

1.1

$$K(x, y) = \exp\left(-\frac{1}{2\sigma^2} \|x - y\|^2\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2} \|x\|^2\right) \exp\left(\frac{1}{\sigma^2} x^T y\right) \exp\left(-\frac{1}{2\sigma^2} \|y\|^2\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2} x^T x\right) \exp\left(\frac{1}{\sigma^2} x^T y\right) \exp\left(-\frac{1}{2\sigma^2} y^T y\right)$$

$$K_0(x, y) = x^T y \quad - \text{Pd kernel}$$

We know that if $K(x, y)$ is a Pd kernel
so is $\lambda K(x, y)$ ~~if~~ if $\lambda > 0$

so scaling $K_0(x, y)$ by $\frac{1}{\sigma^2}$
we get

$$K_1(x, y) = \frac{1}{\sigma^2} K_0(x, y) = \frac{x^T y}{\sigma^2}$$

also ~~know~~ note that if K is a Pd kernel
so $\exp(K(x, y))$ is

$$\exp(K(x, y)) = \lim_{N \rightarrow \infty} \sum_{i=0}^N \frac{1}{i!} (K(x, y))^i$$

and $(K(x, y))^n$ is valid kernel because of
products property that $K(x, y) = K_1(x, y) K_2(x, y)$
is Pd kernel if K_1 and K_2 are Pd kernels

So! we now get $K_2(x, y) = \exp(K_1(x, y))$

$$K_2(x, y) = \exp\left(\frac{x^T y}{\sigma^2}\right)$$

~~similar steps for~~

also note that

if K_1 is a Pd Kernel

then for some $f: X \rightarrow \mathbb{R}$

$$K_2(x, y) = f(x) K_1(x, y) f(y)$$

is also a valid kernel

as we can use a feature map

$$\phi'(x) \mapsto f(x) \phi(x)$$

finally

$$K(x, y) = \exp\left(\frac{-1}{2\sigma^2} \|x\|^2\right) \exp\left(\frac{x^T y}{\sigma^2}\right) \exp\left(\frac{-1}{2\sigma^2} \|y\|^2\right)$$

$$K(x, y) = f(x) K_2(x, y) f(y)$$

$$\text{where } f(x) = \exp\left(\frac{-\|x\|^2}{2\sigma^2}\right)$$

$$\text{and } K_2(x, y) = \exp\left(\frac{x^T y}{\sigma^2}\right)$$

• hence $K(x, y)$ is a valid kernel

1.2 (b) (ii) Number of mistakes is minimised for $\sigma = 1$

1.2 (b) (iii) ~~The nature of~~ Graph depicts an decreasing that increasing trend which is justified as follows
 → for smaller sigma the model is overfitting the data, so higher test error than $\sigma = 1$
 → for larger sigma the model is underfitting the data, so higher test error than $\sigma = 1$

