

1.

1) 2-layer linear neural network :  $\hat{y} = f_w(x) = w^T x$

$$\text{cost function: } J(w) = \frac{1}{2n} \sum_{i=1}^n (\hat{y}^{(i)} - y^{(i)})^2$$

$$\text{Newton's Method: } x^* = x^{(0)} - H(f)(x^{(0)})^{-1} \nabla_w f(x^{(0)})$$

Replace  $x$  with  $w$ :  $w^* = w^{(0)} - H(f)(w^{(0)})^{-1} \nabla_w f(w^{(0)})$

$$J(w) = \frac{1}{2n} \sum_{i=1}^n (x^T w - y^{(i)})^2$$

$$\nabla_w J(w) = \frac{1}{n} (x)(x^T w - y)$$

$$\nabla_w^2 J(w) = \frac{1}{n} (x x^T) \rightarrow \text{Hessian}$$

$$w^* = w^{(0)} - \left( \frac{1}{n} (x x^T) \right)^{-1} \left( \frac{1}{n} (x)(x^T w^{(0)} - y) \right)$$

$$= w^{(0)} - (x x^T)^{-1} (x x^T w^{(0)} - xy)$$

$$= w^{(0)} - w^{(0)} + (x x^T)^{-1} xy$$

$$= (x x^T)^{-1} xy$$

$\therefore$  The Newton's Method converges to the optimal solution  $w^* = (x x^T)^{-1} x y$  in 1 iteration

no matter what the starting point  $w^{(0)}$  of the search is.

2.

$$2) \hat{y}_k = \frac{\exp z_k}{\sum_{k=1}^c \exp z_k}, \quad z_k = x^T w^{(k)} + b_k, \quad f_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c y_i^{(k)} \log \hat{y}_k^{(k)}$$

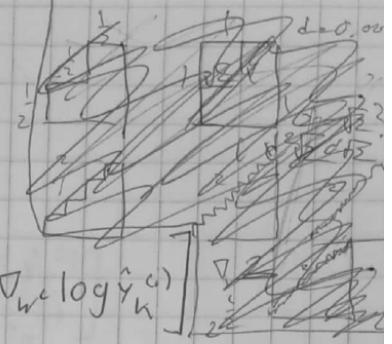
For  $l=k$ :

$$\begin{aligned}\nabla_{W^{(l)}} \hat{y}_k^{(i)} &= \nabla_{W^{(l)}} \left( \frac{\exp z_k}{\sum_{k'=1}^c \exp z_{k'}} \right)^{(i)} \\ &= \nabla_{W^{(l)}} \left( \frac{\exp(x^T w^{(l)} + b_l)}{\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'})} \right)^{(i)} \quad * d \left( \frac{f(x)}{g(x)} \right) = \frac{g(x) df(x) - f(x) dg(x)}{(g(x))^2} \\ &= \frac{\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'}) \cdot x^{(i)} \exp(x^T w^{(l)} + b_l) - \exp(x^T w^{(l)} + b_l) \cdot (x^{(i)} \exp(x^T w^{(k')} + b_{k'}))}{(\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'}))^2} \\ &= \frac{x^{(i)} \exp(x^T w^{(l)} + b_l) \left[ \sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'}) - \exp(x^T w^{(l)} + b_l) \right]}{(\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'}))^2} \\ &= \frac{x^{(i)} \exp(x^T w^{(l)} + b_l)}{\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'})} - \frac{x^{(i)} \exp(x^T w^{(l)} + b_l) \cdot \exp(x^T w^{(l)} + b_l)}{(\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'}))^2} \\ &= x^{(i)} \hat{y}_L^{(i)} - x^{(i)} \hat{y}_L^{(i)} \left( \frac{\exp(x^T w^{(l)} + b_l)}{\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'})} \right) \\ &= x^{(i)} \hat{y}_L^{(i)} - x^{(i)} \hat{y}_L^{(i)} \hat{y}_L^{(i)} \\ &= x^{(i)} \hat{y}_L^{(i)} (1 - \hat{y}_L^{(i)})\end{aligned}$$

