

$$1. \quad W = \begin{bmatrix} W_{00} & W_{01} & W_{02} \\ W_{10} & W_{11} & W_{12} \\ W_{20} & W_{21} & W_{22} \end{bmatrix} \quad X = \begin{bmatrix} X_{00} & X_{01} & X_{02} \\ X_{10} & X_{11} & X_{12} \\ X_{20} & X_{21} & X_{22} \end{bmatrix}$$

$$\text{Stride} = 4 \quad \text{padding} = 2$$

$$\text{output size} = \frac{3 - 3 + 2(2)}{4} + 1 = 2 \Rightarrow 2 \times 2 \text{ matrix}$$

$$X_{\text{padded}} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & X_{00} & X_{01} & X_{02} & 0 & 0 \\ 0 & 0 & X_{10} & X_{11} & X_{12} & 0 & 0 \\ 0 & 0 & X_{20} & X_{21} & X_{22} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$\nwarrow w_1 \quad \nearrow w_2$   
 $\downarrow w_3 \quad \swarrow w_4$

$$w_1 = W_{22} X_{00} \quad w_2 = W_{20} X_{02}$$

$$w_3 = W_{02} X_{20} \quad w_4 = W_{00} X_{22}$$

$$WX = \begin{bmatrix} W_{22} X_{00} & W_{20} X_{02} \\ W_{02} X_{20} & W_{00} X_{22} \end{bmatrix}$$

$$Y = [W_{22} X_{00} \quad W_{20} X_{02} \quad W_{02} X_{20} \quad W_{00} X_{22}]^T$$

$$Y = \underbrace{\begin{bmatrix} W_{22} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & W_{20} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & W_{02} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & W_{00} \end{bmatrix}}_A \begin{bmatrix} X_{00} \\ X_{01} \\ X_{02} \\ X_{10} \\ X_{11} \\ X_{12} \\ X_{20} \\ X_{21} \\ X_{22} \end{bmatrix} = Ax$$

$$2. \quad W = \begin{bmatrix} W_{00} & W_{01} \\ W_{10} & W_{11} \end{bmatrix} \quad X = \begin{bmatrix} X_{00} & X_{01} \\ X_{10} & X_{11} \end{bmatrix}$$

Stride = 2

define  $W$  as

$$W = \begin{bmatrix} W_{00} & W_{01} & \xrightarrow{X_{00}} & \xrightarrow{X_{01}} \\ W_{10} & W_{11} & W_{00} & W_{01} \\ W_{00} & W_{01} & W_{00} & W_{01} \\ W_{10} & W_{11} & W_{10} & W_{11} \end{bmatrix} = \begin{bmatrix} X_{00}W_{00} & X_{00}W_{01} & X_{01}W_{00} & X_{01}W_{01} \\ X_{00}W_{10} & X_{00}W_{11} & X_{01}W_{10} & X_{01}W_{11} \\ X_{10}W_{00} & X_{10}W_{01} & X_{11}W_{00} & X_{11}W_{01} \\ X_{10}W_{10} & X_{10}W_{11} & X_{11}W_{10} & X_{11}W_{11} \end{bmatrix}$$

row-wise  $W = [X_{00}W_{00} \ X_{00}W_{01} \ X_{01}W_{00} \ X_{01}W_{01} \ X_{00}W_{10} \ X_{00}W_{11} \ X_{01}W_{10} \ X_{01}W_{11}]$   
 $\dots \ X_{11}W_{11}]$

$$Y = \begin{bmatrix} W_{00} & 0 & 0 & 0 \\ W_{01} & 0 & 0 & 0 \\ 0 & W_{00} & 0 & 0 \\ 0 & W_{01} & 0 & 0 \\ W_{10} & 0 & 0 & 0 \\ W_{11} & 0 & 0 & 0 \\ 0 & W_{10} & 0 & 0 \\ 0 & W_{11} & 0 & 0 \\ 0 & 0 & W_{00} & 0 \\ 0 & 0 & W_{01} & 0 \\ 0 & 0 & 0 & W_{00} \\ 0 & 0 & 0 & W_{01} \\ 0 & 0 & W_{10} & 0 \\ 0 & 0 & W_{11} & 0 \\ 0 & 0 & 0 & W_{10} \\ 0 & 0 & 0 & W_{11} \end{bmatrix} \begin{bmatrix} X_{00} \\ X_{01} \\ X_{10} \\ X_{11} \end{bmatrix} = Ax$$

