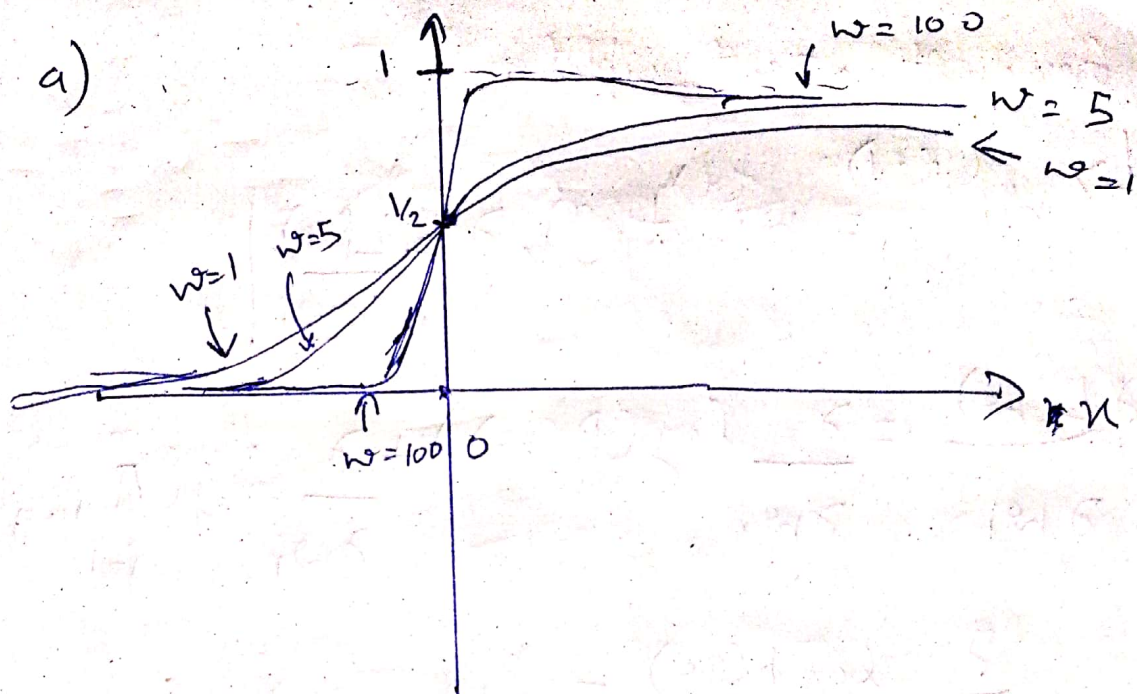A1 a)



As $w$ increases. the curve gets steeper, so the model would be almost sure of the class (with almost 1 probability). Even a little change in $x$ the class ~~probabilitie~~ probabilities may change by a lot (when $x$ near 0) when $w$ is large. So a solution with large weights can cause overfiting.

A1 (b) $W_{MAP} = \arg\max\limits_{w_0, \cdots, w_d} \prod\limits_{i=1}^{n} P(Y_i | X_i, w_0, \cdots w_d) \cdot P(w_0, \cdots, w_d)$

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

$w^* = \arg\max\limits_{w} \log(W_{MAP})$

$= \arg\max \left[ \sum\limits_{j=1}^{n} \log(P(y^j | x^j, w)) - \sum\limits_{i=0}^{d} \frac{w_i^2}{2} \right]$

Gradient ascent update

$\rightarrow \quad w_i^{(t+1)} = w_i^{(t)} + \eta \, \partial \frac{L(w)}{\partial w_i} \Big|_t$

$$\frac{\partial L(w)}{\partial w_i} = \frac{\partial}{\partial w_i} \log(P(w)) + \frac{\partial}{\partial w_i} \log\left(\prod_{j=1}^{n} P(y^j | n^j_{w})\right)$$

$$\frac{\partial}{\partial w_i} \log p(w) = -w_i$$

Update rule:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left(-w_i^{(t)} + \sum_j n_i^j \left(y^j - P(y=1 | n^j, w^{(t)})\right)\right)$$

A'(c) Probabilites of all the classes add up to 1 $\Rightarrow$

$$\Rightarrow P(y=y_k | x) = 1 - \sum_{k=1}^{k-1} P(Y=y_k | \alpha x)$$

$$P(Y=y_k | x) \propto \exp\left(w_{k0} + \sum_{i=1}^{d} w_{ki} X_i\right), \text{ for } k=1-k-1$$

$$P(Y=y_k | \alpha x) = \frac{1}{1 + \sum_{k=1}^{k-1} \exp\left(w_{k0} + \sum_{i=1}^{d} w_{ki} X_i\right)}$$
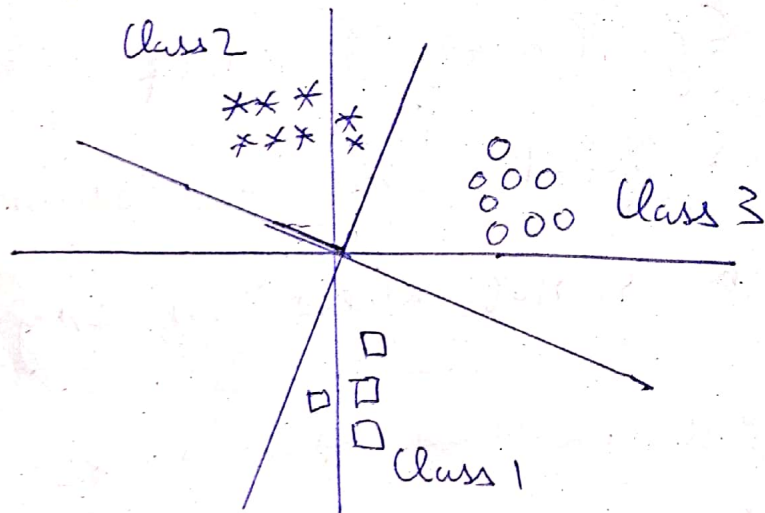
for $k=1, 2 \cdots \cdots k-1$

$$P(Y = y_k | X)$$

$$= \frac{\exp\left(w_{k_0} + \sum_{i=1}^{d} w_{ki} X_i\right)}{1 + \sum_{k=1}^{K-1} \exp\left(w_{k_0} + \sum_{i=1}^{d} w_{ki} X_i\right)}$$

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

A1d) The decision boundary between each class is linear $\Rightarrow$ overall decision boundary is piece wise linear.

A2(a) $\hat{y} = \bar{y} (G + dI)^{-1} \bar{z}$

where $z = \langle \phi(u), \phi(\bar{u}) \rangle = K(u, \bar{u})$

$z$ is a function of $\bar{u}$

$(G + dI)^{-1}$ is not a function of $\bar{u}$

$l(\bar{u}) = (G + dI) \bar{z}$

$\hat{y} = \bar{y}^T \cdot l(u)$ which is a dot product hence a linear smoother.

A2 b) No, it is not a linear smoother.

Counter Eg: Constant input where each point has $x_i = 1$. '$w$' is median of $y$'s.

Clearly, $w$ is not linear in any of $y$'s, the median changes as the rank of $y$'s does.

A2(c) Yes, it's a linear smoother.

$$\ell(n): \ell(\bar{u})_i = \begin{cases} \dfrac{1}{|B_k|} & \text{if } u_i \in B_k \\ 0 & \text{else} \end{cases}$$

Range is divided into bins.

$\ell(u) \cdot \bar{y} \Rightarrow$ Output the avg of value of the bin that $n$ belongs.

Hence this a linear smoother.

γ.

each

*y's.

y's,

rank