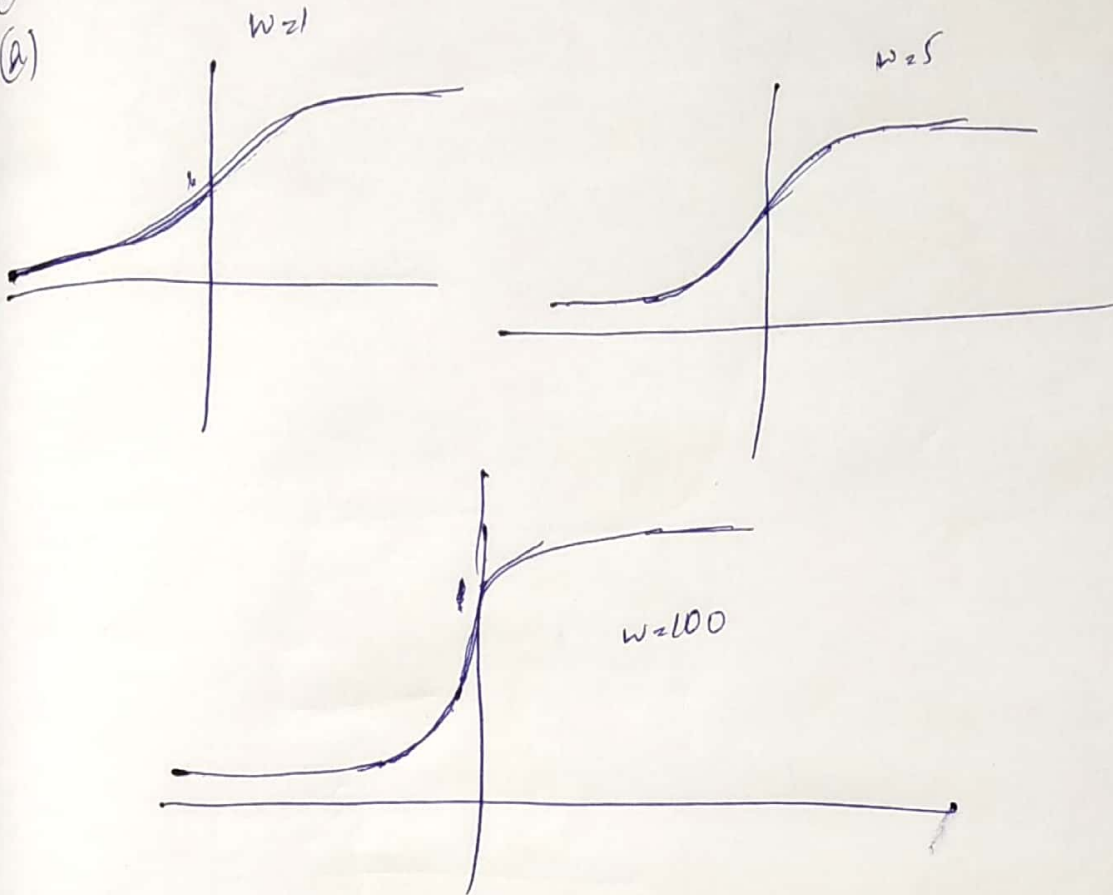


# Homework-3

G. Deepika  
CS16B16EN14015

(1)

(a)



The curve is getting steeper as  $w$  increases. A steeper curve means that the model is nearly completely sure of the class, with larger weights even small changes in input can lead to a large change in probability of class thus leading to easy flipping of the predicted output.  
(This is the reason we claim that it overfits)

(b) Let  $w = [w_0, w_1, \dots, w_d]$

$$\mathcal{L}(w) = \log(p(w)) - \sum_{j=1}^n \log(p(y^j | x^j, w))$$

$$p(w) = \prod_{i=0}^d \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{w_i^2}{2}\right)$$

$$\therefore \text{MAP estimate: } w^* = \arg \max_w \mathcal{L}(w) = \arg \max_w \left( \sum_{j=1}^n \log(p(y^j | x^j, w)) - \sum_i \frac{w_i^2}{2} \right)$$

$$\text{Gradient ascent update rule: } w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left. \frac{\partial \mathcal{L}(w)}{\partial w_i} \right|_t$$

$$\therefore \frac{\partial L(w)}{\partial w_i} = \left( \frac{\partial \log p(w)}{\partial w_i} \right) + \frac{\partial \log \left( \prod_{j=1}^n p(y^j/x^j, w) \right)}{\partial w_i}$$

↓

$$= -w_i$$

⇒ Final update rule is:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left( -w_i^{(t)} + \sum_j x_i^j (y^j - p(y=1/x^j, w^{(t)})) \right)$$

(c) Since all probabilities must sum to 1,

$$P(Y=y_k | X) = 1 - \sum_{k=1}^{K-1} P(Y=y_k | X)$$

We can define,

$$P(Y=y_k | X) = \frac{1}{1 + \sum_{k=1}^{K-1} \exp \left( w_{k0} + \sum_{i=1}^d w_{ki} x_i \right)}$$

and for  $k=1, \dots, K-1$

$$P(Y=y_k | X) = \frac{\exp \left( w_{k0} + \sum_{i=1}^d w_{ki} x_i \right)}{1 + \sum_{k=1}^{K-1} \exp \left( w_{k0} + \sum_{i=1}^d w_{ki} x_i \right)}$$

~~Classification rule~~

Classification rule simply picks the label with highest probability:

$$y = y_{k^*} \text{ where } k^* = \arg \max_{k \in \{1, \dots, K\}} P(Y=y_k | X)$$

(d) The decision boundary between each pair of classes is linear and hence the overall decision boundary is piece-wise linear. Equivalently, since  $\arg \max_i \exp(a_i) = \arg \max_i a_i$ , and  $\max$  of linear functions is piece-wise linear, the overall decision boundary is piece-wise linear.

(2)

(a)  $k(x_i, x) = \exp \frac{\|x_i - x\|_2^2}{\sigma^2}$  (input  $\bar{x}$ )

for linear smoother,

$$\begin{aligned} \hat{y} &= \frac{\sum_{i=1}^n k(x_i, x) \cdot y_i}{\sum_{i=1}^n k(x_i, x)} \\ \hat{y} &= l(x)^T \cdot y \end{aligned}$$

So, kernel regression is linear smoother.

(b) Proof by counter example:

Let's say input:  $x_i = k$  for different values of  $y$   
( $i = 1, 2, 3, \dots, n$ )

An optimal  $w$  makes the same number of +ve & -ve errors.

So that implies,  $w$  is the median of all  $y$

and  $w$  is not linear in any of  $y$ 's and the median changes as the value of  $y$  changes.

$\therefore w$  is not linear it doesn't fit a linear smoother



(c)  $\hat{y} = \frac{1}{|B_k|} \sum_{i: x_i \in B_k} y_i$  for  $x \in B_k$  where  $|B_k| = \text{no. of points in } B_k$

\* Regressogram is a linear smoother &

$$\hat{r}(x) = \frac{\sum_{i=1}^n I(x_i \in B(x)) y_i}{\sum_{i=1}^n I(x_i \in B(x))}$$

where  $|B_k| = \sum_{i=1}^n I(x_i \in B_k)$

$$\hat{r}(x) = \frac{I(x_i \in B_k)}{|B_k|}$$