Problem 2:

a)

$$\nabla l(\theta) = \sum_{i=i}^{n} \{\frac{exp(-\theta x^i)x^i}{1 + exp(-\theta x^i)} + x^i(y^i - 1)\}$$

Pseudo-Code

---
**Algorithm 1** GradientDescent1
---
1: Initialize parameter $\theta^0$
2: **while** $t <$ number of iteration **do**
3:     Calculate $\nabla\theta$ using (2)
4:     Update $\theta$ using $\theta^t = \theta^{t-1} + \eta\nabla\theta$
5: **end while**
---

b) Stochastic Gradient Descent

$$\nabla\theta = \sum_{i \in S_k} \{\frac{exp(-\theta x^i)x^i}{1 + exp(-\theta x^i)} + x^i(y^i - 1)\}$$

---
**Algorithm 3** StochasticGradientDescent
---
1: Initialize parameter $\theta^0$ and the difference $\epsilon$
2: **while** $\epsilon <$ threshold **do**
3:     Randomly sample 10% of the data without replacement
4:     Calculate $\nabla\theta$ using (3)
5:     Update $\theta$ using $\theta^t = \theta^{t-1} + \eta\nabla\theta$
6:     Calculate the old objective value using $\theta^{t-1}$ and new objective value using $\theta^t$, and do a difference to update $\epsilon$
7: **end while**
---

Hessian Matrix

$$\nabla^2 f(x) = \begin{bmatrix} \dfrac{\partial^2 f(x)}{\partial x_1^2} & \dfrac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \dfrac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \dfrac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \dfrac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \dfrac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \dfrac{\partial^2 f(x)}{\partial x_n \partial x_1} & \dfrac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \dfrac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}$$

Second Derivative of Log-likelihood function

$$\frac{\partial^2 l(\theta)}{\partial \theta_p \partial \theta_q} = \sum_{i=1}^{n} -\frac{x_p^{(i)} x_q^{(i)} \exp(\theta X^{(i)})}{(1 + \exp(\theta X^{(i)}))^2}$$

$$= \sum_{i=1}^{n} -x_p^{(i)} x_q^{(i)} S(-\theta X) S(\theta X)$$

$$= -X X^T S(-\theta X) S(\theta X)$$

Where

$$S = \frac{1}{1+\exp(-x)} = \frac{\exp(x)}{\exp(x)+1}$$

As S is the sigmoid function, S is always positive, as is $XX^T$. Therefore, the training problem is concave. Thus, it does not has local minima, and with sufficiently small training rate, the gradient descent will converges to global minima.