

Juhi Patel - JPP2464

(1.1) $y = f(x) = w \cdot x + b = 2x_1 + x_2 + b$

$$\frac{\partial y}{\partial x_1} = w_1 = 2$$

$$\frac{\partial y}{\partial x_2} = w_2 = 1$$

(1.2) $MSE = \frac{1}{4}((0-0.1)^2 + (1-0.95)^2 + (1-1.10)^2 + (0-0.2)^2)$
 $= 0.015625$

(1.3) i) $y = w \cdot x + b = 2x_1 + x_2 + 3$
 $\hat{y}_1 = 0 + 0 + 3 = 3$
 $\hat{y}_2 = 0 + 1 + 3 = 4$
 $\hat{y}_3 = 2 + 0 + 3 = 5$
 $\hat{y}_4 = 2 + 1 + 3 = 6$

$MSE = \frac{1}{4}((3-0)^2 + (4-1)^2 + (5-1)^2 + (6-0)^2)$
 $= 17.5$

ii) $\hat{y}_i = w \cdot x_i + b = w_1 x_{1i} + w_2 x_{2i} + b$

$$MSE = \frac{1}{n} \sum_{i=1}^n (w_1 x_{1i} + w_2 x_{2i} + b - y_i)^2$$

$$\frac{d(MSE)}{d(w_1)} = \frac{2}{n} \sum_{i=1}^n (w_1 x_{1i} + w_2 x_{2i} + b - y_i) x_{1i}$$

$$= \frac{2}{4} [(3-0) \times 0 + (4-1) \times 0 + (5-1) \times 1 + (6-0) \times 1]$$
$$= 5$$

$$\frac{d(MSE)}{d(w_2)} = \frac{2}{n} \sum_{i=1}^n (w_1 x_{1i} + w_2 x_{2i} + b - y_i) x_{2i}$$

$$= \frac{2}{4} [(3-0) \times 0 + (4-1) \times 1 + (5-1) \times 0 + (6-0) \times 1]$$
$$= 4.5$$

direction $(-5, -4.5, -8)$

$$\text{iii) } \frac{d(\text{MSE})}{dw_1} = \frac{d(\text{MSE})}{dw_2} = \frac{d(\text{MSE})}{db} = 0$$

$$f(x) = w \cdot x + b = w_1 x_1 + w_2 x_2 + b = w_1 x_1 + w_2 x_2 + b$$

$$\hat{y}_1 = b \rightarrow \hat{y}_1 - y_1 = b$$

$$\hat{y}_2 = w_2 + b \rightarrow \hat{y}_2 - y_2 = w_2 + b - 1$$

$$\hat{y}_3 = w_1 + b \rightarrow \hat{y}_3 - y_3 = w_1 + b - 1$$

$$\hat{y}_4 = w_1 + w_2 + b \rightarrow \hat{y}_4 - y_4 = w_1 + w_2 + b$$

$$\frac{d(\text{MSE})}{dw_1} = 0 \rightarrow \sum_{i=1}^n (\hat{y}_i - y_i) x_{1i} = 0$$

$$w_1 + b - 1 \quad w_1 + w_2 + b = 0 \rightarrow 2w_1 + w_2 + 2b - 1 = 0$$

$$\frac{d(\text{MSE})}{dw_2} = 0 \rightarrow \sum_{i=1}^n (\hat{y}_i - y_i) x_{2i} = 0$$

$$\rightarrow w_2 + b + w_1 + w_2 + b = 0 \rightarrow w_1 + 2w_2 + 2b - 1 = 0$$

$$\frac{d(\text{MSE})}{db} = 0 \rightarrow \sum_{i=1}^n (\hat{y}_i - y_i) = 0$$

$$\rightarrow 2w_1 + 2w_2 + 4b - 2 = 0 \rightarrow w_1 + w_2 + 2b - 1 = 0$$

$$w_1 = 0, w_2 = 0, w_3 = \frac{1}{2}$$

Extra Credit on python

(2.1)

$$x = [1, -1, 3, 4, 4]$$

$$\text{kernel} = [1, 1]$$

$$\text{for stride} = 1$$

$$\begin{aligned} \text{conv}(k, x) &= [1-1, -1+3, 3+4, 4+4] \\ &= [0, 2, 7, 8] \end{aligned}$$

$$\text{for stride} = 2,$$

$$\text{conv}(k, x) = [1-1, 3+4] = [0, 7]$$

(2.2)