For  $l \neq k$ 

$$\begin{aligned}\nabla_{W^{(l)}} \hat{y}_k^{(i)} &= \nabla_{W^{(l)}} \left( \frac{\exp z_k}{\sum_{k'=1}^c \exp z_{k'}} \right)^{(i)} \\ &= \nabla_{W^{(l)}} \left( \frac{\exp(x^T w^{(k)} + b_k)}{\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'})} \right)^{(i)} \\ &= 0 - \frac{\exp(x^T w^{(k)} + b_k) \cdot x^{(i)} \exp(x^T w_L + b_L)}{(\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'}))^2} \\ &= -x^{(i)} \hat{y}_k^{(i)} \hat{y}_L^{(i)}\end{aligned}$$

log



$$\nabla_{W^{(l)}} f_{CE}(W, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c Y_k^{(i)} \nabla_{W^{(l)}} \log \hat{Y}_k^{(i)}$$

$$\textcircled{1} \quad = -\frac{1}{n} \sum_{i=1}^n \left[ \cancel{\sum_{k \neq l} Y_k^{(i)} \nabla_{W^{(l)}} \log \hat{Y}_k^{(i)}} + \sum_{k \neq l} Y_k^{(i)} \nabla_{W^{(l)}} \log \hat{Y}_k^{(i)} \right]$$

$$\textcircled{2} \text{ Subst. } \nabla_{W^{(l)}} \log \hat{Y}_k^{(i)} = \frac{\nabla_{W^{(l)}} \hat{Y}_k^{(i)}}{\hat{Y}_k^{(i)}} : = -\frac{1}{n} \sum_{i=1}^n \left[ \frac{Y_k^{(i)} \nabla_{W^{(l)}} \hat{Y}_k^{(i)}}{\hat{Y}_k^{(i)}} + \sum_{k \neq l} \left[ \frac{Y_k^{(i)} \nabla_{W^{(l)}} \hat{Y}_k^{(i)}}{\hat{Y}_k^{(i)}} \right] \right]$$

$$\textcircled{3} \text{ Subst. } \nabla_{W^{(l)}} \hat{Y}_k^{(i)} = x^{(i)} \cdot \hat{Y}_k^{(i)} (1 - \hat{Y}_k^{(i)})$$

for  $k=1$

and

$$\text{SubstL } \nabla_{W^{(l)}} \hat{Y}_k^{(i)} = -x^{(i)} \hat{Y}_k^{(i)} \hat{Y}_l^{(i)} : = -\frac{1}{n} \sum_{i=1}^n \left[ \frac{Y_k^{(i)} (x^{(i)} \hat{Y}_l^{(i)} (-\hat{Y}_k^{(i)}))}{\hat{Y}_k^{(i)}} + \sum_{k \neq l} \left[ \frac{Y_k^{(i)} - x^{(i)} \hat{Y}_k^{(i)} \hat{Y}_l^{(i)}}{\hat{Y}_k^{(i)}} \right] \right]$$

Then cancel terms

$$= -\frac{1}{n} \sum_{i=1}^n \left[ Y_k^{(i)} x^{(i)} - \hat{Y}_k^{(i)} x^{(i)} \hat{Y}_l^{(i)} + \sum_{k \neq l} \left[ Y_k^{(i)} (-x^{(i)}) \hat{Y}_l^{(i)} \right] \right]$$

$$\textcircled{4} \text{ use } \sum_{k \neq l} a_k = a_1 + \sum_{k \neq l} a_k \text{ to combine}$$

$$Y_k^{(i)} x^{(i)} \hat{Y}_l^{(i)} \text{ with } \sum_{k \neq l} \left[ Y_k^{(i)} x^{(i)} \hat{Y}_l^{(i)} \right] : = -\frac{1}{n} \sum_{i=1}^n \left[ Y_k^{(i)} x^{(i)} - \sum_{k=1}^c \left[ x^{(i)} \hat{Y}_k^{(i)} \hat{Y}_l^{(i)} \right] \right]$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[ a_1 Y_k^{(i)} - x^{(i)} \hat{Y}_l^{(i)} \sum_{k=1}^c Y_k^{(i)} \right]$$

$$\textcircled{5} \text{ use } \sum_{k \in L} Y_k^{(i)} = 1 : \\ = -\frac{1}{n} \sum_{i=1}^n \left[ x^{(i)} \hat{Y}_l^{(i)} - x^{(i)} \hat{Y}_l^{(i)} \right]$$

$$\nabla_{W^{(l)}} f_{CE}(W, b) = -\frac{1}{n} \sum_{i=1}^n x^{(i)} \left( Y_k^{(i)} - \hat{Y}_l^{(i)} \right)$$

