e^K$  is a kernel

Theorem 7: If  $K$  is a kernel,  $K^n$  is a kernel

Proof: For  $n=1$ ,  $K^n = K$ , which is a kernel

Let for  $n=m-1$ ,  $K^n$  is a kernel

$$K^m = K \cdot (K^{m-1})$$

By ④,  $K^m$  is a kernel

$\therefore$  By induction,  $K^n$  is a kernel

Theorem 6: If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , and  $K$  is a kernel,  $f(x)K(x,y)f(y)$  is also a kernel

Proof: If  $\phi$  is a feature map of  $K$ ,  $f(x)\phi(x)$  is the feature map for the new kernel

$$\text{Now, consider } K(x,y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}} = e^{-\frac{\|x\|^2}{2\sigma^2}} e^{(\frac{x^T y}{\sigma^2})} e^{-\frac{\|y\|^2}{2\sigma^2}}$$

By ④ and ①,  $\frac{1}{\sigma^2} x^T y$  is a kernel by taking  $r$  as  $\frac{1}{\sigma^2}$

By ⑤,  $e^{(\frac{x^T y}{\sigma^2})}$  is a kernel

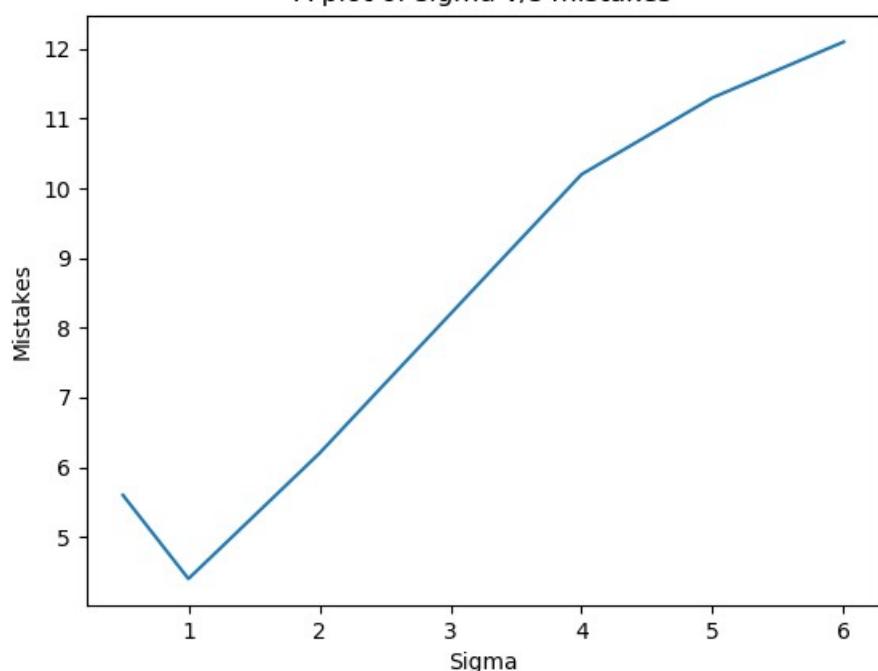
By ⑥,  $e^{-\frac{\|x\|^2}{2\sigma^2}} e^{(\frac{x^T y}{\sigma^2})} e^{-\frac{\|y\|^2}{2\sigma^2}}$  is a kernel by taking  
 $f(x) = e^{-\frac{\|x\|^2}{2\sigma^2}}$

$\therefore K(x,y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$  is a kernel

1.2

b) (ii)