$\underbrace{\phantom{W_{00} \ 0 \ 0 \ 0 \ W_{01} \ 0 \ 0 \ 0 \ 0 \ W_{10} \ 0 \ 0 \ 0 \ 0 \ W_{11} \ 0 \ 0 \ 0 \ 0 \ W_{10} \ 0 \ 0 \ 0 \ 0 \ W_{11} \ 0 \ 0 \ 0 \ 0 \ W_{00} \ 0 \ 0 \ 0 \ 0 \ W_{01} \ 0 \ 0 \ 0 \ 0 \ W_{10} \ 0 \ 0 \ 0 \ 0 \ W_{11} \ 0 \ 0 \ 0 \ 0 \ W_{10} \ 0 \ 0 \ 0 \ 0 \ W_{11} \ 0 \ 0 \ 0 \ 0 \ W_{00} \ 0 \ 0 \ 0 \ 0 \ W_{01} \ 0 \ 0 \ 0 \ 0 \ W_{10} \ 0 \ 0 \ 0 \ 0 \ W_{11} \ 0 \ 0 \ 0 \ 0 \ W_{10} \ 0 \ 0 \ 0 \ 0 \ W_{11}}_A$

3. Kernel  $(o \times r^2, i, k, k) \equiv \text{Kernel}(o, i, k_{xr}, k_{xr})$

$$o=1 \quad r=2 \quad k=1$$

$$\Rightarrow (4, 1, 1, 1) \equiv (1, 1, 2, 2)$$

For  $w$  of size  $(4, 1, 1, 1)$

$$w \Rightarrow w_1 = [w_1]$$

$$w_2 = [w_2]$$

$$w_3 = [w_3]$$

$$w_4 = [w_4]$$

$$x = \begin{bmatrix} x_{00} & x_{01} \\ x_{10} & x_{11} \end{bmatrix}$$

For  $w$  of size  $(1, 1, 2, 2)$

$$w = \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix} \quad x = \begin{bmatrix} x_{00} & x_{01} \\ x_{10} & x_{11} \end{bmatrix}$$

From  $q^2$  we can write:

$$Y = \begin{bmatrix} w_1 & 0 & 0 & 0 \\ w_2 & 0 & 0 & 0 \\ 0 & w_1 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ w_3 & 0 & 0 & 0 \\ 0 & w_3 & 0 & 0 \\ 0 & 0 & w_1 & 0 \\ 0 & 0 & w_2 & 0 \\ w_4 & 0 & 0 & 0 \\ 0 & w_4 & 0 & 0 \\ 0 & 0 & w_3 & 0 \\ 0 & 0 & 0 & w_4 \end{bmatrix} \begin{bmatrix} x_{00} \\ x_{01} \\ x_{10} \\ x_{11} \end{bmatrix} = \begin{bmatrix} w_1 x_{00} \\ w_2 x_{00} \\ 0 \\ w_1 x_{01} \\ 0 \\ w_2 x_{01} \\ w_3 x_{00} \\ w_3 x_{01} \\ w_4 x_{00} \\ w_4 x_{01} \\ w_3 x_{10} \\ w_4 x_{10} \\ w_4 x_{11} \\ w_4 x_{11} \end{bmatrix}$$

$$= \begin{bmatrix} x_{00} \\ x_{01} \\ x_{10} \\ x_{11} \\ w_1 x_{00} \\ w_2 x_{00} \\ w_1 x_{01} \\ w_2 x_{01} \\ w_3 x_{00} \\ w_3 x_{01} \\ w_4 x_{00} \\ w_4 x_{01} \\ w_4 x_{10} \\ w_4 x_{11} \end{bmatrix}$$