(i)

$$I = \begin{bmatrix} 0.1 & -0.6 & 0.4 & 0.8 \\ -0.4 & 0.3 & 0.9 & 0.2 \\ 0.5 & 0.2 & 0.8 & -0.7 \\ 0.3 & 0.7 & -0.4 & 0.1 \end{bmatrix}$$

$$K_1 = \begin{bmatrix} 1 & 6 \\ 0 & 1 \end{bmatrix}, K_2 = \begin{bmatrix} 6 & 1 \\ 1 & 0 \end{bmatrix}$$

$$\text{Conv}(K_1, I), \text{stride} = [1, 1]$$

$$= \begin{bmatrix} 0.1+0.3 & -0.6+0.9 & 0.4+0.2 \\ -0.4+0.2 & 0.3+0.8 & 0.9-0.7 \\ 0.5+0.7 & 0.2-0.4 & 0.8+0.1 \end{bmatrix} = \begin{bmatrix} 0.4 & 0.3 & 0.6 \\ -0.2 & 1.1 & 0.2 \\ 1.2 & -0.2 & 0.9 \end{bmatrix}$$

$$\text{Conv}(K_2, I), \text{stride} = [1, 1]$$

$$\begin{bmatrix} -0.6-0.4 & 0.4+0.3 & 0.8+0.9 \\ 0.3+0.5 & 0.9+0.2 & 0.2+0.8 \\ 0.2+0.3 & 0.8+0.7 & -0.7-0.4 \end{bmatrix} = \begin{bmatrix} -1.0 & 0.7 & 1.7 \\ 0.8 & 1.1 & 1.0 \\ 0.5 & 1.5 & -1.1 \end{bmatrix}$$

- ii) each channel has dimensions 3×3 in $2.2i$
 if the stride is $[2, 2]$ channel would be 2×2
- iii) filter should 4×1 , stride 1 or 2 or any other
- iv) maxpool (I) with filters 2×2 stride = 1
- $$\begin{bmatrix} 0.3 & 0.9 & 0.9 \\ 0.5 & 0.9 & 0.9 \\ 0.7 & 0.8 & 0.8 \end{bmatrix} = 3 \times 3 \text{ dimensions}$$

Extra Credit

advantages of max pooling

- (i) easier to do computationally
- (ii) easier extraction when it comes to extracting most prominent features

When do you think max-pooling operation may not be advantageous?

* You ~~are~~ would not want to extract prominent features only while still keeping subtle features to remain in network - the max-pool would discard them

3.1

- i) False
- ii) True
- iii) False
- iv) False

3.2

$$\text{pr}(x = \text{cat}) = \text{sigmoid}(w'x)$$

$$= \frac{1}{1 + \exp(-w_1 x_1 - w_2 x_2)}$$

$$\text{given } x_i = [2, 1], \text{ \& } w = [1, 1]$$

$$i) \text{pr}(x = \text{dog}) =$$

$$1 - \frac{1}{1 + \exp(-3)} = 1 - 0.952$$

$$\text{Pr}(x = \text{cat}) = 0.952 = \underline{95.2\%}$$

ii)

$$\|x - x_1\|_{\text{inf}} \leq 0.1$$

$$x = [a, b] \quad \text{then } x - x_1 = [a-2, b-1]$$

$$\text{Case 1: } a > 2 \quad \& \quad b > 1$$

$$\|x - x_1\|_{\text{inf}} = \begin{cases} a-2; & a-2 > b-1 \\ b-1; & b-1 > a-2 \end{cases} = \begin{cases} a-2; & a-b > 1 \\ b-1; & a-b < 1 \end{cases}$$