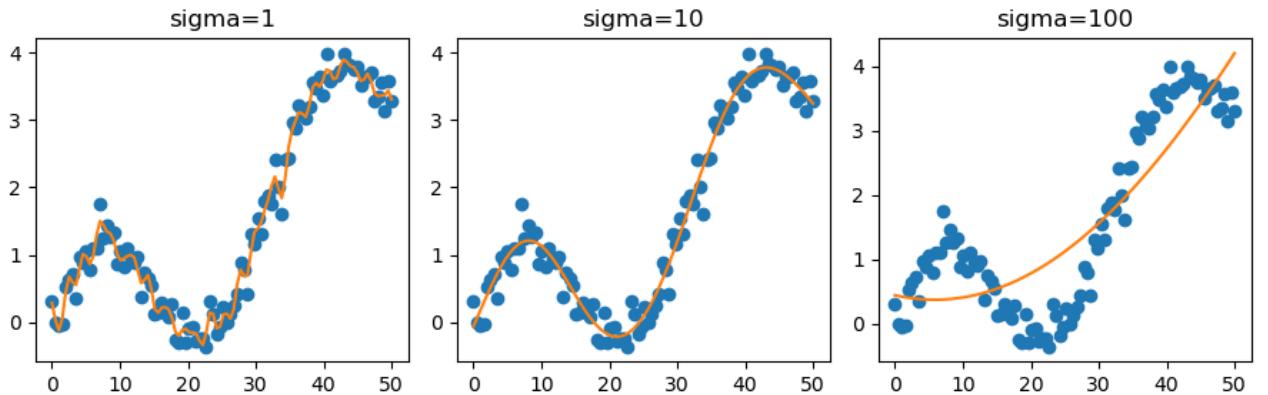
A plot of sigma v/s mistakes



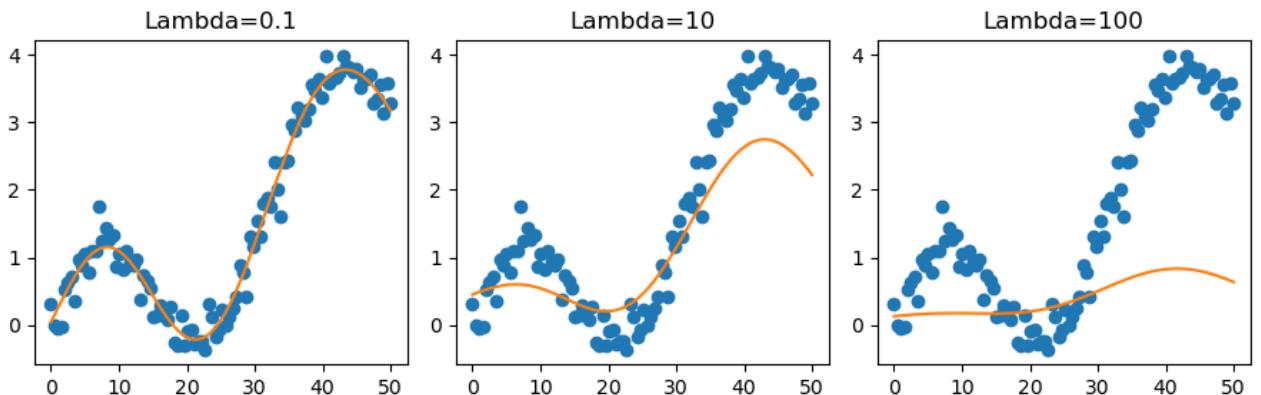
From the plot, the best value of sigma is 1.

(iii) The sigma in gaussian kernel controls the tolerance for similarity between two points. If sigma is low, the similarity between two points has to be very high or  $\|x-y\|^2$  must be low to make the kernel function value high. If sigma is high, even if similarity is less or  $\|x-y\|^2$  is high, the kernel function value becomes higher. So, as we increase sigma, we accomodate more and more points for a given point as being similar to it. So, until we reach the ideal value of sigma, number of mistakes reduces because initially we were overfitting, but once we cross the ideal value, the number of mistakes increase because underfitting occurs.

c) (ii)



The sigma in gaussian kernel controls the tolerance for similarity between two points. If sigma is low, the similarity between two points has to be very high or  $\|x-y\|^2$  must be low to make the kernel function value high. If sigma is high, even if similarity is less or  $\|x-y\|^2$  is high, the kernel function value becomes higher. So, as we increase sigma, we accomodate more and more points for a given point as being similar to it. As seen in the graph, until we reach the ideal value of sigma, we were overfitting because we considered only the closest points, but once we cross the ideal value, underfitting occurs because now we are also considering those points also which are not properly related.