$$Y = \begin{bmatrix} w_1 & 0 & 0 & 0 \\ 0 & w_1 & 0 & 0 \\ 0 & 0 & w_1 & 0 \\ 0 & 0 & 0 & w_1 \\ w_2 & 0 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ 0 & 0 & w_2 & 0 \\ w_2 & 0 & 0 & 0 \\ 0 & w_3 & 0 & 0 \\ 0 & 0 & w_3 & 0 \\ w_3 & 0 & 0 & 0 \\ 0 & w_3 & 0 & 0 \\ 0 & 0 & w_3 & 0 \\ w_4 & 0 & 0 & 0 \\ 0 & w_4 & 0 & 0 \\ 0 & 0 & w_4 & 0 \\ 0 & 0 & 0 & w_4 \end{bmatrix} \begin{bmatrix} x_{00} \\ x_{01} \\ x_{10} \\ x_{11} \end{bmatrix} = \begin{bmatrix} w_1 x_{00} \\ w_1 x_{01} \\ w_1 x_{10} \\ w_1 x_{11} \\ w_2 x_{00} \\ w_2 x_{01} \\ w_2 x_{10} \\ w_2 x_{11} \\ w_3 x_{00} \\ w_3 x_{01} \\ w_3 x_{10} \\ w_3 x_{11} \\ w_4 x_{00} \\ w_4 x_{01} \\ w_4 x_{10} \\ w_4 x_{11} \end{bmatrix}$$

We can see that both  $Y$ 's have same elements, just in different order

$$4. f(x) = \begin{cases} 1 & \text{if } w^T x + b \geq 0 \\ 0 & \text{if } w^T x + b < 0 \end{cases}$$

| $x_1$ | $x_2$ | $f_{\text{AND}}(x)$ |
|-------|-------|---------------------|
| 0     | 0     | 0                   |
| 0     | 1     | 0                   |
| 1     | 0     | 0                   |
| 1     | 1     | 1                   |

$$W_{\text{AND}} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad b_{\text{AND}} = -2$$

$$f\left(W_{\text{AND}}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b_{\text{AND}}\right) = f_{\text{AND}}(x)$$

$$\begin{array}{l} x_1=0 \\ x_2=1 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 2 = -1 \Rightarrow 0$$

$$\begin{array}{l} x_1=1 \\ x_2=0 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 2 = -1 \Rightarrow 0$$

$$\begin{array}{l} x_1=1 \\ x_2=1 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 2 = 0 \Rightarrow 1$$

$$\begin{array}{l} x_1=0 \\ x_2=0 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 2 = -2 \Rightarrow 0$$

$$W_{\text{OR}} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad b_{\text{OR}} = -0.1$$

$$\begin{array}{l} x_1=0 \\ x_2=1 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 0.1 = 0.9 \Rightarrow 1$$

$$\begin{array}{l} x_1=1 \\ x_2=0 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 0.1 = 0.9 \Rightarrow 1$$

$$\begin{array}{l} x_1=1 \\ x_2=1 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 0.1 = 1.9 \Rightarrow 1$$

$$\begin{array}{l} x_1=0 \\ x_2=0 \end{array} \Rightarrow \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 = -0.1 \Rightarrow 0$$

| 5. | $x_1$ | $x_2$ | $f_{XOR}(x)$ |
|----|-------|-------|--------------|
|    | 0     | 0     | 0            |
|    | 0     | 1     | 1            |
|    | 1     | 0     | 1            |
|    | 1     | 1     | 0            |

proof by contradiction:

If  $f(x) = w^T x + b$  can output a linear model function

$$\downarrow \text{Let } x_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, x_2 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\text{Then } f(x_1) < f(x_2)$$

$$f\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) < f\left(\begin{bmatrix} 0 \\ 1 \end{bmatrix}\right)$$

$$w\begin{bmatrix} 1 \\ 0 \end{bmatrix} + b < w\begin{bmatrix} 0 \\ 1 \end{bmatrix} + b$$

$$\left( \text{Let } w = [w_1 \ w_2] \right)$$

$$[w_1 \ w_2] \begin{bmatrix} 1 \\ 0 \end{bmatrix} + b < [w_1 \ w_2] \begin{bmatrix} 0 \\ 1 \end{bmatrix} + b$$

$$w_1 + w_2 + b < w_2 + b$$

$$w_1 < 0$$

Now let's look at for  $x_1 = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, x_2 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ :

