

Que 2)

(a) Given,

$$\text{kernel is } k(x_i, x) = \frac{\exp \|x_i - x\|_2^2}{\sigma^2} \quad \text{--- (1)}$$

In kernel regression, given a kernel k the output \hat{y} is approximated as :-

$$\hat{y} = \frac{\sum_{i=1}^n k\left(\frac{x_i - x}{n}\right) \cdot y_i}{\sum_{i=1}^n k\left(\frac{x_i - x}{n}\right)} \quad \text{--- (2)}$$

as the value of kernel is given in the question, let's substitute that value in the equation (2)

then, we have

$$\hat{y} = \frac{\sum_{i=1}^n \frac{\exp \|x_i - x\|_2^2}{\sigma^2 n} \cdot y_i}{\sum_{i=1}^n \frac{\exp \|x_i - x\|_2^2}{\sigma^2}}$$

do in
d like
of class
ecision
n.

according to the question

$$\hat{y} = l(x)^T y$$

where $l(x) = (l_1(x), \dots, l_n(x))^T$

and also y is a $n \times 1$ vector

$$\rightarrow \text{let } w_i = \frac{\exp \|x_i - x\|_2^2}{\sigma^2}$$

So,

$$\hat{y} = \frac{\sum_{i=1}^n \left(\frac{\exp \|x_i - x\|_2^2}{\sigma^2 n} \right) y_i}{\sum_{i=1}^n \left(\frac{\exp \|x_i - x\|_2^2}{\sigma^2} \right)}$$

\nwarrow
 $l_i(x)$

$$l_i(x) = \frac{w_i}{\sum_{i=1}^n w_i}$$

So, $\hat{y} = \sum_{i=1}^n l_i(x) y_i$

Thus, kernel regression is linear smoother.
because it is in the form of $l(x)^T y$.

Hence proved

P.T.O

Que 2) b)

To prove :-

If we fit a linear regression model by minimizing the sum of absolute values of residuals i.e. $\|Hw - Y\|_1$ instead of $\|Hw - Y\|_2^2$ then, it is not a linear

Smoother -

Proof - Minimizing L_1 norm is less stable than minimizing L_2 norm.

Example :-

Let, in a model if one point is having a error of 0 and other point is having a error of 10, if we take L_1 norm but if we take L_2 norm it will give an error of 100

→ If other model has 2 points and both having a error of 5 and 5. So, total error given by L_1 norm is 10 and total error by L_2 norm is $25 + 25 = 50$

So, if we take L_2 norm then it will suggest second model to be best, which is true because it is considering both the points and also penalising the the points having larger error.

→ Minimizing sum of absolute value of errors can be seen as finding median,

by
values of
of
linear

than

is
point is
L1 norm
give an

both

al

suggest
e

and
exger

lian,

where an optimal median, makes same number
of +ve errors and -ve errors
and median minimizes the L1 norm.

But, median is a non-linear function because,
it violates additivity.

Que2) c)

Given,

$$\hat{y} = \frac{1}{|B_k|} \sum_{i: x_i \in B_k} y_i \quad \text{--- (1)}$$

if we write the above eqn in the form
of Indicator Random variable.

→ Indicator random variable is

$$\begin{cases} 1, & \text{if } x_i \in B_k \\ 0, & \text{if } x_i \notin B_k \end{cases}$$

→ Also B_1, \dots, B_k are the bins

$$\Rightarrow \hat{y} = \frac{\sum_{i=1}^n I(x_i \in B_k) y_i}{\sum_{i=1}^n I(x_i \in B_k)} \quad \text{--- (2)}$$

eqn (1) in terms of Indicator
random variable.

Here,
$$l_i(x) = \frac{I(x_i \in B_k)}{\sum_{j=1}^n I(x_j \in B_k)}$$

as, in (a) part of this question we got to
know that linear smoother is defined

as
$$\hat{y} = \sum_{i=1}^n l_i(x) y_i$$

$$= l(x)^T y$$

as eqn ② can be written in the form

$$\sum_{i=1}^n y_i l(x_i)$$

$$\hat{y} = \frac{\sum_{i=1}^n y_i I(x_i \in B_x) l(x_i)}{\sum_{i=1}^n I(x_i \in B_x)}$$

Hence it is a linear smoother