Since lambda is the parameter for imposing penalty on weights to avoid overfitting, as lambda increases, the penalty imposed increases which makes the values of weights smaller and closer to

zero. As seen in the graph, as lambda increases after the optimal value, we tend to underfit the data. Also, decreasing lambda beyond the optimal value may cause overfitting due to large weights.

## 2.1

(i)

since we already know that  $K(x,y)$  is a kernel, there must exist some  $\phi$  such that  $K(x,y) = \langle \phi(x), \phi(y) \rangle$

~~that~~  $\phi: \mathbb{R}^m \rightarrow H$ , where  $H$  is the Hilbert space

Now consider  $\phi_g(x) = \phi(g(x))$ , where,  $g: \mathbb{R}^m \rightarrow \mathbb{R}^m$

If we consider ~~the kernel~~  $K(g(x), g(y))$ , we can write it as

$K(g(x), g(y)) = \langle \phi_g(x), \phi_g(y) \rangle = K_{\text{new}}(x, y)$

As  $\phi_g: \mathbb{R}^m \rightarrow H$  is a feature map for  $K_{\text{new}}$ , we can say that its a kernel.

(ii)

Theorem 1: If  $K$  is a kernel,  $rK$  is a kernel for  $r > 0$ .

Proof: If  $\phi$  is feature map for  $K$ ,  $r\phi$  is feature map for  $rK$ .

Theorems from lectures:

• 2: If  $K_1, K_2$  are kernels,  $K_1 + K_2$  is a kernel

3: If  $K_1, K_2$  are kernels,  $K_1 K_2$  is a kernel

Theorem 7: If  $K$  is a kernel,  $K^n$  is a kernel.

Proof: For  $n=1$ ,  $K^1 = K$ , which is a kernel.

Let for  $n=m-1$ ,  $K^m$  is a kernel.

$$K^m = K \cdot (K^{m-1})$$

By ③,  $K^m$  is a kernel.

$\therefore$  By induction,  $K^n$  is a kernel.

Theorem 8: If  $g$  is a polynomial with non-negative coefficients,  $g(K)$  is a kernel given  $K$  is a kernel.

Proof: Every term of  $g(K)$  is of the form  $\alpha K^n$  where,  $\alpha > 0$ .

By ⑦ and ⑥,  $\alpha K^n$  is a kernel.

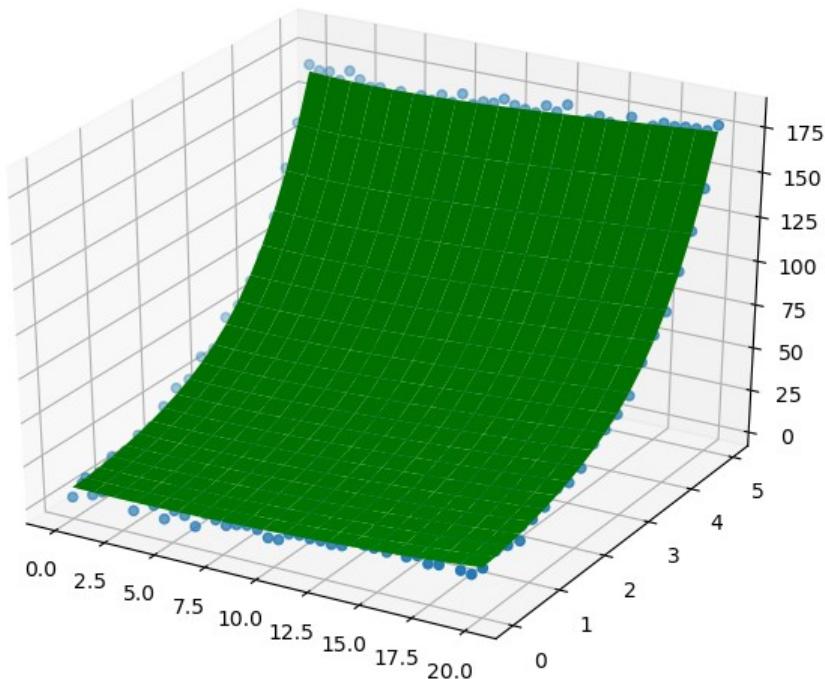
$$g(K) = \sum \alpha K^n$$

As each term is a kernel, by applying ② repetitively, the whole summation is a kernel.

