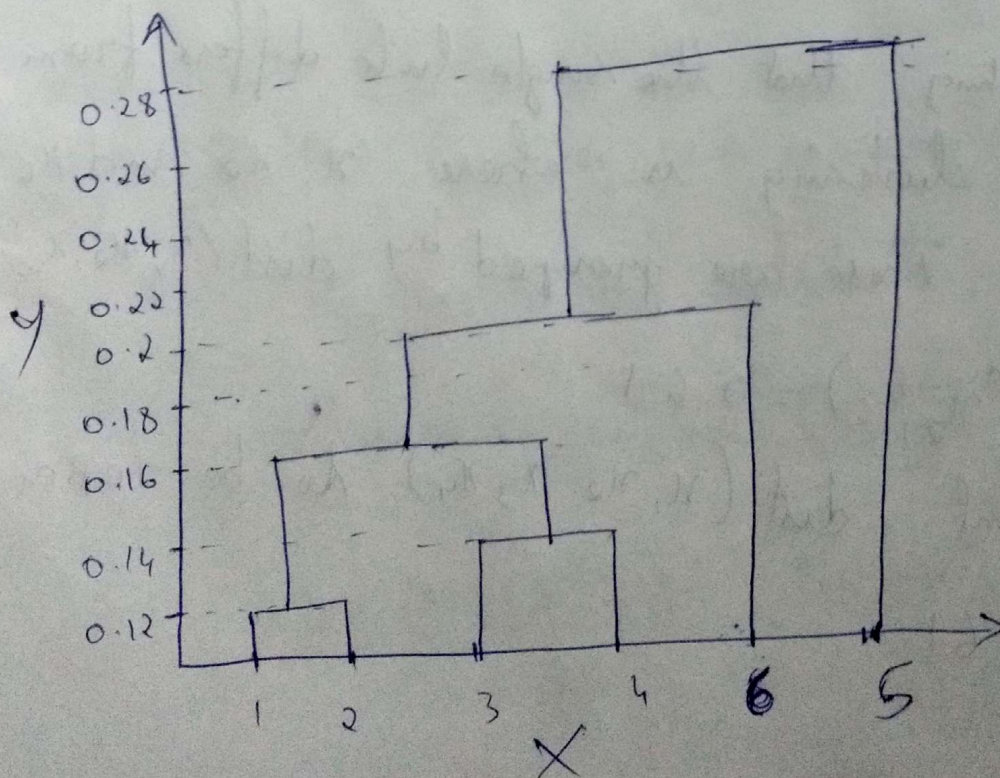


1) Hierarchical clustering

Distance matrix is as follows

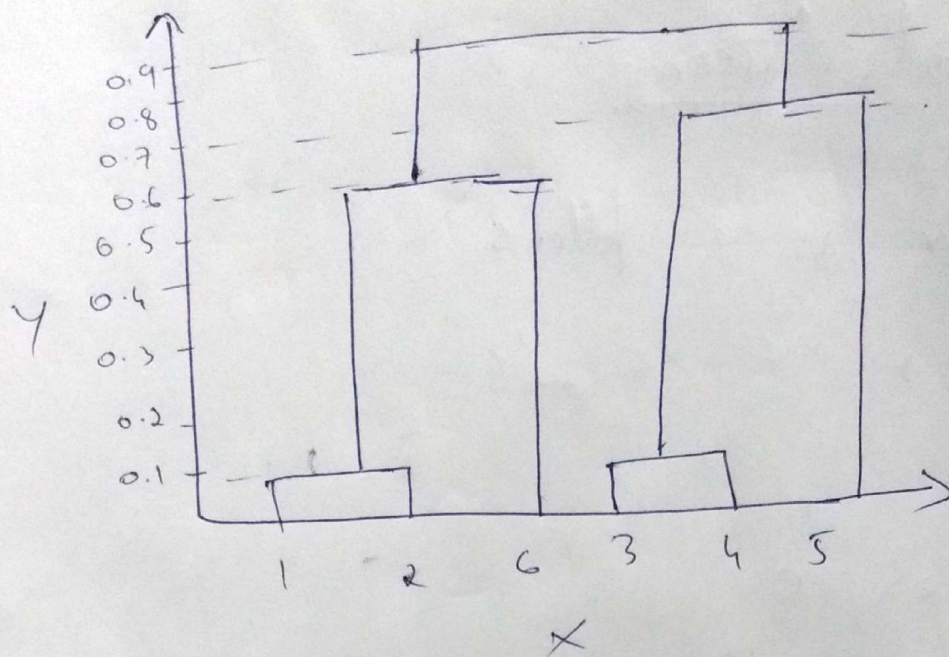
	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
$x_1$	0					
$x_2$	0.12	0				
$x_3$	0.51	0.25	0			
$x_4$	0.84	0.16	0.14	0		
$x_5$	0.28	0.77	0.70	0.45	0	
$x_6$	0.34	0.61	0.93	0.2	0.67	0

a) Constructing Dendrogram:





b) Dendrogram for final result of hierarchical clustering with complete link



c) We notice that only 2 values change, so the answer remains the same, i.e. it is same as previous answer above.

