1) a)

$\frac{1}{1+e^{-x}}$ 

$\frac{1}{1+e^{-5x}}$ 

$\frac{1}{1+e^{-100x}}$ 

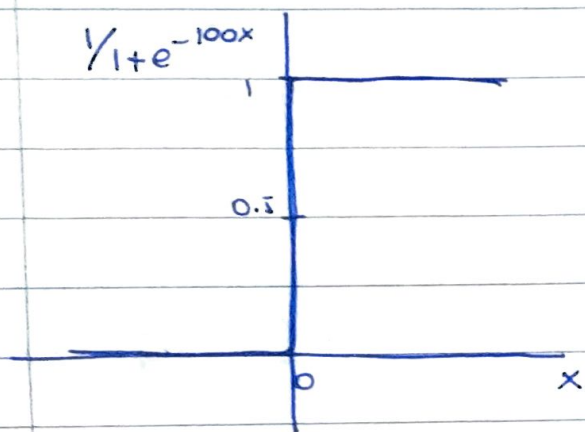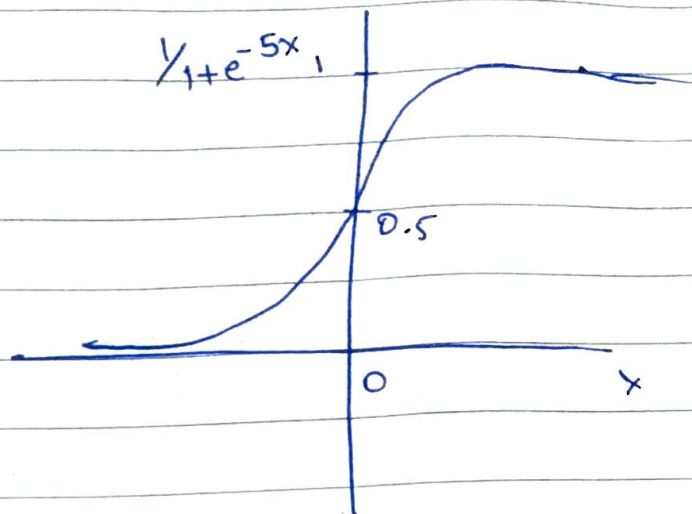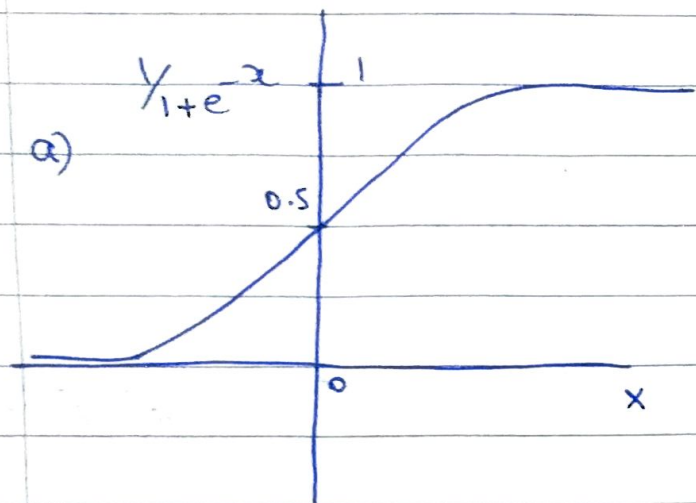As $\omega$ increases, the curve becomes more ||el to y-axis and begins to look like step function.

The reason why it overfits as $\omega$ increases is because with large value of $\omega$, small change in input leads to change in output label. ~~The point of probability never tends to affect the class.~~

1) (b)  MLE $\Rightarrow$  $w = \max\limits_{w_0,..w_d} \prod\limits_{i=1}^{n} P\left(Y_i / X_i, w_0, ... w_d\right)$

MAP $\Rightarrow$  $w = \max\limits_{w_0,...w_d} \prod\limits_{i=1}^{n} P\left(Y_i / X_i, w_0, ... w_d\right) P\left(w_0, ... w_d\right)$

Assume  Standard Gaussian prior $N(0, I)$.  Gradient Ascent?

let  $w = [w_0, ... ... w_d]$

for Gaussian distribution wkt
$$P(X) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(x - \mu/\sigma)^2}$$

$P(w) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(w - \mu/\sigma)^2}$   for  $N(0, I)$.