As the whole summation is  $g(K)$ ,  $g(K)$  is a kernel.

$\therefore g(K)$  is a kernel.

2.2 The kernel function used is  $(\langle x, y \rangle + 1)^4$ .



Let  $\bar{x}_1$  be the first cluster mean (cluster centre)

$$\bar{x}_1 = \frac{x_1 + x_2 + \dots + x_m}{m}$$

Let  $\bar{x}_2$  be the second cluster mean (cluster centre)

$$\bar{x}_2 = \frac{x_{m+1} + x_{m+2} + \dots + x_n}{n-m}$$

Let the separating hyperplane be  $a \cdot x + b = 0$

Claim:  $a = \bar{x}_1 - \bar{x}_2$

$$b = \frac{\|\bar{x}_2\|^2 - \|\bar{x}_1\|^2}{2} \quad \text{and for } x_1, \dots, x_m, \text{ as}$$

Proof: Let  $x$  belong to the first cluster ~~but  $a \cdot x + b > 0$~~   
As  $x$  belongs to first cluster,

$\|x - \bar{x}_1\|^2 < \|x - \bar{x}_2\|^2$  ( $\because$  no point is equidistant from two centres)

$$(x - \bar{x}_1)^T (x - \bar{x}_1) < (x - \bar{x}_2)^T (x - \bar{x}_2)$$

$$\|x\|^2 - 2\bar{x}_1^T x + \|\bar{x}_1\|^2 < \|x\|^2 - 2\bar{x}_2^T x + \|\bar{x}_2\|^2$$

$$2(\bar{x}_1^T - \bar{x}_2^T)x > \|\bar{x}_1\|^2 - \|\bar{x}_2\|^2$$

$$2(\bar{x}_1^T - \bar{x}_2^T)x > \|\bar{x}_1\|^2 - \|\bar{x}_2\|^2 \quad \text{---} \textcircled{1}$$

$$a \cdot x + b \leq 0$$

$$(\bar{x}_1 - \bar{x}_2) \cdot x + \left( \frac{\|\bar{x}_2\|^2 - \|\bar{x}_1\|^2}{2} \right) \leq 0$$

$$(\bar{x}_1 - \bar{x}_2)^T x \leq \frac{\|\bar{x}_1\|^2 - \|\bar{x}_2\|^2}{2}$$

$$2(\bar{x}_1 - \bar{x}_2)^T x \leq \|\bar{x}_1\|^2 - \|\bar{x}_2\|^2 \quad \text{---} \textcircled{2}$$

But  $\textcircled{1}$  and  $\textcircled{2}$  are contradictions of each other

$\therefore$  By proof by contradiction, if  $x \in \{x_1, \dots, x_m\}$ ,  $a \cdot x + b > 0$  and similarly, if  $x \in \{x_{m+1}, \dots, x_n\}$ ,  $a \cdot x + b < 0$

$\therefore$  The two clusters are separated by the hyperplane  $\alpha x + b = 0$ , where,