Now,

The first thing that the single link differs from complete link clustering is where  $x$ ,  $x_2$  and  $x_6$  are grouped. These are grouped by  $\text{dist}(x, x_2, x_6)$

$$\Rightarrow \text{dist}(x_2, x_6) = 0.61$$

$\therefore$  We want  $\text{dist}(x, x_2, x_3, x_4)$  to be lesser than 0.61.



$\Rightarrow$  Then we want

$$\text{dist}(x_1, x_2, x_3, x_4, x_6) = \text{dist}(x_3, x_6)$$

$$= 0.93$$

to be lesser than this so that

$x_1, x_2, x_3, x_4$  and  $x_6$  are grouped together.

After these changes both become identical.

$$2) \text{ Covariance Matrix} = \frac{1}{n} \left[ (E[x] - \bar{x})(E[x] - \bar{x})^T \right]$$

let  $n = (\text{some } k) \times m$  (multiple)

$\Rightarrow C$  is

$$1) E[a]^2 \frac{(M-1)^2}{M^2} \text{ for } (i, i).$$

$$2) E[a]^2 \frac{(M-1)}{M^2} \quad \forall (x, i) \text{ or } (i, x) \quad x < n \text{ samples.}$$

$$3) E[a]^2 / M \quad \text{otherwise}$$

Taking  
if we  
get

$\Rightarrow$  (a) Consider non-diagonal elements.

$$2 \text{ sample will contribute} = E[a]^2 \frac{(M-1)}{M^2}$$

$$\text{others will contribute (3)} \rightarrow \frac{E[a]^2}{M}$$

$$\Rightarrow \left( \frac{M-2}{M^2} \right) E[a]^2 - \frac{2 E[a]^2 (M-1)}{M^2}$$

$$\Rightarrow - E[a]^2 / M$$



### Diagonal elements

one will contribute to (1) for given  $(i, i)$   
rest of them  $\rightarrow$  other case

$$\Rightarrow \frac{M-1}{M^2} f(a)^2 - \frac{f(a)^2 (M-1)^2}{M^2}$$

$$= f(a)^2 - \frac{f(a)^2}{M}$$

$$\therefore C_{ij} = \begin{cases} -\frac{1}{M} f(a)^2, & i \neq j \\ \left(\frac{M-1}{M}\right) f(a)^2, & i = j \end{cases}$$

b) We can verify that

$$T = (1, 1, \dots, 1)$$

is given eigen vector of

$$C \times T = \left(\frac{M-1}{M}\right) f(a)^2 + \frac{(M-1) f(a)^2}{M^2}$$

$\rightarrow$  Let difference b/w diagonal and non-diagonal be  $\delta$ .