$$f(x_1) < f(x_2)$$

$$w\begin{bmatrix} 0 \\ 1 \end{bmatrix} + b < w\begin{bmatrix} 1 \\ 0 \end{bmatrix} + b$$

$$b < w_1 + b$$

$$0 < w_1$$

$\Rightarrow$  we have a contradiction, so XOR can not be represented using Linear Model  $\blacksquare$

$$h(x) = w^{(3)} \max \{0, w^{(2)} \{0, w^{(1)}x + b^{(1)}\} + b^{(2)}\} + b^{(3)}$$

$$w^{(1)} = \begin{bmatrix} 1.5 \\ 0.5 \end{bmatrix} \quad w^{(2)} = \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \quad w^{(3)} = \begin{bmatrix} 1 & 1 \end{bmatrix}$$

$$b^{(1)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad b^{(2)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad b^{(3)} = -1$$

$$\underline{x_0 = 2}$$

$$\begin{aligned} & [1 \ 1] \max \{0, \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \max \{0, \begin{bmatrix} 1.5 \\ 0.5 \end{bmatrix} [2] + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} - 1 \\ &= [1 \ 1] \max \{0, \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} - 1 \\ &= [1 \ 1] \max \{0, \begin{bmatrix} 7 \\ 8 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} - 1 \\ &= [1 \ 1] \begin{bmatrix} 7 \\ 9 \end{bmatrix} - 1 = 7+9-1 = 15 \end{aligned}$$

$$\left. \begin{array}{l} x=2 \\ w=7 \\ b=1 \end{array} \right\} \Rightarrow y = 7x+1 = 7(2)+1 = 15$$

$$\frac{dh}{dx} = w = 7$$

$$\underline{x_0 = -1}$$

$$\begin{aligned} & [1 \ 1] \max \{0, \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \max \{0, \begin{bmatrix} 1.5 \\ 0.5 \end{bmatrix} [-1] + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} - 1 \\ &= [1 \ 1] \max \{0, \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \max \{0, \begin{bmatrix} -1.5 \\ 0.5 \end{bmatrix}\} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} - 1 \\ &= [1 \ 1] \max \{0, \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} - 1 \\ &= [1 \ 1] \max \{0, \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix}\} - 1 \\ &= [1 \ 1] \begin{bmatrix} 1 \\ 1.5 \end{bmatrix} - 1 = 1+1.5-1 = 1.5 \end{aligned}$$

$$\left. \begin{array}{l} x=-1 \\ w=-1 \\ b=0.5 \end{array} \right\} y = -1x + 0.5 = (-1)(-1) + 0.5 = 1.5$$

$$\frac{dh}{dx} = w = -1$$

$$X_0 = 1$$

$$\begin{aligned} & [1, 1] \max \{ 0, [1, 2] \max \{ 0, [1.5] [1, 1] + [0] \} + [0] \} - 1 \\ &= [1, 1] \max \{ 0, [1, 2] [1.5] + [0] \} - 1 \\ &= [1, 1] \max \{ 0, [4.5] + [0] \} - 1 \\ &= [1, 1] [4.5] - 1 = 4.5 + 5.5 - 1 = 9 \end{aligned}$$

$$\left. \begin{array}{l} x=1 \\ w=8 \\ b=1 \end{array} \right\} y = 8x + 1 = 8(1) + 1 = 9$$

$$\frac{dh}{dx} = w = 8$$

7. Given :  $f(x) = \mathbf{W}x + b$  divides input space into 2 regions for value  $y$

For each linear region, the region of the input is a bijection

$$w_{ij}^{(1)} = \begin{cases} 2 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases} \leftarrow W \text{ is a diagonal matrix}$$

$$b_i^{(1)} = -1$$

$$f_i(x) = |W^{(1)}x + b_i|$$

Proof : Take hypercube and partition it into halves with hyperplanes in each dimension  
then for each region  $R_i$ , there is a bijection

$$x_i = \begin{cases} \frac{-y_i+1}{2} & \text{if } y_i=1, x_i = \begin{cases} -1 \\ +1 \end{cases} \\ \frac{+y_i+1}{2} & \end{cases}$$

$$\begin{aligned} f_i(x_i) &= |W^{(1)}x + b_i| = \sum_j w_{ij} x_i + b_i \\ &= |2x_i + b_i| \\ &= |2\left(\frac{\pm y_i+1}{2}\right) - 1| \\ &= |\pm y_i + 1 - 1| = y_i \end{aligned}$$