$$\alpha = \frac{\sum_{i=1}^m x_i}{m} - \frac{\sum_{j=m+1}^n x_j}{n-m}$$

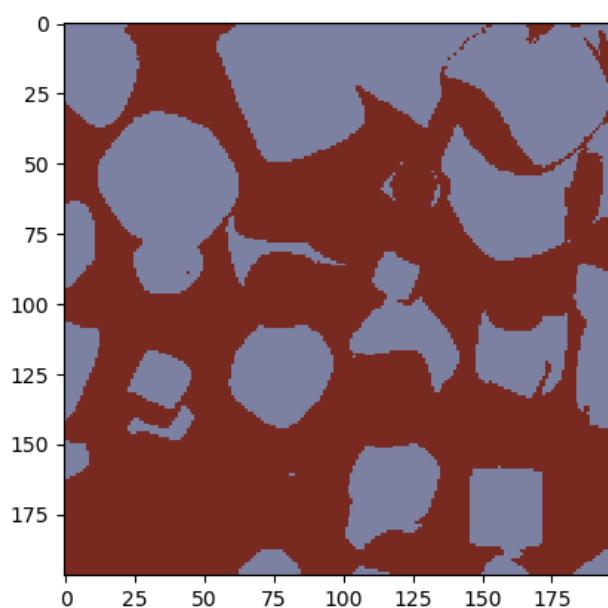
$$b = \frac{\left\| \frac{\sum_{j=m+1}^n x_j}{n-m} \right\|^2 - \left\| \frac{\sum_{i=1}^m x_i}{m} \right\|^2}{2}$$

### 3.2

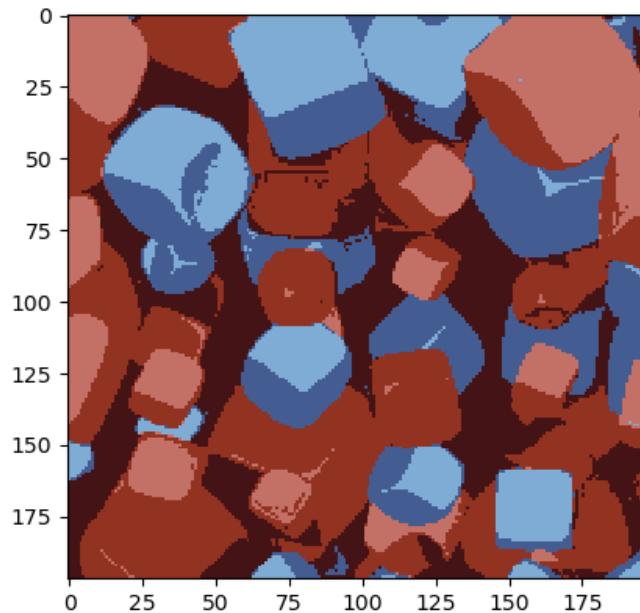
The images are shown in the order of increasing number of clusters 2,5,10.

(ii)

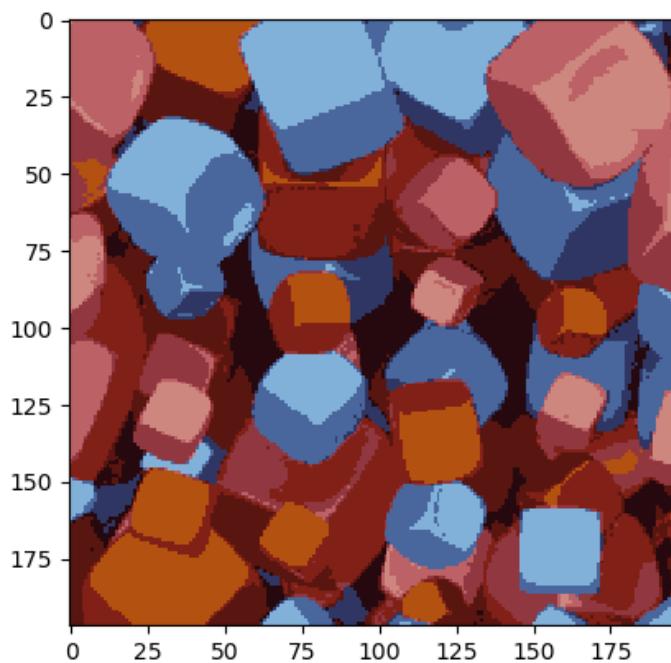
I.



For number of clusters = 2, the generated image is not even close to the actual image, since the number of clusters is low compared to the actual number of clusters in the image.