$$\therefore \text{Rank}(C - \delta I) = 1$$

$$\lambda(\text{geometrical}) = M - \text{rank}(C - \delta I)$$

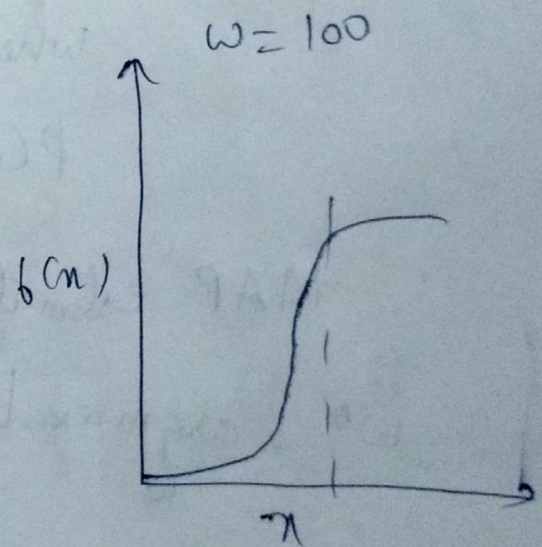
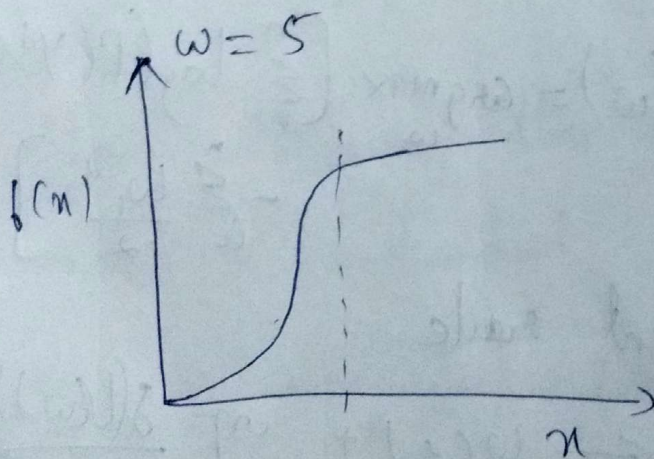
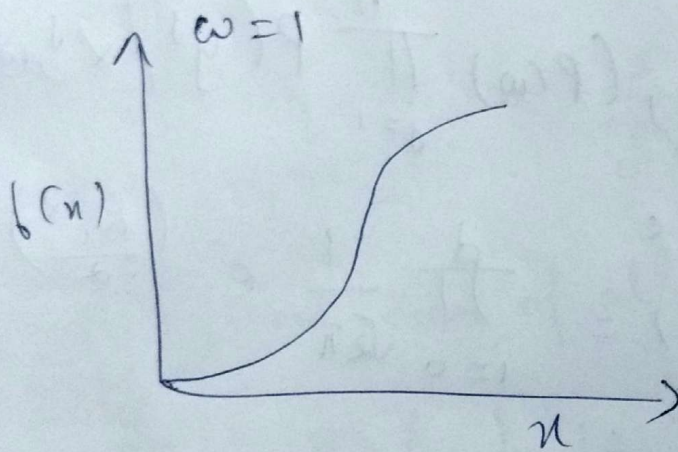


Hence geometric multiplicity of  $\delta = M-1$ ,

$\therefore$  Rest of  $M-1$  eigen vectors have eigen values  $\delta$ .

c) PCA is not suited as vectors  $\delta$  to be picked by Nltk, can have only linear relation between them.

3) (a)  $f(n) = \frac{1}{1 + e^{-\omega n}}$



We can observe that as  $\omega$  increases  $f(n)$  gets steeper and the curve gets more steeper. It means model is completely sure of the class.



As heights get high, with large heights, small changes in  $n$  can lead to large changes in probability leading to misclassification and overfitting.

$$b) \quad \omega = [\omega_0, \omega_1, \dots, \omega_d]^T$$

We take here log conditional posterior instead of log conditional likelihood.

$$\text{i.e.} \quad L(\omega) = \log(P(\omega)) \prod_{j=1}^n P(y^j / x^j, \omega)$$

$$\text{where } P(\omega) = \prod_{i=0}^d \frac{1}{\sqrt{2\pi}} e^{-\frac{\omega_i^2}{2}}$$

$\therefore$  MAP estimate is

$$\omega^* = \underset{\omega}{\operatorname{argmax}} L(\omega) = \underset{\omega}{\operatorname{argmax}} \left[ \sum_{j=1}^n \log(P(y^j / x^j, \omega)) - \sum_{i=0}^d \frac{\omega_i^2}{2} \right]$$

$\therefore$  The gradient ascent rule

$$\Rightarrow \omega_i(t+1) \leftarrow \omega_i(t) + \eta \left. \frac{\partial L(\omega)}{\partial \omega_i} \right|_t$$

For log conditional posterior, it is

$$\frac{\partial L(\omega)}{\partial \omega_i} = \frac{\partial}{\partial \omega_i} \log P(\omega) + \frac{\partial}{\partial \omega_i} \log \left( \prod_{j=1}^n P(y^j / x^j, \omega) \right)$$



Here the second term is same as derived in before for unregularized case.

First term leads to extra factor of

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

Final update is as follows

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \left( -w_i^{(t)} + \sum_{j=1}^n x_j^i (y_j - P(y=1 | x, w^{(t)})) \right)$$

(c) We know sum of all probabilities = 1

$$\therefore P(y = y_k | x) = 1 - \sum_{n=1}^{K-1} P(y = y_n | x)$$

Here introducing this set of weights will make it redundant.

$$\therefore P(y = y_k | x) = \frac{e^{w_{k0}} + \sum_{i=1}^d w_{ki} x_i}{\left( 1 + \sum_{k=1}^{K-1} e^{w_{k0}} + \sum_{i=1}^d w_{ki} x_i \right)}$$



2) Classification rule simply picks ~~the~~ label with highest probability.

$$y = y_k, \quad k^* = \underset{k}{\operatorname{argmax}} P(y = y_k / X)$$