Assume  $x_1 \neq x_2$

$$\Rightarrow 2x_1 \neq 2x_2$$

$$\Rightarrow |f(x_1)| \neq |f(x_2)|$$

$\Rightarrow$  we have a bijection and  $|R| = |(0,1)^d| = 2^d$  possible choices of input regions  $R_i$

Alternatively, we can see that for  $d=1$

$$O = (0,1)^1$$

input space  $R = \{(-1,0), (0,1)\} \Rightarrow 2$  input regions

then it is intuitive to see that

$$O = (0,1)^d \Rightarrow 2^d \text{ input regions}$$

8.  $g$  has  $n_g$  regions of  $(0,1)^d$   
 $f$  has  $n_f$  regions of  $(0,1)^d$

$$f \circ g(\cdot) = f(g(\cdot))$$

Let  $G$  be input region of  $g(\cdot)$

Let  $F$  be input region of  $f(\cdot)$

① Each subregion in  $G$  bijects onto  $f$   
because  $g(G)$  bijects onto  $(0,1)^d$

②  $f$  bijects between  $F$  and  $O$

① + ②  $\Rightarrow f \circ g$  bijects between  $G$  and  $O$

$$g(G) \Rightarrow |n_g|$$

$$f(g(G)) \Rightarrow |n_f n_g| \text{ regions}$$

$$9. h_1 = |w_1 \cdot x + b_1|$$

$$h_2 = |w_2 \cdot h_1 + b_2|$$

:

$$h_L = |w_L \cdot h_{L-1} + b_L|$$

$$x \in (0,1)^d \quad f(x) = h_L$$

Each  $h_i$  is implicitly a function of  $x$ :

$$h_1(x) = |w_1 \cdot x + b_1|$$

$$h_2(x) = |w_2 \cdot (w_1 \cdot x + b_1) + b_2|$$

:

$$h_L \cdot h_{L-1} \cdots h_1(x) = |w_L \cdot h_{L-1} \cdots h_1(x) + b_L|$$

Proof by induction:

Base case:  $L=1$  we know from q4 that  $h_1(x)$  identifies

$2^d$  input region space

$$2^{Ld} = 2^1 d = 2^d$$

Inductive Step:

Given true for  $L = l-1 \Rightarrow 2^{(l-1)d}$  input region space, prove

true for  $L+1 = l$

Let  $g(x) = h_{l-1} \cdot h_{l-2} \cdots h_1(x) + b_{l-1}$  identify  $2^{(l-1)d}$  input region space

Let  $f(x) = h_l(g(x))$

By q8 we know  $f \circ g(\cdot) = |n_f n_g|$  input region space

$h_l$  identifies  $2^d$  input region space

$\Rightarrow f(x) = h_l(g(x)) = h_l \circ g \Rightarrow |n_{h_l} n_g| = |2^d \cdot 2^{(l-1)d}| = 2^{ld}$   
input region space



$$10. \quad \min_{w \in \mathbb{R}^d} f(w) = \sum_{i=1}^n (y_i - w^\top x_i)^2 \quad \textcircled{1}$$

$$w^{(t+1)} = w^{(t)} - \eta \nabla f(w^{(t)})$$

$$\overset{\vdots}{w^{g_d}} = w^{(g_d-1)} - \eta \nabla f(w^{(g_d-1)})$$

11

$$w^{g_d} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \|w\|_2^2 \quad \text{s.t. } Xw = y \quad \textcircled{2}$$

$$\|w\|_2^2 = \sum_i^d w_i^2$$

For \textcircled{2}:

Introduce Lagrange Multipliers

$$L(x, \lambda) = \|w\|_2^2 + x^\top (w^\top x - y)$$

$$\text{s.t. } \nabla_x L = 2x + w^\top \lambda = 0$$

$$\nabla_y L = w^\top x - y = 0$$

$$\Rightarrow \lambda = -2(w^\top)^{-1}y$$

$$\Rightarrow x = w^\top (ww^\top)^{-1}y \quad \textcircled{3}$$

$$\text{For } \textcircled{1}: \quad w^{g_d} = (w^\top w)^2 + 2(y - w^\top x)$$

$$x = (w^\top w + I)^{-1} w^\top y$$

$$= w^\top (ww^\top)^{-1}y \quad \textcircled{4}$$

\textcircled{3} = \textcircled{4} \Rightarrow \text{solutions found by both are same}

11. (25) means sol'n will have smallest norm of any sol'n  
(24) is the optimization problem  
→ gradient descent will find the optimal sol'n for the loss function

12. Least squares regression can always find global min of linear system if it has a soln.

### **Key contributions**

The paper presents findings on double descent, where when model size is increased, performance first decreases then increases. They present empirical evidence that challenges and reconciles some of the conventional wisdoms of statistical theories and machine learning. They show that when model complexity is small compared to the number of samples, test error as a function of model complexity behaves like the classical bias/variance tradeoff. However, when model complexity large enough to interpolate, increasing complexity only decreases test error.

### **Strengths**

The paper thoroughly presents experiments and finds that back their hypothesis and findings such as when the model complexity is at the transition from under to over parameterization, increasing the number of samples shifts the peak to the right and actually increases test error. The paper's findings provide useful ways of thinking and explanation for model complexity and performance.

### **Weaknesses**

Many of the paper's findings cease to exist with optimal early stopping. Thus, models with good early stopping will not find much benefit from the findings of the paper. The paper does not provide a full understanding of optimal early stopping and double descent.

## **My Takeaways**

This fact that most of the paper's findings are negated in the presence of optimal early stopping makes me wonder how to achieve optimal early stopping. With different data and different architectures, optimal early stopping should be related to model complexity in relation to number of samples.

Collaborators:  
ML Theory textbook

# Assignment 3

Your name: Bojun Yang  
Your GTID: byang301

# Visualization

## Saliency Map (1 point)

hay



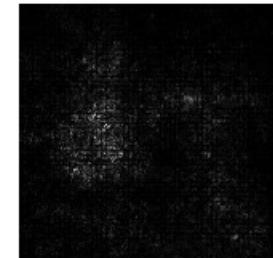
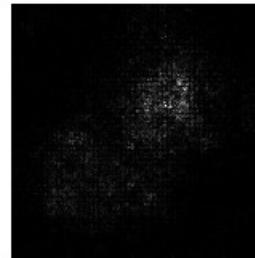
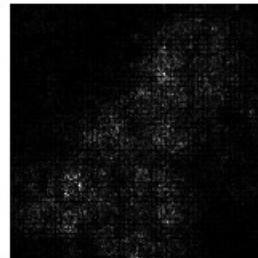
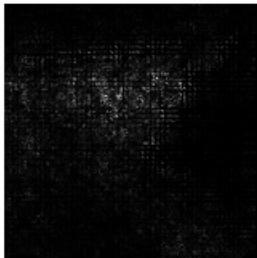
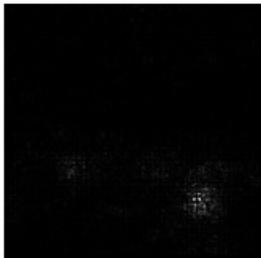
quail



Tibetan mastiff



Border terrier brown bear, bruin, Ursus arctos



## Saliency Map Captum (1 point)



hay



quail



Tibetan mastiff

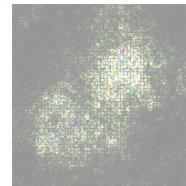
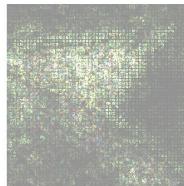
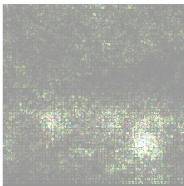