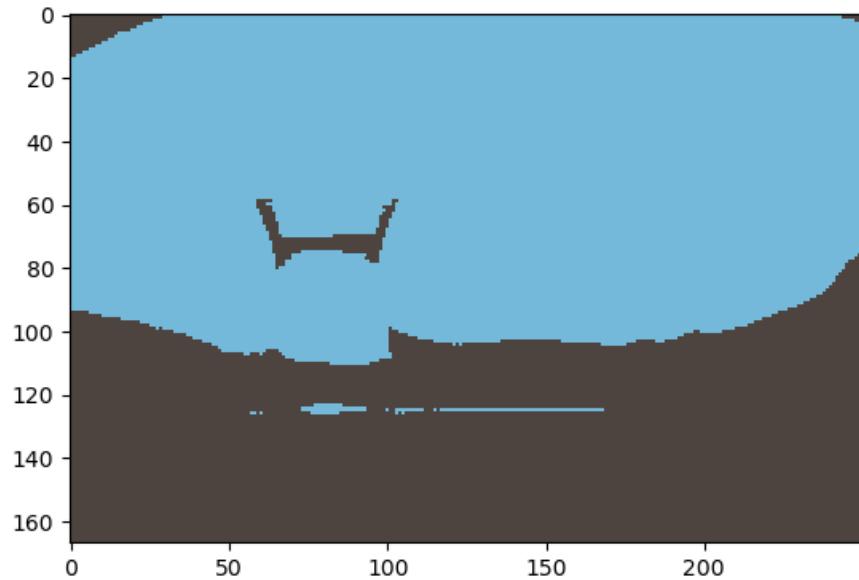


For number of clusters = 5, the generated image is somewhat close to the actual image, since the number of clusters is comparable to the actual number of clusters in the image. Here the structure of the image is captured but colours are not captured properly.

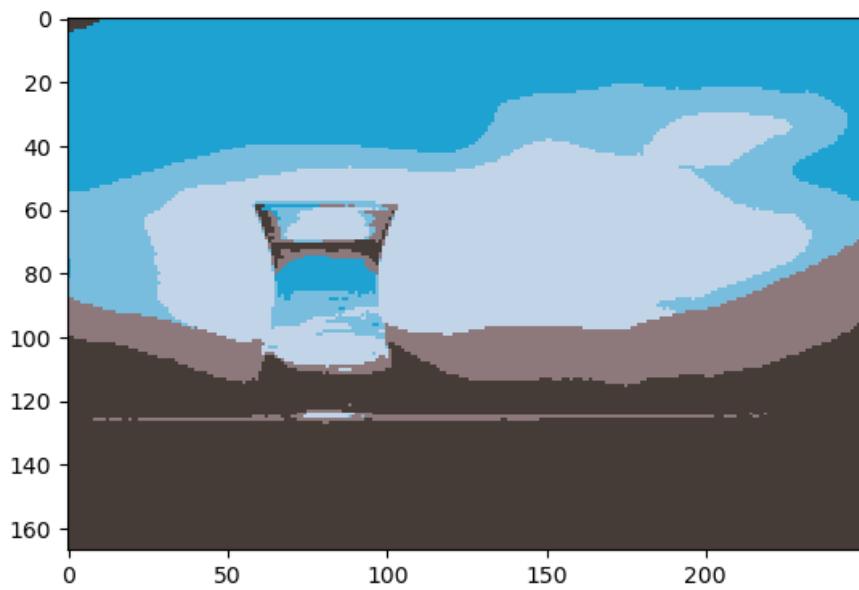


For number of clusters = 10, the generated image is very close to the actual image, since the number of clusters is almost equal to the actual number of clusters in the image.