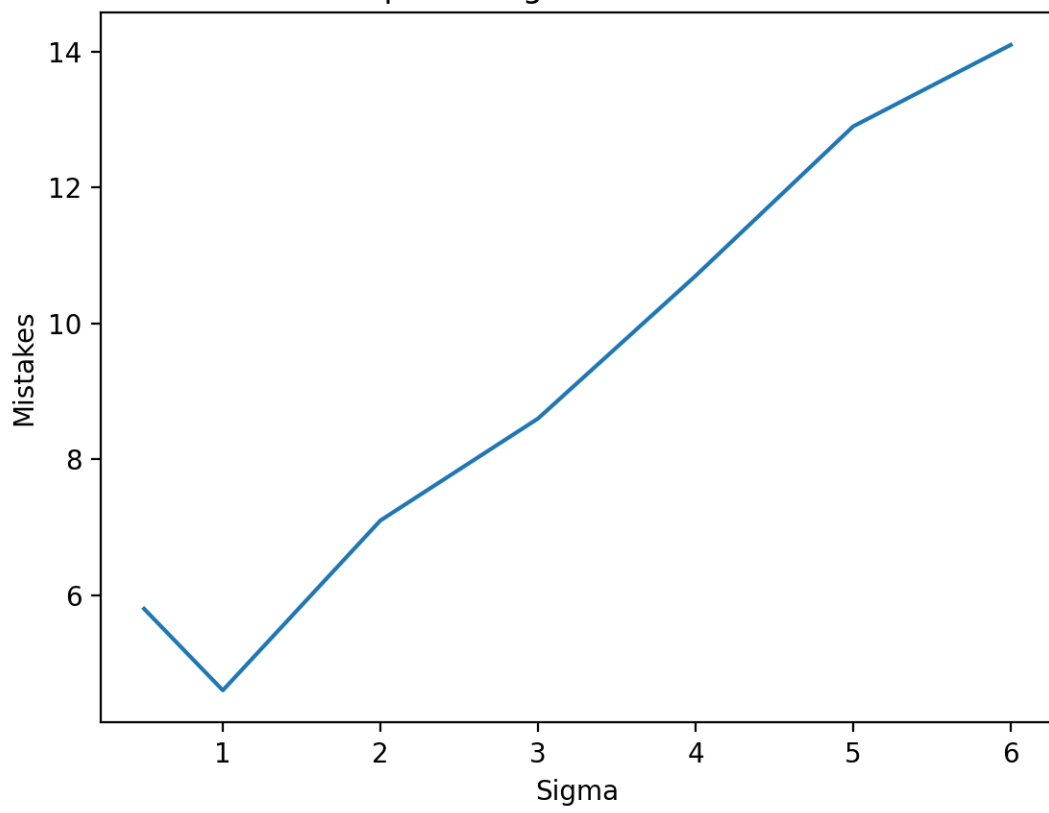
1.2 (c) (ii) for sigma variation

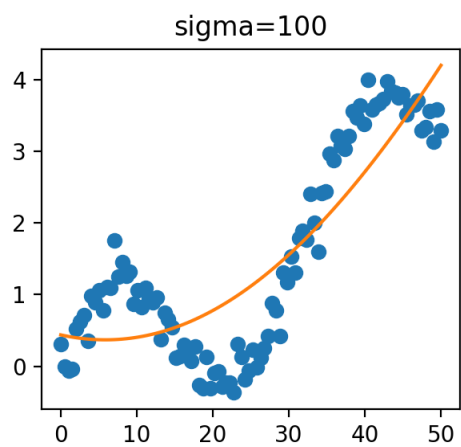
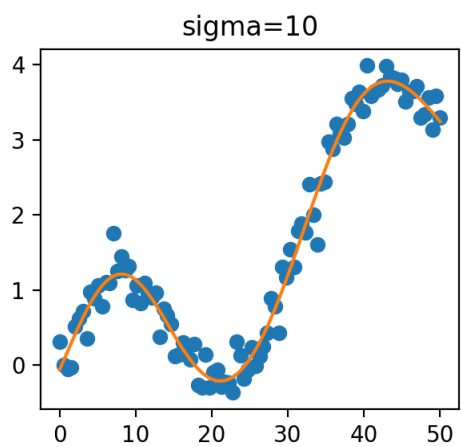
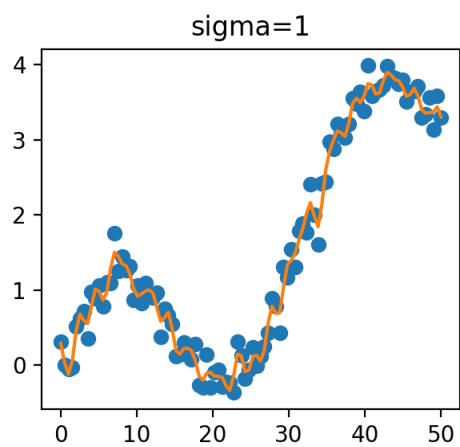
- small sigma overfitting
- large sigma underfitting
- sigma = 10 Good fit

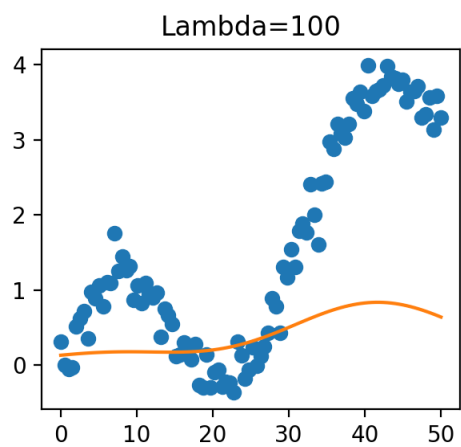
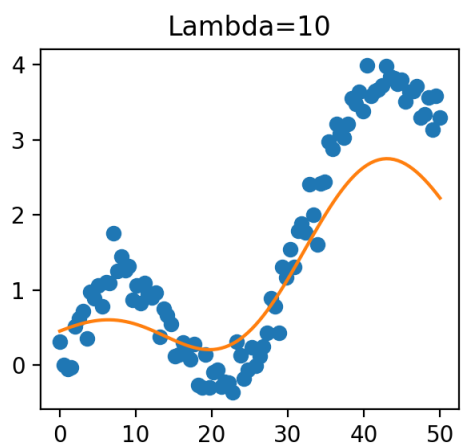
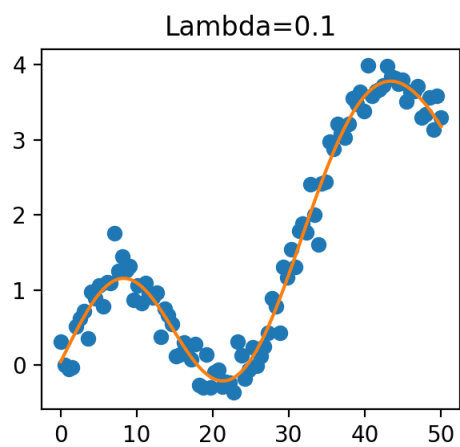
for lambda variation