Proved

Show that  $\nabla_b F_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})$

$$1. \nabla_b F_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c Y_k^{(i)} \nabla_b \log \hat{Y}_k^{(i)}$$

$$2. \text{ rewrite } \nabla_b \log \hat{Y}_k^{(i)} \text{ as } \left( \frac{\nabla_b \hat{Y}_k^{(i)}}{\hat{Y}_k^{(i)}} \right) : -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^c Y_k^{(i)} \left[ \frac{\nabla_b \hat{Y}_k^{(i)}}{\hat{Y}_k^{(i)}} \right]$$

3. Substitute  $\hat{Y}_k^{(i)}$  for  $\frac{\exp(x^T w^{(i)} + b_i)}{\sum_{k=1}^c \exp(x^T w^{(k)} + b_k)}$  and simplify  $\nabla_b \hat{Y}_k^{(i)}$

$$\text{a) } \nabla_b \hat{Y}_k^{(i)} = \nabla_b \left[ \frac{\exp(x^T w^{(k)} + b_k)}{\sum_{k'=1}^c \exp(x^T w^{(k')} + b_{k'})} \right] \quad * \text{using } \nabla_b \hat{Y}_k^{(i)} = \hat{Y}_k^{(i)} (1 - \hat{Y}_k^{(i)})$$

$$\text{i. if } k=1: \nabla_b \hat{Y}_{k=1}^{(i)} = \frac{\exp(x^T w_1 + b_1)}{\sum_{k'=1}^c \exp(x^T w_{k'} + b_{k'})} \cdot \left[ \frac{1 - \exp(x^T w_1 + b_1)}{\sum_{k'=1}^c \exp(x^T w_{k'} + b_{k'})} \right] \quad \text{since } \frac{\partial \cdot b}{\partial \hat{Y}_k^{(i)}} \cdot \frac{\partial \hat{Y}_k^{(i)}}{\partial z_n^{(i)}} \cdot \frac{\partial z_n^{(i)}}{\partial x_n^{(i)}} = \exp(z_n^{(i)}) = \exp(b)$$

$$f(b) = x^T w + b$$

$$g(y) = \exp(b)$$

$$h(w) = \frac{1}{5}$$

$$\text{ii. } \nabla_b \hat{Y}_{k \neq 1}^{(i)} = \hat{Y}_1^{(i)} (1 - \hat{Y}_1^{(i)})$$

$$\text{iii. if } k \neq 1: \nabla_b \hat{Y}_k^{(i)} = \nabla_b \left[ \frac{\exp(x^T w^{(k)} + b^{(k)})}{\sum_{k'=1}^c \exp(x^T w^{(k')} + b^{(k')})} \right] \quad * \text{Solving with chain rule}$$

$$\text{iv. } = 0 \cdot \frac{\exp(x^T w + b) \cdot 1}{\exp(x^T w + b)} - \left[ \frac{\exp(x^T w + b) \cdot \exp(x^T w + b)}{\sum_{k'=1}^c (\exp(x^T w + b))} \right] \quad \text{since } \frac{\partial \cdot b}{\partial \hat{Y}_k^{(i)}} \cdot \frac{\partial \hat{Y}_k^{(i)}}{\partial z_n^{(i)}} \cdot \frac{\partial z_n^{(i)}}{\partial b} = g'(f(b)) \cdot f'(b) \cdot h(g(f(b)))$$

$$V_i = 0 - \frac{\exp(x^T w_k + b_k) \cdot \exp(x^T w_k + b_k)}{\sum_{k'=1}^c \exp(x^T w_{k'} + b_{k'})} \quad \text{since } g(f(b)) \cdot h'(g(f(b))) = g'(f(b)) \cdot f'(b)$$

$$\text{v. } \nabla_b \hat{Y}_{k \neq 1}^{(i)} = \hat{Y}_k^{(i)} \hat{Y}_1^{(i)}$$