$$\text{Case 2 } a < 2 \quad \text{and } b > 1:$$

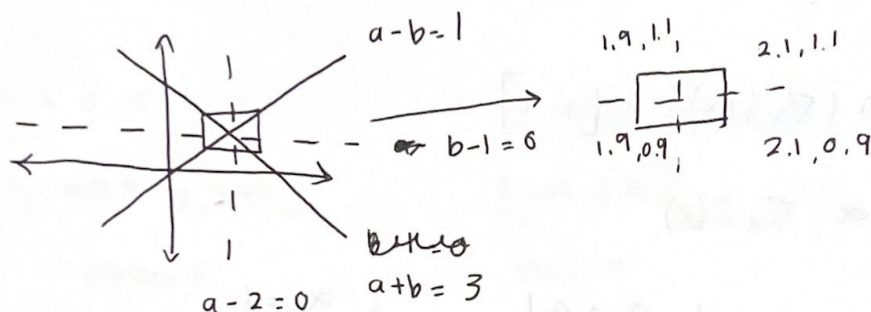
$$\|x - x_1\|_{\text{inf}} = \begin{cases} a-2; & 2-a > b-1 \\ b-1; & 2-a < b-1 \end{cases} = \begin{cases} a-2; & a+b < 3 \\ b-1; & a+b > 3 \end{cases}$$

$$\text{Case 3 } a > 2 \quad \text{and } b < 1$$

$$\|x - x_1\|_{\text{inf}} = \begin{cases} a-2; & a-2 > 1-b \\ b-1; & a-2 < 1-b \end{cases} = \begin{cases} a-2; & a+b > 3 \\ b-1; & a+b < 3 \end{cases}$$

$$\text{Case 4 } a < 2 \quad \& \quad b < 1$$

$$\|x - x_1\|_{\text{inf}} = \begin{cases} a-2; & 2-a > 1-b \\ b-1; & 2-a < 1-b \end{cases} = \begin{cases} a-2; & a-b < 1 \\ b-1; & a-b > 1 \end{cases}$$



iii)

$$x_{\text{optimiz}} = \max_x \left(1 - \frac{1}{1 + \exp(-w'x)} \right)$$

$$\text{s.t. } \|x - x_1\|_{\text{inf}} < 0.1$$

$$\text{given } x_1 = [2, 1]$$

iv)

$$f(x) = \text{sigmoid}(w'x) = \frac{1}{1 + \exp(-w'x)}$$

$$f(x) = \frac{d}{dx} \left[\frac{1}{1 + \exp(-w'x)} \right]$$

$$= \frac{-1}{(1 + \exp(-w'x))^2} \times \exp(-w'x) \times -w'$$

$$= \frac{w' \exp(-w'x)}{(1 + \exp(-w'x))^2}$$

for $w = [1, 1]$ and $x_1 = [2, 1]$

$$\nabla_x f(x) = \frac{\exp(-5)}{(1 + \exp(-3))^2} [1, 1] = [0.045, 0.045]$$

$$\text{sign}(\nabla_x f(x)) = [1, 1]$$

$$x_{\text{attack}} = x_1 - \alpha \nabla_x f(x)$$

$$\alpha = 0.001$$

$$[1.99995, 0.99995]$$

threat model

$$\alpha = 0.1$$

$$[1.995, 0.995]$$

threat model

$$\alpha = 1$$

$$[1.995, 0.995]$$

threat model

$$x_{\text{attack}} = x_1 - \alpha \text{sign}(\nabla_x f(x))$$

$$\alpha = 0.001$$

$$[1.999, 0.999]$$

threat model

$$\alpha = 0.1$$

$$[1.9, 0.9]$$

threat model

$$\alpha = 1$$

$$[1, 0]$$

not in threat model

ii)