- As lambda increases model is underfitting

A plot of sigma v/s mistakes







2.2 (i) $K(x_1, x_2)$ is a pd kernel

so ϕ exists such that
 $\phi: R^n \rightarrow H$ such that $K(x_1, x_2) = \langle \phi(x_1), \phi(x_2) \rangle$

and we know that $g(x): R^n \rightarrow R^n$

so consider $\phi_g: R^n \rightarrow H$ such that

$$\phi_g(x) = \phi(g(x))$$

$$\text{then } K_{\text{new}}(x_1, x_2) = K(g(x_1), g(x_2))$$

$$K_{\text{new}}(x_1, x_2) = \langle \phi_g(x_1), \phi_g(x_2) \rangle$$

$K_{\text{new}}(x_1, x_2)$ is a pd kernel

2.1 (ii) Let $q(x) = \sum_i a_i x^i$

where $a_i > 0$ for all i

Now $K(x_1, x_2)$ is a pd kernel

~~it is sufficient to show that~~

$$\text{Now } q(K(x_1, x_2)) = \sum_i a_i (K(x_1, x_2))^i$$

it is a known property that linear combination of valid kernel is also valid kernel so it is sufficient to show that $\lambda (K(x_1, x_2))^n$ is a valid kernel where $\lambda > 0$ and $n \in \mathbb{Z}^+$

Sajal
Date: / / 20
Page:

Now we know that
product of two pd kernel is also a
pd kernel so generalising this
property of pd kernels we can
say that

$$K_{\text{new}}(x, x_2) = K_1(x_1, x_2) K_2(x_1, x_2) \dots K_n(x_1, x_2)$$

K_{new} is a pd kernel if K_i is pd
for all $i \in \{1, n\}$

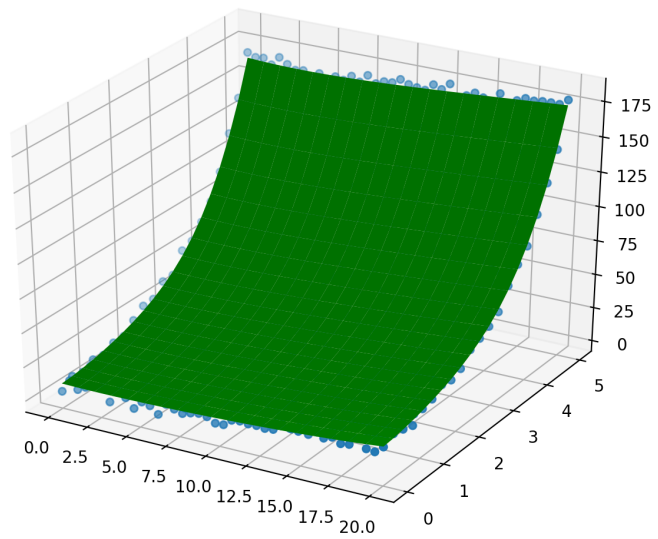
hence $(K(x_1, x_2))^n$ is a valid
kernel from above stated property
if $K(x_1, x_2)$ is valid pd kernel

and using the other property
that $K(x, y) = \lambda K_1(x, y)$
 K is a pd kernel if K is a pd
kernel and $\lambda > 0$

so hence $\lambda (K(x_1, x_2))^n$ is a
pd kernel for $\lambda > 0$ and $n \in \mathbb{Z}^+$

so $g(K(x_1, x_2))$ is a valid
pd kernel

$$\underline{\underline{2.2}} \quad \text{my-kernel}(x, y) = (1 + x^T y)^4$$



3.2Centroid of first cluster, $\alpha_1 = \sum_{i=1}^m x_i / m$ " " second " , $\alpha_2 = \sum_{i=m+1}^n x_i / (n-m)$

