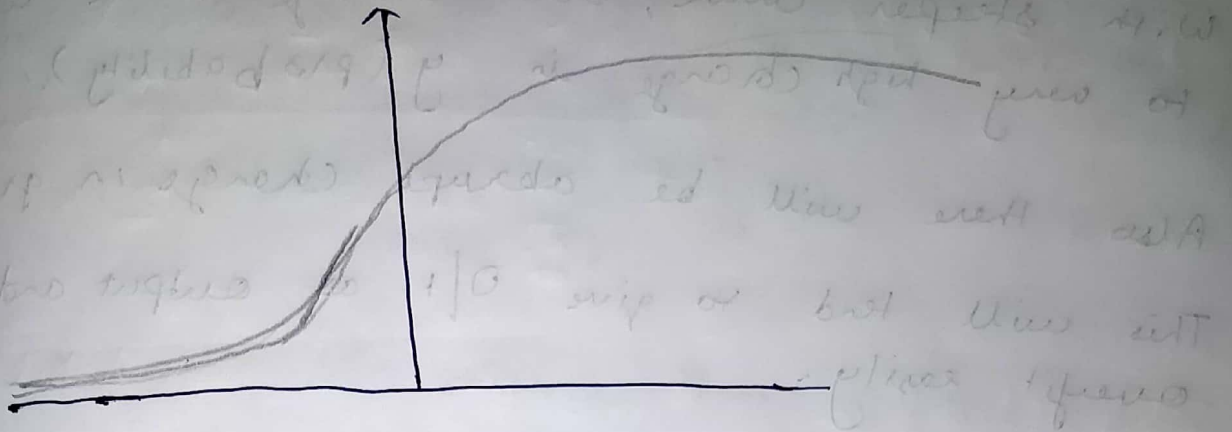
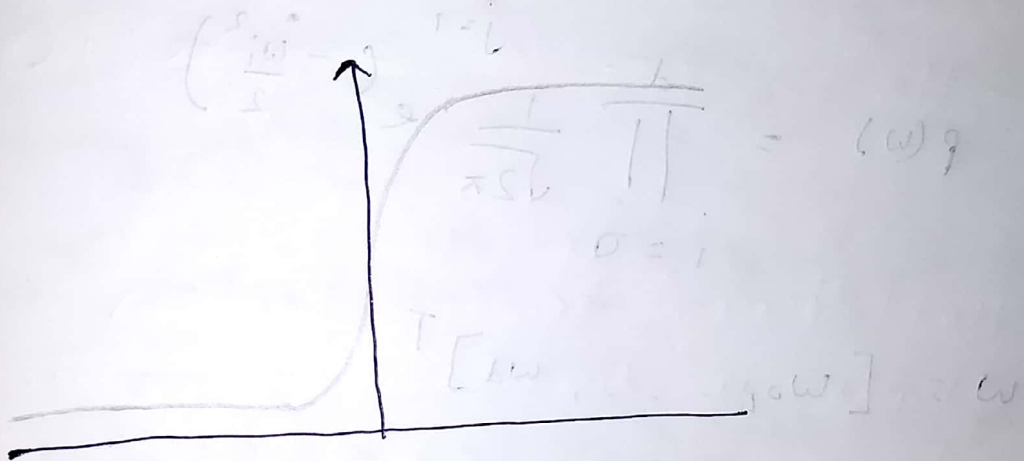


