

CS6140 Assignment 3

1.4 ① $\sigma(a) = \frac{1}{1 + e^{-w^T x}}$

$$\frac{\partial \sigma(x)}{\partial x} = \frac{\partial \sigma}{\partial x} \left((1 + e^{-w^T x})^{-1} \right)$$

$$= \frac{\partial \sigma}{\partial x} \left(- \frac{1}{(1 + e^{-w^T x})^2} \right) (-w^T e^{-w^T x})$$

$$\boxed{\frac{\partial \sigma(x)}{\partial x} = \frac{w^T e^{-w^T x}}{(1 + e^{-w^T x})^2}}$$

② $P(y=1|x, w) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}}$

$$P(y=-1|x, w) = \sigma(w^T x) \stackrel{?}{=} \frac{1}{1 + e^{w^T x}}$$

$$P(y=-1|x, w) = 1 - P(y=1, x, w)$$

$$= 1 - \frac{1}{1 + e^{-w^T x}}$$

$$= \frac{1 + e^{-w^T x} - 1}{1 + e^{-w^T x}}$$

$$= \frac{e^{-w^T x}}{1 + e^{-w^T x}} \times \frac{1/e^{-w^T x}}{1/e^{-w^T x}}$$

$$= \frac{1}{\frac{1}{e^{-w^T x}} + 1} = \frac{1}{1 + e^{w^T x}} \quad \checkmark$$

$$\boxed{\therefore P(y=\pm 1|x, w) = \sigma(w^T x) = \frac{1}{1 + e^{-y w^T x}}}$$

$$3) L_{\log} = \sum_{i=1}^N \log(1 + e^{-y_i w^T x_i})$$

$$\begin{aligned} L_{\log} &= \sum_{i=1}^N y_i \log(\sigma(w^T x_i)) + (1 - y_i) \log(1 - \sigma(w^T x_i)) \\ &= \sum_{i=1}^N y_i \left(\log\left(\frac{1}{1 + e^{-y_i w^T x_i}}\right) + \log\left(\frac{1}{1 + e^{y_i w^T x_i}}\right) - y_i \left(\log\left(1 - \frac{1}{1 + e^{y_i w^T x_i}}\right) \right) \right) \\ &\quad \text{from part 2} \quad \nearrow \\ &= \sum_{i=1}^N -y_i \left(\log(1 + e^{-y_i w^T x_i}) \right) + \log(1 + e^{-y_i w^T x_i}) + y_i \left(\log(1 + e^{y_i w^T x_i}) \right) \end{aligned}$$

$$L_{\log} = \sum_{i=1}^N \log(1 + e^{-y_i w^T x_i}) \quad \checkmark$$

if (y_i) & $(w^T x_i)$ have the same sign, the error recorded (log loss) is small and close to 0, causing little correction during gradient descent. If their signs are different, error (log loss) is high, and a larger correction is made to the weight creating a model during future iterations with less log loss.

2.6

$$① p(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

Bernoulli likelihood:

$$L(\pi, \theta | D) = \prod_{i=1}^N \prod_{k=1}^K \prod_{j=1}^M (\pi_k \theta_{jk}^{i,j} (1-\theta_{jk}^{k,j})^{y_{i,j,k}})$$

$$L(\theta) \times p(\theta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1} \times \left(\prod_{i=1}^N \prod_{k=1}^K \prod_{j=1}^M \pi_k \theta_{jk}^{i,j} (1-\theta_{jk}^{k,j})^{y_{i,j,k}} \right)$$

$$\frac{\partial (L(\theta) \times p(\theta))}{\partial \theta} = 0 = \frac{\partial}{\partial \theta} \left(\downarrow \right)$$

I have no idea how to solve this...
but solving for θ , will give us the MAP

$$(2.6) (2) \quad P(\theta) = \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{j=1}^K \theta_j^{\alpha_j - 1}$$

$$L(\theta) = \prod_{i=1}^N \prod_{k=1}^C \prod_{j=1}^K (\pi_k \theta_{jk}^{i_j} (1 - \theta_{jk}^{x_{ij}}))^{y_{ijk}}$$

$$\frac{\partial (L(\theta) \times P(\theta))}{\partial \theta} = \frac{\partial}{\partial \theta} \left(\prod_{i=1}^N \prod_{k=1}^C \prod_{j=1}^K (\pi_k \theta_{jk}^{i_j} (1 - \theta_{jk}^{x_{ij}}))^{y_{ijk}} \right) \times \frac{\Gamma(\sum_{i=1}^K \alpha_i)}{\prod_{i=1}^K \Gamma(\alpha_i)} \prod_{j=1}^K \theta_j^{\alpha_j - 1} = 0$$

Solve for θ for MAP estimate.

(2.6)

(3) The only thing that differs between MLE and MAP estimates is the inclusion of the prior, $P(\theta)$ in MAP. This means that the MAP is now weighted with some prior.

If the prior is uniform, NOT a beta or Dirichlet distribution, then we could ignore the constant and $\theta_{MAP} = \theta_{MLE}$.