For points in first cluster

$$\|x - \alpha_1\|^2 < \|x - \alpha_2\|^2$$

$$2(\alpha_2 - \alpha_1) \cdot x + (\alpha_1 \cdot \alpha_1 - \alpha_2 \cdot \alpha_2) < 0$$

For points in second cluster

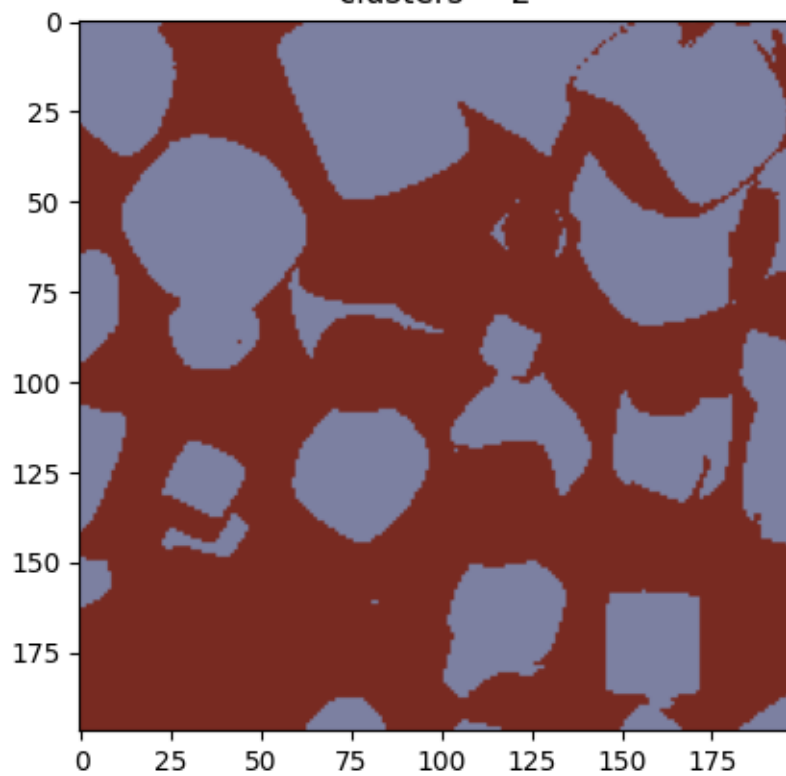
$$\|x - \alpha_1\|^2 > \|x - \alpha_2\|^2$$

$$2(\alpha_2 - \alpha_1) \cdot x + (\alpha_1 \cdot \alpha_1 - \alpha_2 \cdot \alpha_2) > 0$$

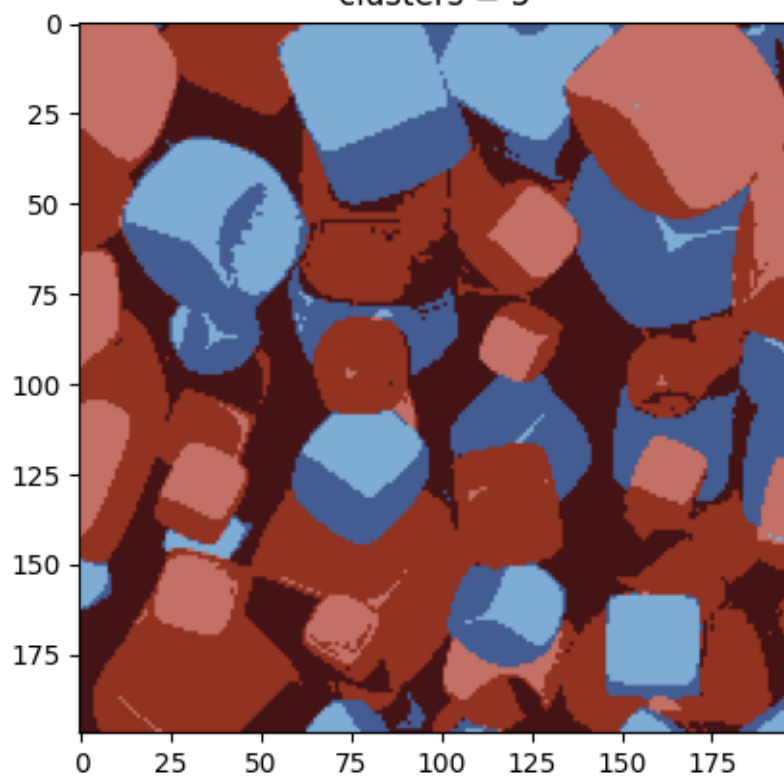
hence $2(\alpha_2 - \alpha_1) \cdot x + (\alpha_1 \cdot \alpha_1 - \alpha_2 \cdot \alpha_2) = 0$