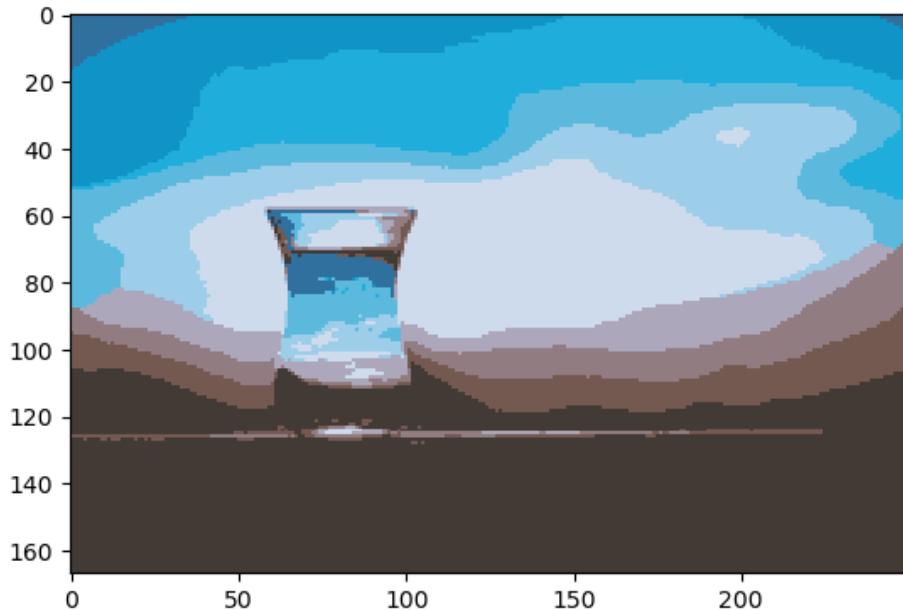
II.



For number of clusters = 2, the generated image is not even close to the actual image, since the number of clusters is low compared to the actual number of clusters in the image.

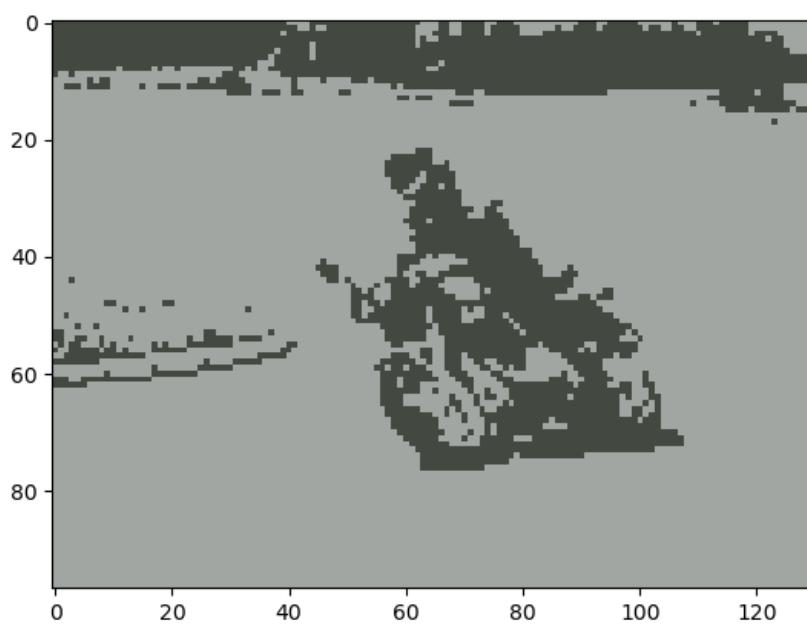


For number of clusters = 5, the generated image is very close to the actual image, since the number of clusters is almost equal to the actual number of clusters in the image.