$\quad = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}w^2}$

$\dfrac{\partial \ln(P(w))}{\partial w} = -w$

from MLE

$$L(w) = \sum\limits_{i} \left[ y^i \ln P(y^i = 1 / x_i, w) + (1 - y^i) \ln P(y^i = 0 / x^i, w)\right]$$

$$+ \ln P(w)$$

$$\frac{\partial}{\partial w} L(w) = \sum\limits_{i} \left[ \frac{y^i}{\sigma(w^T x_i)} \sigma(w^T x_i)(1 - \sigma(w^T x_i)) \cdot x_i \right.$$
$$\left. - \left(\frac{1 - y_i}{1 - \sigma(w^T x_i)} (1 - \sigma(w^T x_i)) \cdot \sigma(w^T x_i) - x_i \right) \right]$$
$$- w$$

$$\frac{\partial}{\partial w} L(w) = \sum\limits_{i} x_i \left[ y_i - P(Y = 1 / X, w)\right] - w$$

Gradient Ascent formula:  $w_{t+1} \leftarrow w_t + \eta \left[ \frac{\partial}{\partial w} L(w)\right]$

$$w_{t+1} \leftarrow w_t + \eta \left[ -w_t + \sum\limits_{i} x_i \left[ y_i - P(Y = 1 / X, w)\right]\right]$$

1) (i) $\Rightarrow P(Y = y^k / X) \propto \exp\left( w_{k_0} + \sum_{i=1}^{d} w_{k_i} X_i \right)$

From odd's ratio wkt $\ln\left( \frac{Pr(Y_i = 1 / X_i)}{Pr(Y_i = K / X_i)} \right) = w_1 x_i$ where K is set of classes excluding 1.

$\ln\left( \frac{Pr(Y_i = 2 / X_i)}{Pr(Y_i = K / X_i)} \right) = w_2 x_i$

⋮

In General $\ln\left( \frac{Pr(Y_i = k / X_i)}{Pr(Y_i = K / X_i)} \right) = w_k x_i$

$\Rightarrow Pr(Y_i = k / X_i) = Pr(Y_i = K / X_i) \cdot e^{w_k x_i}$

$Pr(Y_i = K / X_i) = 1 - \sum_{R=1}^{K-1} Pr(Y_i = R / X_i)$

$= 1 - \sum_{k=1}^{K-1} Pr(Y_i = K / x_i) \cdot e^{w_k x_i}$

$= 1 - \sum_{k=1}^{K-1} \frac{e^{w_k x_i}}{1 + \sum_{R=1}^{K-1} e^{w_R x_i}}$ 　　[Expand $Pr(Y_i = K / x_i)$]

$= 1 - \frac{\left[ e^{w_1 x_i} + e^{w_2 x_i} + \cdots + e^{w_{K-1} x_i} \right]}{1 + \left[ e^{w_1 x_i} + e^{w_2 x_i} + \cdots + e^{w_{K-1} x_i} \right]}$

$Pr(Y_i = K / x_i) = \frac{1}{1 + \sum_{R=1}^{K-1} e^{w_R x_i}}$

$Pr(Y_i = k / X_i) = \frac{e^{w_k x_i}}{1 + \sum_{k=1}^{K-1} e^{w_R x_i}}$ 　　[combined $w_p + \sum_{i=1}^{d} w_{k_i} x_i$ to $w_k x_i$]

The classification rule implies that we have the label with the highest probability.

$y = y_k^* \quad$ where $k^* = \underset{R \in \{1 \cdots K\}}{\arg\max} \ P(Y = y_k / X)$

1) (d)

label.

Training data