is the separating plane

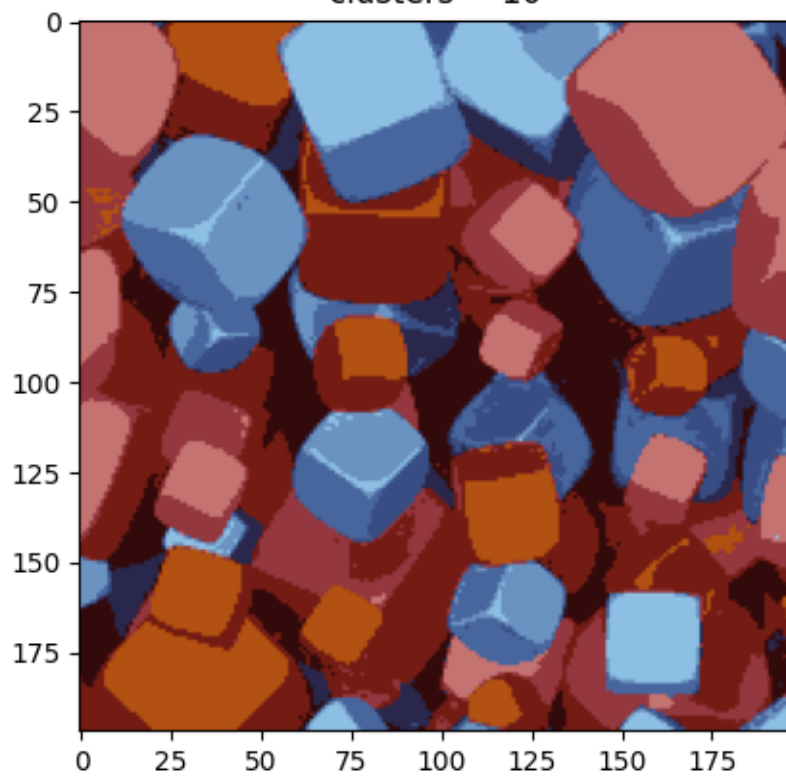
clusters = 2



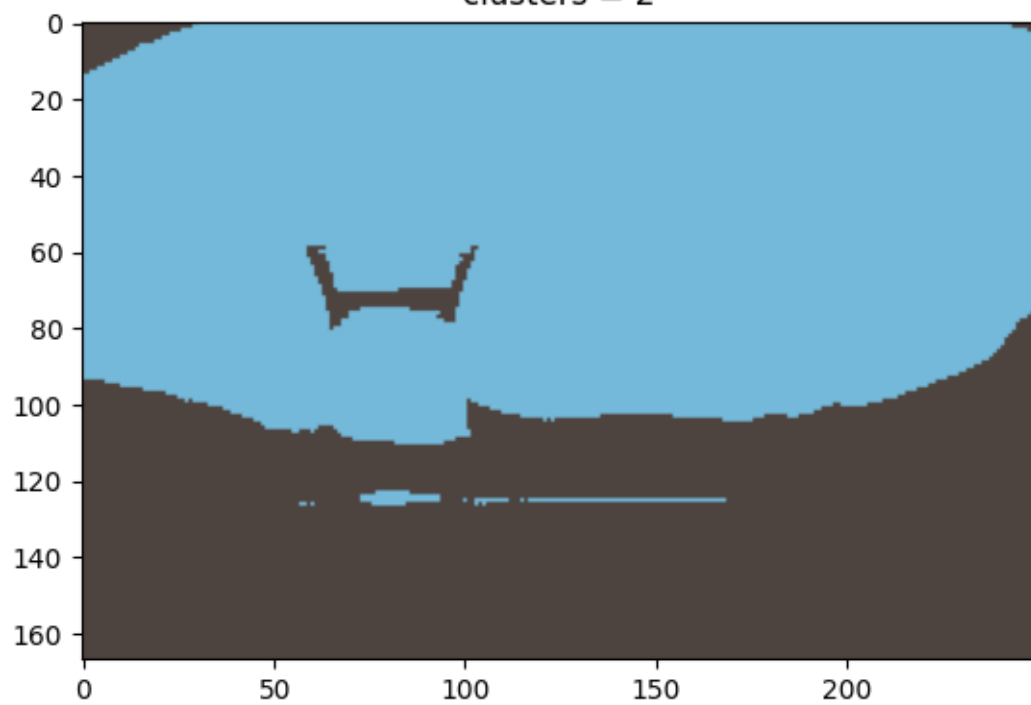
clusters = 5