Border terrier



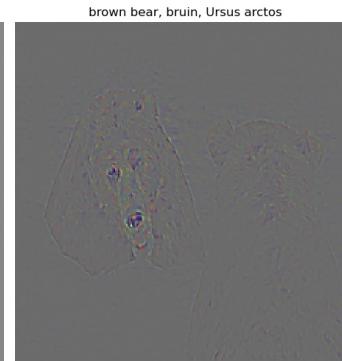
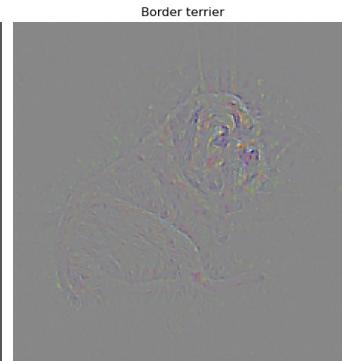
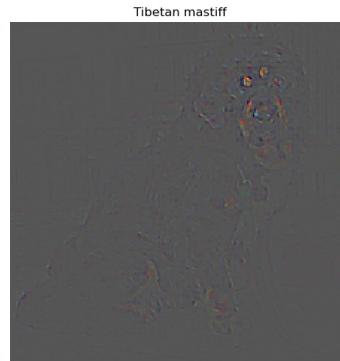
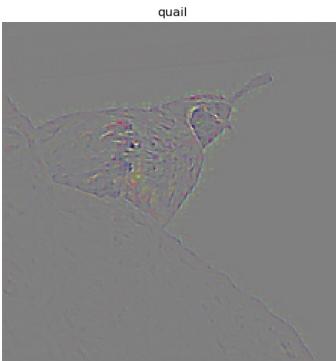
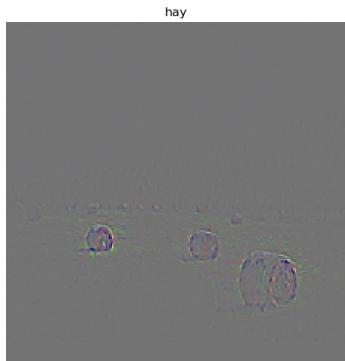
brown bear, bruin, Ursus arctos

Original Image



Saliency

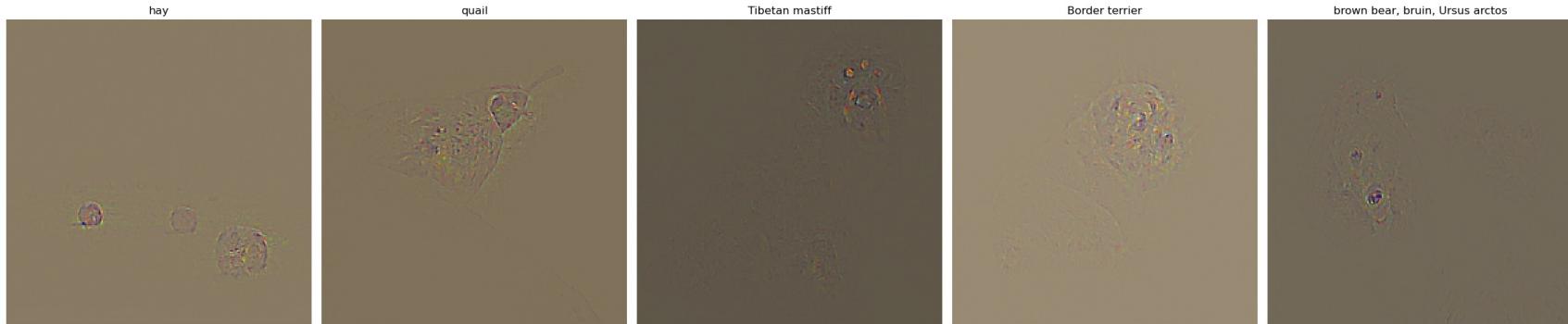
## Guided Backprop (1 point)



## GradCam (1 point)



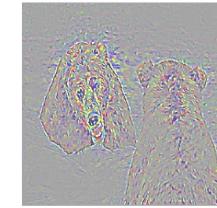
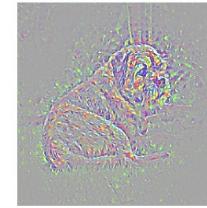
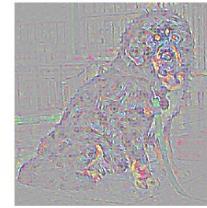
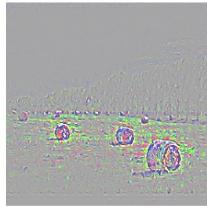
## Guided GradCam (1 point)



