

Quiz 2

1.

The cross entropy loss function is $L(w) = -\frac{1}{n} \sum_{i=1}^n y_i \log(g_i) + (1-y_i) \log(1-g_i)$ where $g_i = f(w^T x_i) = 1/(1+e^{-w^T x_i})$ for $(x_1, y_1) \dots (x_n, y_n)$ in $\mathbb{R}^n \times \{0, 1\}$.

To show that this is convex, set $t = w^T x_i$.

$$\frac{\partial}{\partial t} g_i = \frac{\partial}{\partial t} (1+e^{-t})^{-1} = e^{-t} (1+e^{-t})^{-2} = g_i (1-g_i)$$

$$\frac{\partial \log(g_i)}{\partial w^T} = \frac{1}{g_i} \frac{\partial g_i}{\partial w^T} = \frac{1}{g_i} \frac{\partial g_i}{\partial t} \frac{\partial t}{\partial w^T} = (1-g_i) x_i$$

$$\frac{\partial \log(1-g_i)}{\partial w^T} = \frac{1}{1-g_i} \frac{\partial (1-g_i)}{\partial w^T} = -g_i x_i$$

$$\text{Summation component } l_i(w) = -y_i \log(g_i) - (1-y_i) \log(1-g_i)$$

$$\nabla l_i(w) = -y_i x_i (1-g_i) + (1-y_i) x_i g_i = x_i (g_i - y_i)$$

$$\nabla^2 l_i(w) = x_i x_i^T g_i (1-g_i) = \frac{1}{n} \sum_{i=1}^n \nabla^2 l_i(w) = \frac{1}{n} \sum_{i=1}^n x_i x_i^T g_i (1-g_i) \equiv X D X^T$$

where D is a diagonal matrix with all entries $D_{ii} = g_i (1-g_i) > 0$.

So $\nabla^2 L(w)$ is positive semidefinite and $L(w)$ is convex.

$L(w)$ is not strongly convex. So min of a convex function must achieve global min.

For gradient descent method $x_{k+1} = x_k - \eta_k \nabla L_{B_k}(x_k)$, $w_{k+1} = w_k - \eta_k \nabla L(w_k)$.

$L(w) \in C^2$, $|x_i| \leq M$ for all x_i of $\nabla^2 L(w)$. Use backtracking line search to find η_k . # iterations = $O(1/\epsilon)$ for gradient descent to reach $|L(w_k) - L(w^*)| < \epsilon$.

Adding an L_2 regularization term we get cross-entropy loss $L(w)$ due to convexity. So convergence rate becomes $O(\log(1/\epsilon))$ for reaching $|L(w_k) - L(w^*)| < \epsilon$.

2. (b) With $L(w) = \frac{1}{n} \sum_{i=1}^n l_i(w)$,

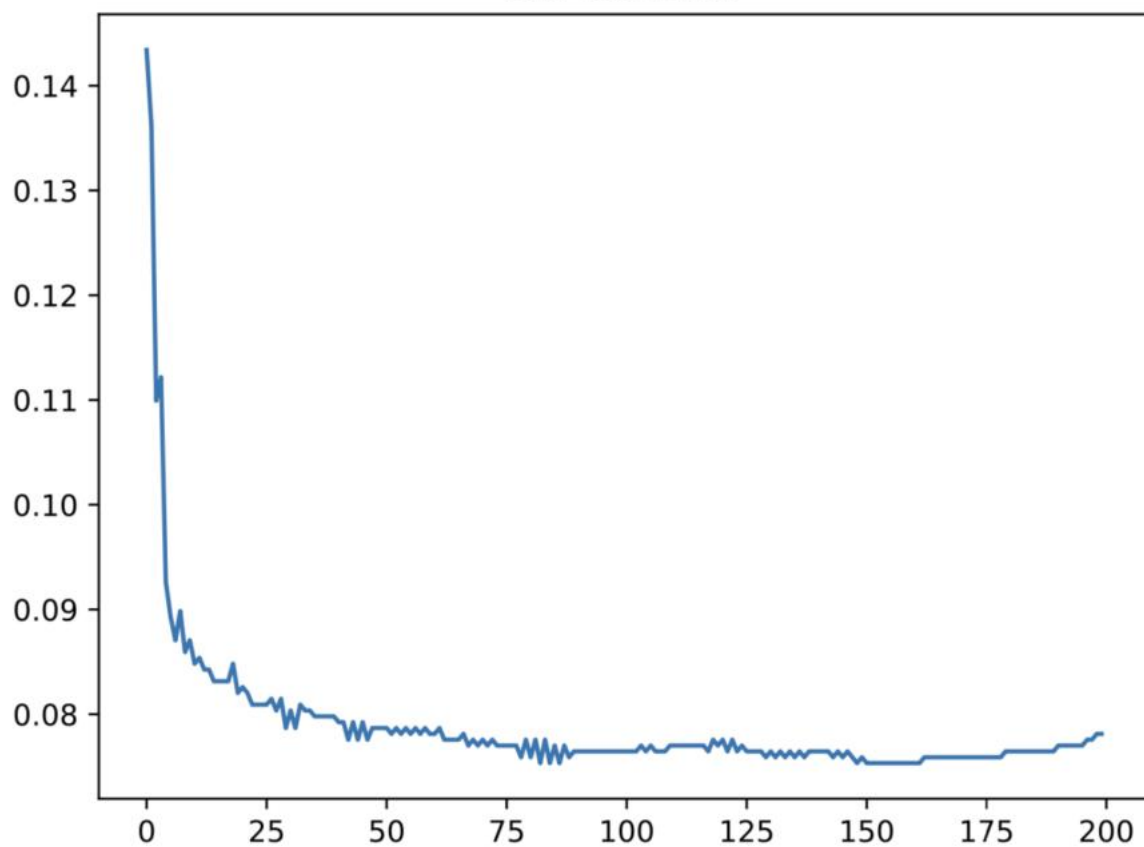
$$\nabla l_i(w) = -y_i x_i (1 - g_i) + (1 - y_i) x_i g_i = x_i (g_i - y_i)$$