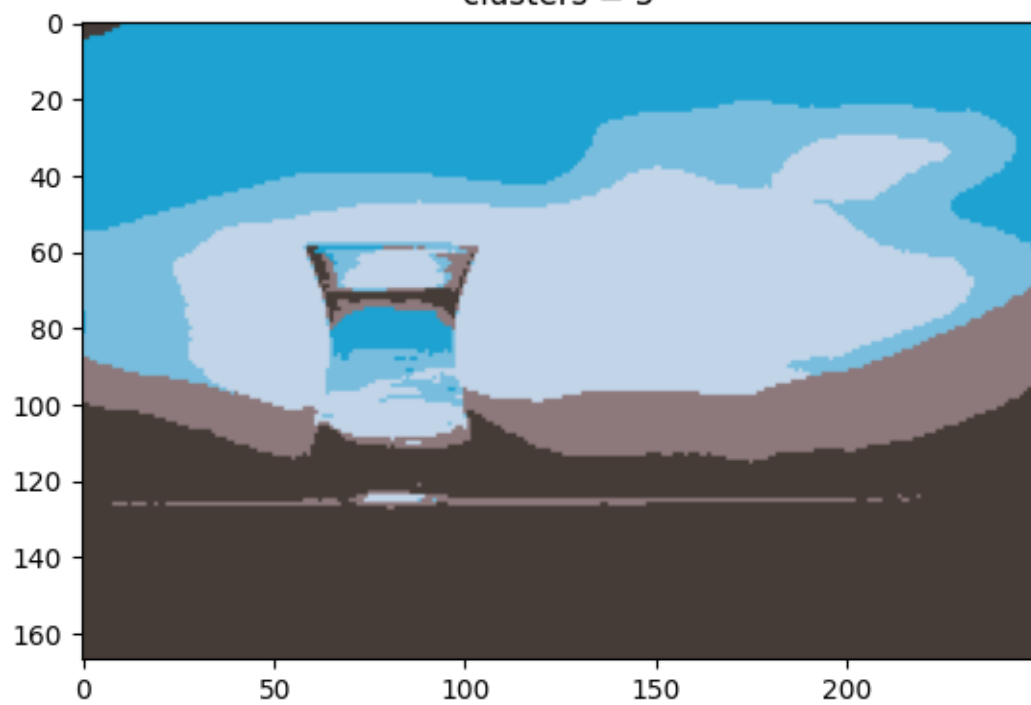
clusters = 10



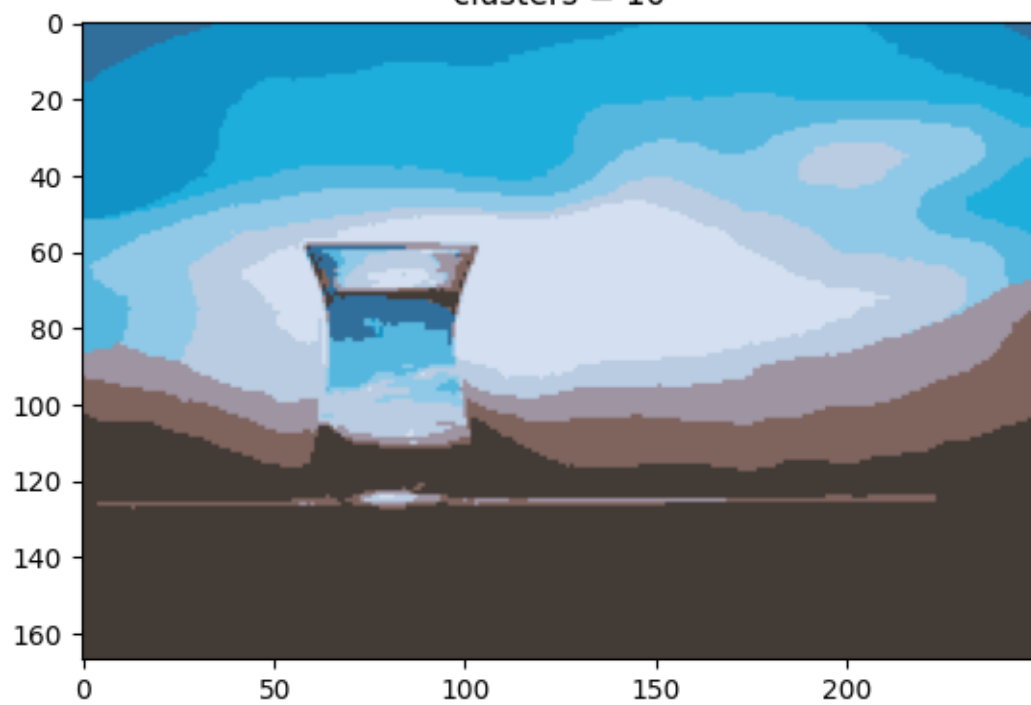
clusters = 2



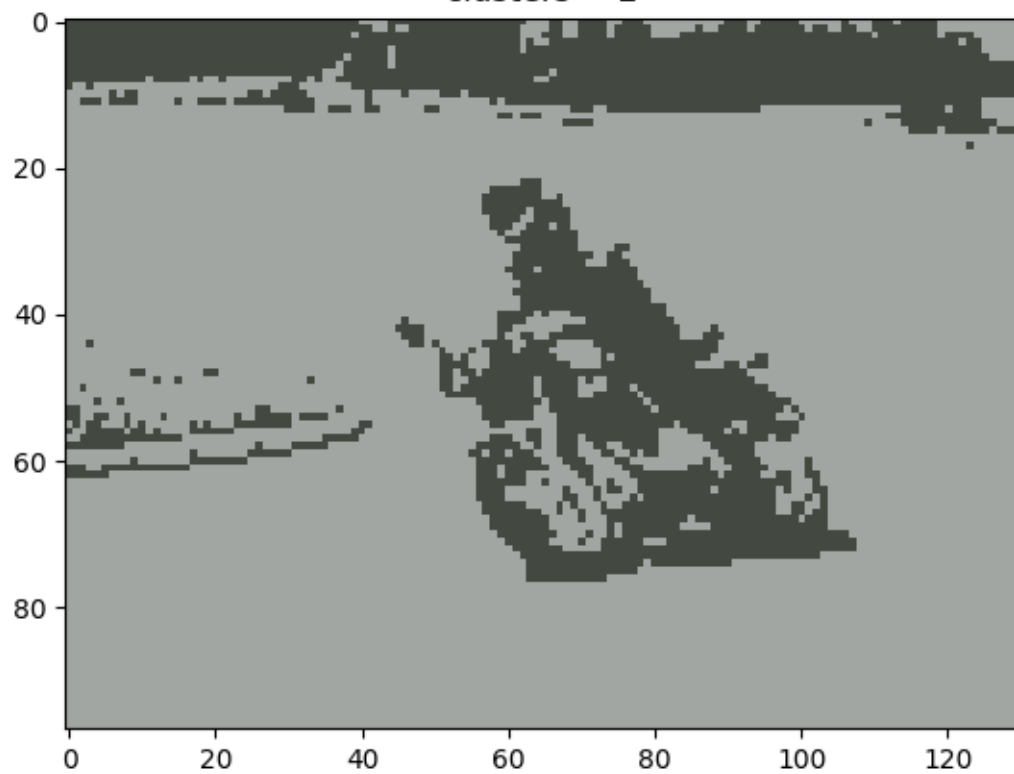
clusters = 5