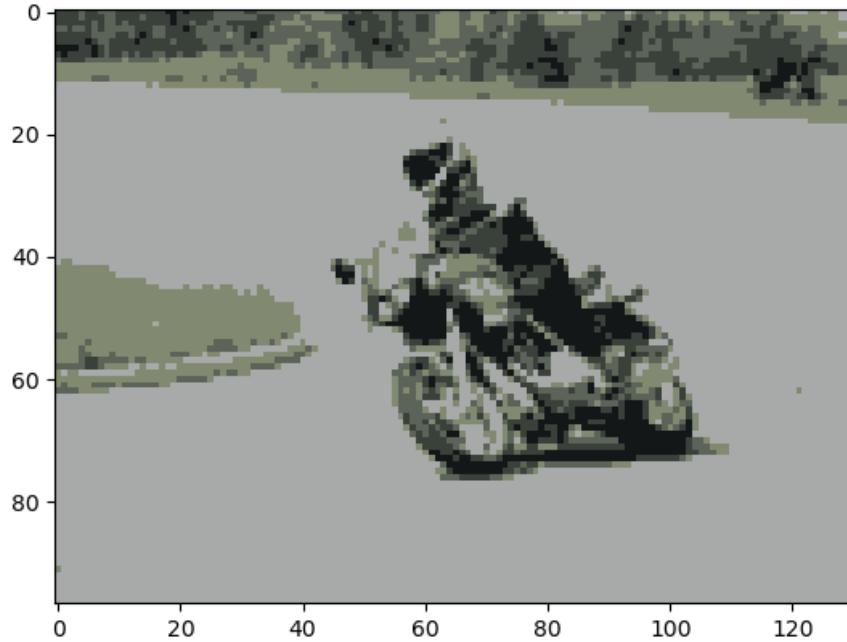


For number of clusters = 10, the generated image adds more detailing to the image with number of clusters as 5, since it has more clusters than the actual number of clusters needed for the image.

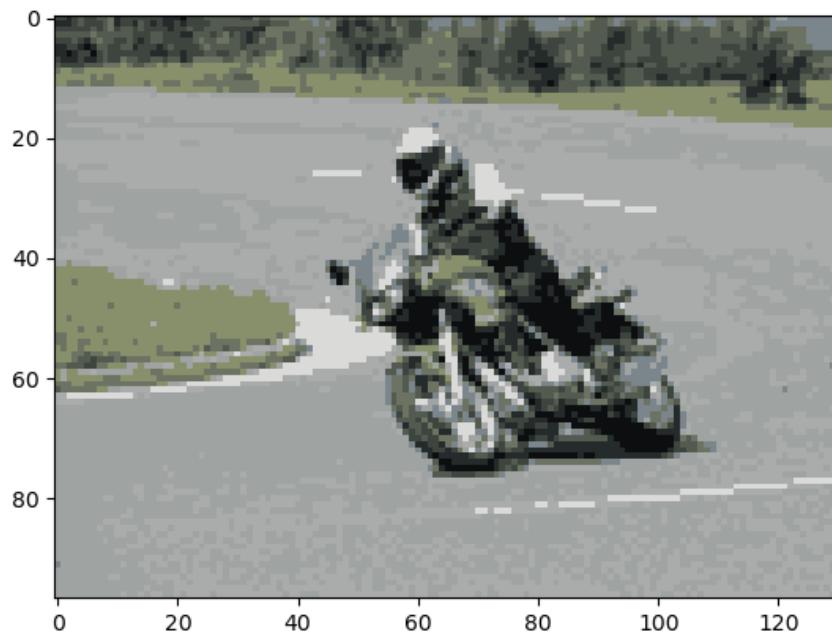
III.



For number of clusters = 2, the generated image is not even close to the actual image, since the number of clusters is low compared to the actual number of clusters in the image.



For number of clusters = 5, the generated image is somewhat close to the actual image, since the number of clusters is less compared to the actual number of clusters in the image.



For number of clusters = 10, the generated image is close but not very close to the actual image, since the number of clusters is less than the actual number of clusters in the image. Here the structure of the image is captured but colours are not captured properly.

(iii) The number of clusters needed to represent the image is almost represented by the actual number of colours and detail in the image(that is by the number of unique (R,G,B) pixels). In the first image, we have less number of colours but more detail, so we required 10 clusters to clearly represent the image. The second image has less colours and less detailing, so only 5 clusters were enough. In the last image, we have more colours and more detailing, so, even 10 clusters are not sufficient. Note that detailing also demands more clusters because, with increasing clusters, points become more closer to the cluster centers and with many clusters, many groups of points with different colours whose cluster center very closely resembles them can be formed.