$$= \frac{1}{n} \sum_{i=1}^n x_i (g_i - y_i) \text{ where } g_i = \frac{1}{1 + e^{-w^T x_i}} - 1$$

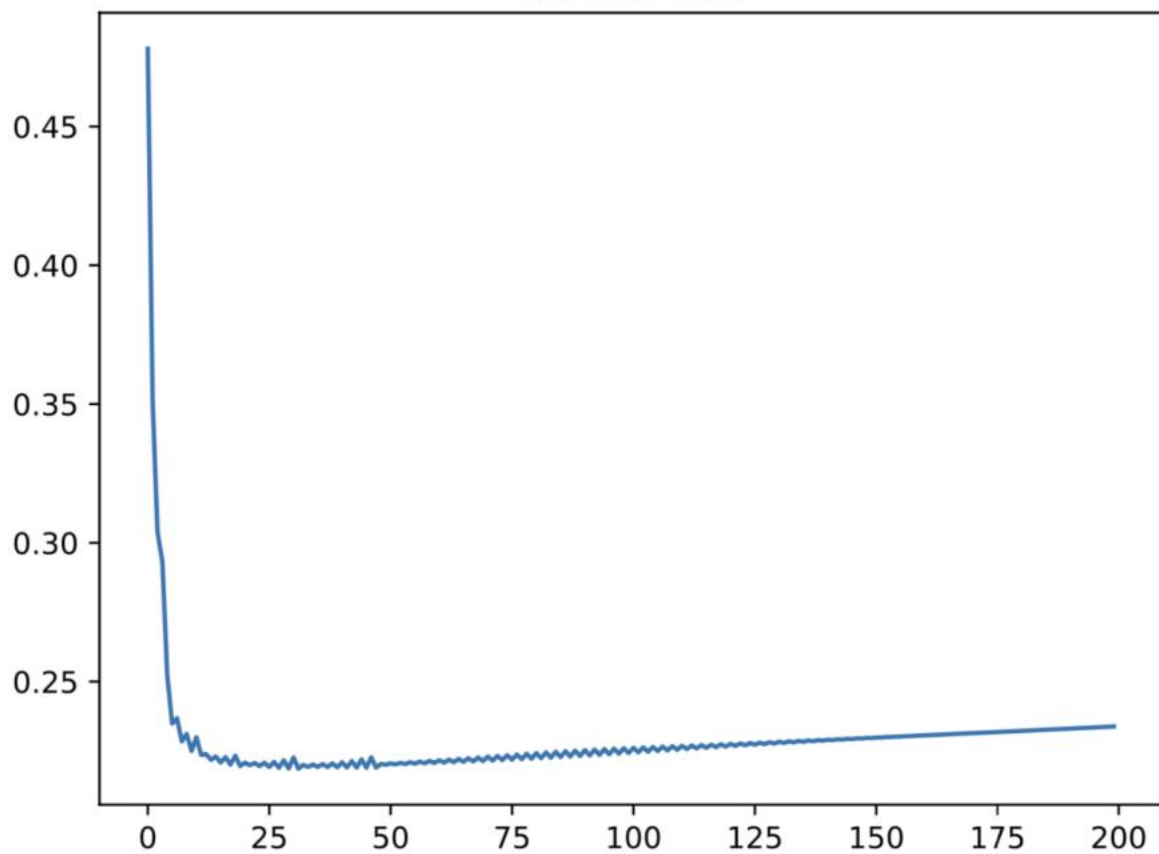
(c) From the ~~at~~ output data, the resulting loss ~~stable~~ becomes stable or converges at ~~the~~ roughly 50 iterations. This is in accordance with $O(1/\epsilon)$ for $|L(w_k) - L(w^*)| < \epsilon$.

Problem 2 Outputs:

GD Test Error



GD Test Loss



GD Training Loss