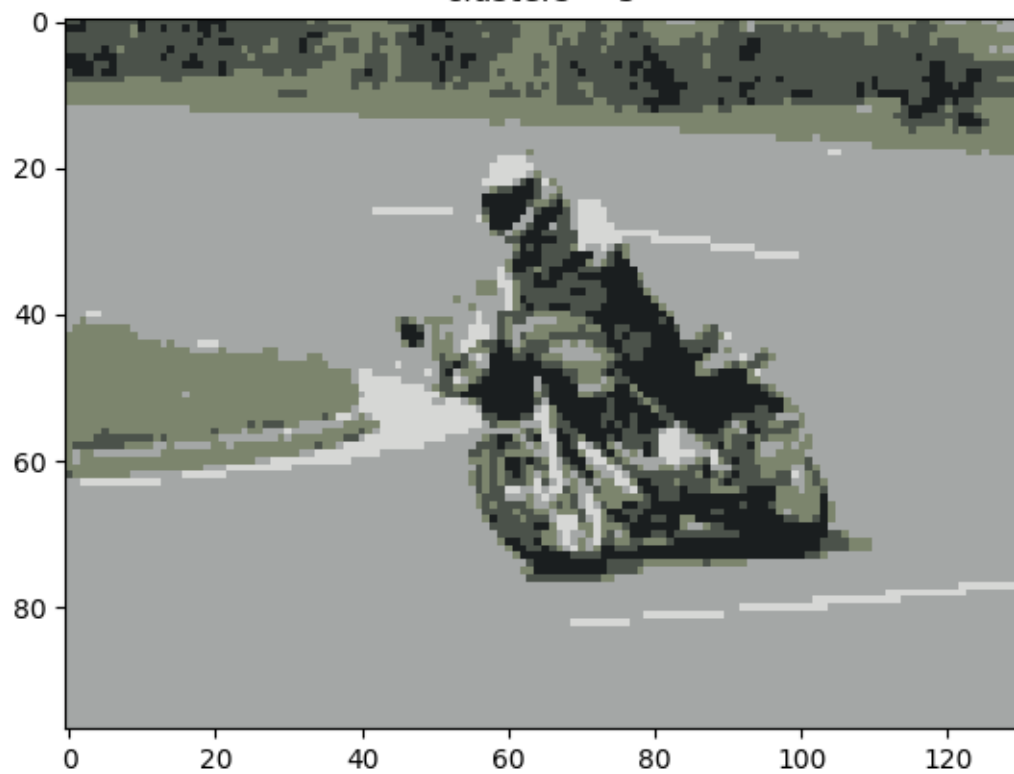
clusters = 10



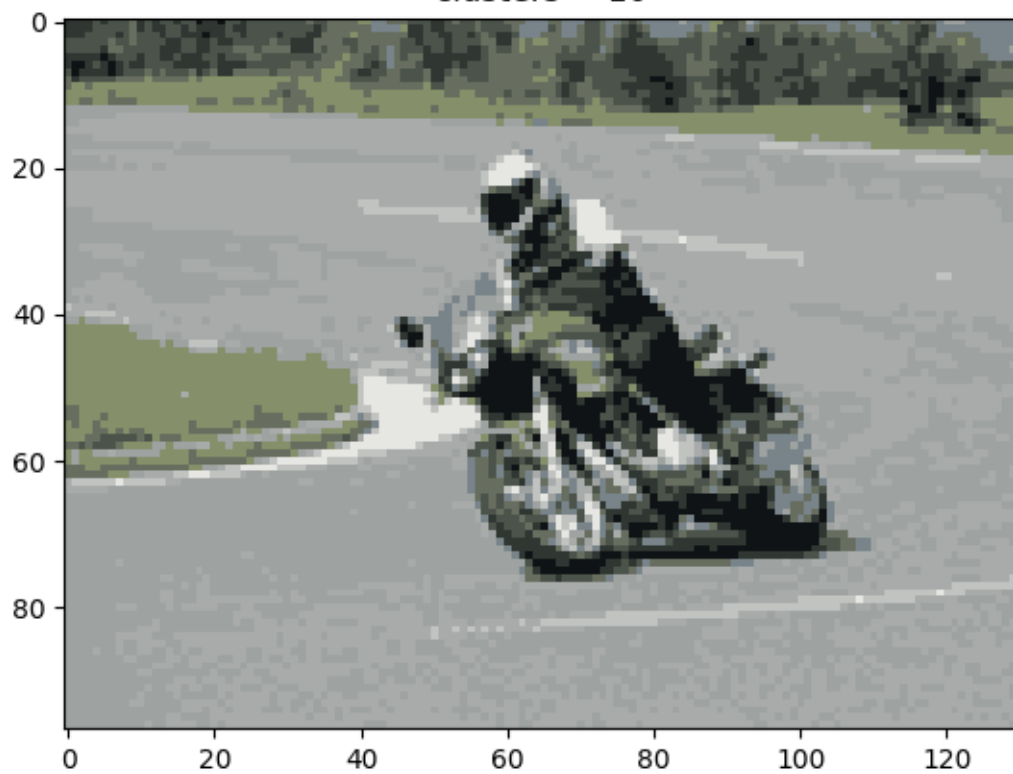
clusters = 2



clusters = 5



clusters = 10



3.2 (iii) The amount of details present in the original image ~~has~~ like number of colors, textures govern the minimum number of clusters required to retain most of details