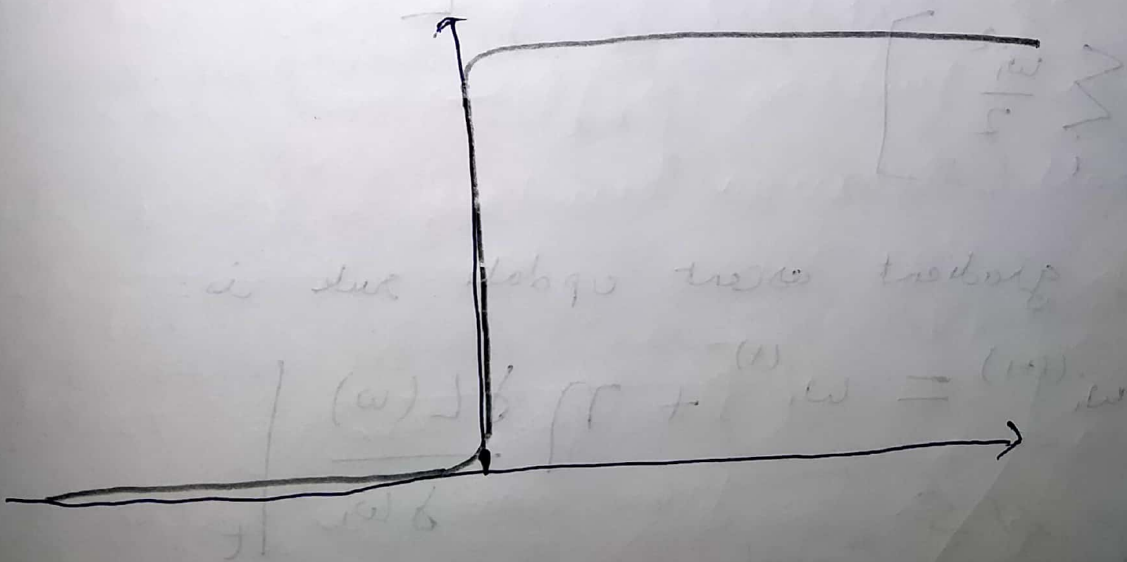
(a)  $\omega = 1$



(b)  $\omega = 5$  (d)  $L(\omega) = \log \left( \prod_{i=1}^n (1 + \frac{\omega^2}{\omega_i^2}) \right)$



(b)  $\omega = 100$  (d)  $L(\omega) = \log \left( \prod_{i=1}^n (1 + \frac{\omega^2}{\omega_i^2}) \right)$



As  $w$  increases, the curve gets steeper.

With steeper curve, same change in  $x$  will lead to very high change in  $y$  (probability).

Also there will be abrupt change in gradients.

This will tend to give 0/1 as output and will overfit easily.

$$(b) \quad L(w) = \log(p(w) \prod_{j=1}^n P(y^j/x^j, w))$$

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

$$w = [w_0, \dots, w_d]^T$$

MAP estimate is

$$w^* = \arg \max_w L(w) = \arg \max_w \left[ \sum_{j=1}^n \log(P(y^j/x^j, w)) - \sum_i \frac{w_i^2}{2} \right]$$

The gradient ascent update rule is:

$$w_i^{(t+1)} = w_i^{(t)} + \eta \left. \frac{\partial L(w)}{\partial w_i} \right|_t$$

The gradient of the log conditional posterior is:

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

2nd term in above eqn is:

For unregularized

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

Thus finally update rule

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

$$(c) P(Y = y_k | X) \propto \exp(w_{k0} + \sum_{i=1}^d w_{ki} x_i) \text{ for}$$

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

$P(Y = y_k | X)$  is for same model

Since all the probabilities sum up to 1

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

Adding another set of weights is redundant

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



and for  $k = 1, 2, \dots, K-1$

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

The classifier will simply pick the label with highest probability

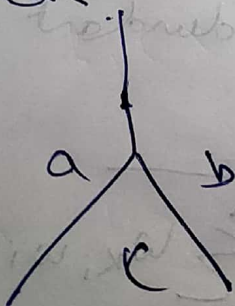
$$y = y_{k^*} \quad \text{where } k^* = \operatorname{argmax}_k (P(Y = y_k | X))$$

$k \in \{1, \dots, K\}$

(d) The decision boundary between each pair of classes is linear and hence overall decision boundary is piece wise linear

Equivalently since

$\operatorname{argmax}_i \exp(a_i) = \operatorname{argmax}_i (a_i)$ , and max of linear function is piece wise linear the overall (decision) boundary is piece wise linear



(Decision boundary resulting from multiclass logistic regression)

2 (a) The given kernel function is a semi gaussian kernel function

kernel regression

$$\hat{y} = \sum_{i=1}^n w_i y_i \quad \Bigg| \quad \sum_{i=1}^n w_i \quad \text{Avg weight}$$

Now, for a point  $x_i$ , the  $w_i$  will be replaced by our kernel

$$w_i = \exp \left( \frac{||x_i - x||^2}{\sigma^2} \right)$$

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

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

The given kernel regression is linear smoother

(b) We are fitting a linear regression model

sum of  $||Hw - y||^2$ , we know absolute value of residuals

Proof of not a linear smoother.

There is no closed form solution for  $w$  that minimizes the sum of absolute value of error

Yet, solution can be seen to be similar to median. An optimal  $w$  make the same



number of positive & negative error

Counter example

Consider a constant input where each training point has  $x_i = 1$  for different  $y$  value so  $w$  is the median of all  $y$ .

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

(c) Dividing  $a, b$  in  $m$  equally spaced bins  $B_1, B_2, \dots, B_m$

$$\hat{y} = \frac{1}{|B_k|} \sum_{i: x_i \in B_k} y_i$$

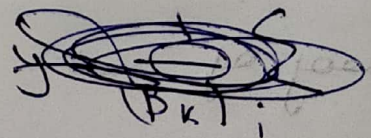
$|B_k|$  = no. of points in  $B_k$

Definition of linear smoother  $[\hat{y} = L(X^T | y)]$

we have  $\hat{y} = \frac{1}{|B_k|} \sum_{i: x_i \in B_k} y_i$

By analysing func. we have

$$Lg(x) = \frac{\sum_{i: x_i \in B_i} y_i}{|B_i|}$$



Then linear smoother condition is satisfied.

Hence our regression is linear smoother.

$$Lg(x) = \frac{\sum_{i: x_i \in B_k} y_i}{|B_k|}$$