4. Since  $\nabla_b \hat{Y}_k^{(i)} = \nabla_b \hat{Y}_k^{(i)}_{(k=1)} + \nabla_b \hat{Y}_k^{(i)}_{(k \neq 1)}$ , substitute into eq (2):

$$\nabla_b F_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n \left[ \sum_{k=1}^c \left[ \frac{\hat{Y}_k^{(i)}}{\hat{Y}_1^{(i)}} \hat{Y}_1^{(i)} (1 - \hat{Y}_1^{(i)}) \right] + \sum_{k \neq 1} \left[ \frac{\hat{Y}_k^{(i)}}{\hat{Y}_1^{(i)}} (-\hat{Y}_k^{(i)} \hat{Y}_1^{(i)}) \right] \right] = -\frac{1}{n} \sum_{i=1}^n \left[ \frac{\hat{Y}_1^{(i)}}{\hat{Y}_1^{(i)}} \hat{Y}_1^{(i)} (1 - \hat{Y}_1^{(i)}) \right] = 1 - \frac{\hat{Y}_1^{(i)}}{\hat{Y}_1^{(i)}} = 1$$

$$= -\frac{1}{n} \sum_{i=1}^n \left[ \hat{Y}_1^{(i)} - \hat{Y}_1^{(i)} \hat{Y}_1^{(i)} - \sum_{k \neq 1} \hat{Y}_k^{(i)} \hat{Y}_1^{(i)} \right] = -\frac{1}{n} \sum_{i=1}^n \left[ \hat{Y}_1^{(i)} - \hat{Y}_1^{(i)} \sum_{k=1}^c \hat{Y}_k^{(i)} \right]$$

$$\therefore \nabla_b F_{CE}(w, b) = -\frac{1}{n} \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})$$

Proved!

3.

$$-\log P(D|W, b) = -\log \prod_{i=1}^n P(y^{(i)}|x^{(i)}, W, b)$$

$$= -\sum_{i=1}^n \log P(y^{(i)}|x^{(i)}, W, b)$$

$$= -\sum_{i=1}^n \log \prod_{k=1}^C \hat{y}_k^{(i) y_k^{(i)}}$$

$$= -\sum_{i=1}^n \sum_{k=1}^C \log \hat{y}_k^{(i) y_k^{(i)}}$$

log power  
property

$$= -\sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \log \hat{y}_k^{(i)}$$

$$\boxed{-\log P(D|W, b) = -\sum_{i=1}^n \sum_{k=1}^C y_k^{(i)} \log \hat{y}_k^{(i)}}$$

4.

```
● Eris-MacBook-Pro:DS541 eri$ python3 homework3_ekim4_ahorn_amevans.py
Optimized hyperparameters:
    Mini-batch size: 200
    Learning rate: 0.01
    # of epochs: 15
    Alpha: 1e-05
    Weights: [[ 0.3028434  0.38584837  0.96888511 ...  0.85339476  0.91509533
    0.86142884]
[ 0.48925388  0.38130835  0.47666787 ...  0.55585233  0.67352697
  0.82948217]
[ 0.54769323  0.73651123  0.66853498 ...  0.19282263  0.82017287
  0.18343668]
...
[ 0.25532291  0.39493413  0.88493968 ...  0.70837728 -0.72493666
  0.78279963]
[ 0.52411611  0.368282     1.46181658 ...  0.15141629  0.47389966
  1.20168119]
[ 0.45373107  0.87966073  0.65457577 ...  0.51667297  0.72426747
  0.1427419 ]
[[ 0.3028434  0.38584837  0.96888511 ...  0.85339476  0.91509533
  0.86142884]
[ 0.48925388  0.38130835  0.47666787 ...  0.55585233  0.67352697
  0.82948217]
[ 0.54769323  0.73651123  0.66853498 ...  0.19282263  0.82017287
  0.18343668]
...
[ 0.25532291  0.39493413  0.88493968 ...  0.70837728 -0.72493666
  0.78279963]
[ 0.52411611  0.368282     1.46181658 ...  0.15141629  0.47389966
  1.20168119]
[ 0.45373107  0.87966073  0.65457577 ...  0.51667297  0.72426747
  0.1427419 ]
    Bias:
        [[0.04171745 0.62786595 0.06841105 0.08477843 0.4107784 0.17955993
  0.05757667 0.05119002 0.11661289 0.73491587]]
    min_cost on validation data: 8.718590722681503

-----TESTING-----
Cross Entropy: 8.554876947412696
% of correctly classified images: 81.51%
```