$$\|x - x_1\|_{\text{int}} \leq 0.1$$

$$x = [a, b] \quad \text{then } x - x_1 = [a-2, b-1]$$

$$\text{Case 1: } a > 2 \quad \& \quad b > 1$$

$$\|x - x_1\|_{\text{int}} = \begin{cases} a-2; a-2 > b-1 \\ b-1; b-1 > a-2 \end{cases} = \begin{cases} a-2; a-b > 1 \\ b-1; a-b < 1 \end{cases}$$

$$\text{Case 2} \quad a < 2 \quad \text{and} \quad b > 1$$

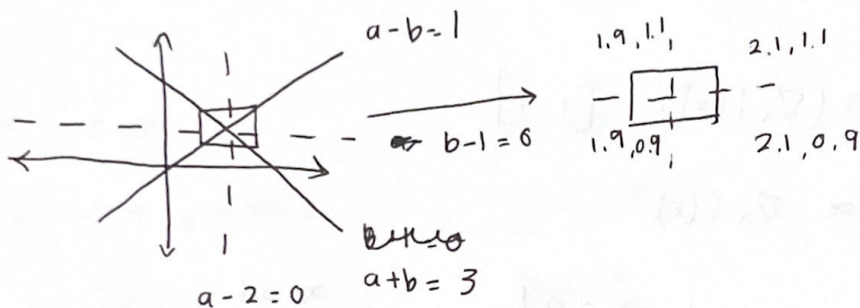
$$\|x - x_1\|_{\text{int}} = \begin{cases} a-2; 2-a > b-1 \\ b-1; 2-a < b-1 \end{cases} = \begin{cases} a-2; a+b < 3 \\ b-1; a+b > 3 \end{cases}$$

$$\text{Case 3} \quad a > 2 \quad \text{and} \quad b < 1$$

$$\|x - x_1\|_{\text{int}} = \begin{cases} a-2; a-2 > 1-b \\ b-1; a-2 < 1-b \end{cases} = \begin{cases} a-2; a+b > 3 \\ b-1; a+b < 3 \end{cases}$$

$$\text{Case 4} \quad a < 2 \quad \& \quad b < 1$$

$$\|x - x_1\|_{\text{int}} = \begin{cases} a-2; 2-a > 1-b \\ b-1; 2-a < 1-b \end{cases} = \begin{cases} a-2; a-b < 1 \\ b-1; a-b > 1 \end{cases}$$



iii)

$$x_{\text{optimiz}} = \max_x \left(\frac{1}{1 + \exp(-w'x)} \right)$$

$$\text{s.t. } \|x - x_1\|_{\text{int}} < 0.1$$

$$\text{given } x_1 = [2, 1]$$

$$(v) \quad f(x) = 1 - \text{sigmoid}(w'x) = 1 - \frac{1}{(1 + \exp(-w'x))^2}$$

$$f'(x) = \frac{df(x)}{dx} = \frac{-w' \exp(-w'x)}{(1 + \exp(-w'x))^2}$$

$$\nabla_x f(x) = [-0.045, -0.045], \quad \text{sign}(\nabla_x f(x)) = [-1, -1]$$

$$x_{\text{attack}} = x_1 + \alpha \nabla_x f(x)$$

$$x_{\text{attack}} = x_1 + \alpha \nabla_x f(x)$$

$$\alpha = 0.001$$

$$[1.99995, 0.99995]$$

threat

$$\alpha = 0.1$$

$$[1.995, 0.995]$$

threat

$$\alpha = 1$$

$$[1.955, 0.955]$$

threat

$$\alpha = 0.001$$

$$[1.999, 0.999]$$

threat

$$\alpha = 0.1$$

$$[1.9, 0.9]$$

threat

$$\alpha = 1$$

$$[1, 0]$$

not threat

vi) we see that new points lie in the threat model when we use gradient values. This means the new images aren't very different. But when we use the signs of the gradient it can have a stronger attack. This can allow the new image to be close to the target image. In this case, targeted and untargeted attack make no difference to the result since there are only two cases.