## Guided Backprop + GradCam (Captum) (1 point)



Original Image



Guided Backprop

hay



quail



Tibetan mastiff



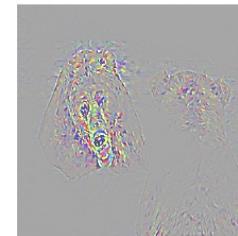
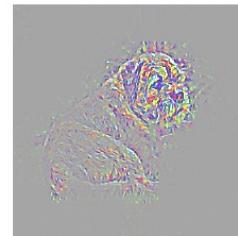
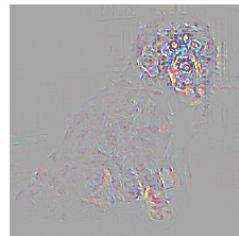
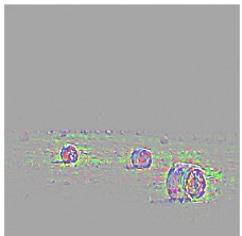
Border terrier



brown bear, bruin, Ursus arctos



Original Image

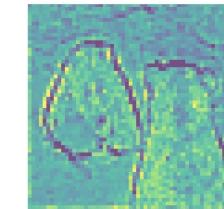
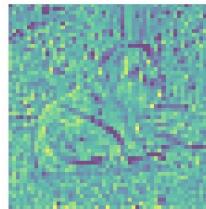
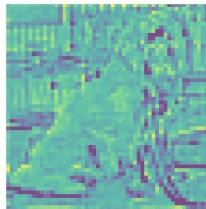
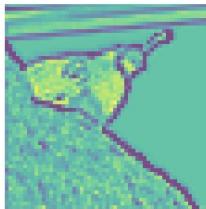
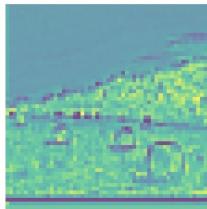


Guided GradCam

## Layers and neurons using Captum (1 point)



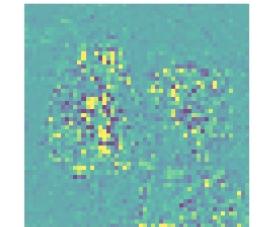
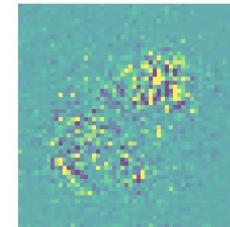
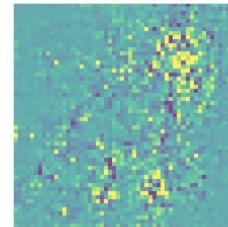
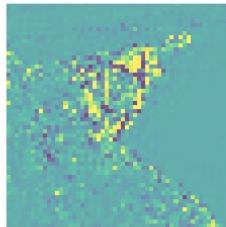
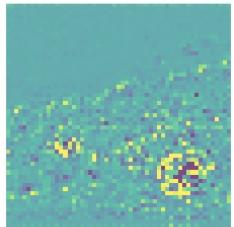
Original Image



Layer GradCam



Original Image



Layer Conductance

What do saliency map and Gradcam tell you? (1 point)

- Saliency map tells use the degree to which each pixel in the image affects the classification score of that image. Brighter pixels correspond to higher effect on the classification.
- Gradcam tells use where the network is looking when it is classifying an input image. It highlights the important regions within the input image for classification.

## Fooling Image (1 point)

hay



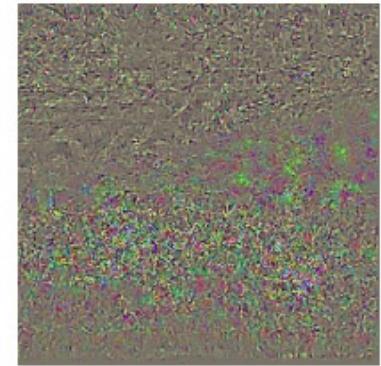
stingray



Difference



Magnified difference (10x)



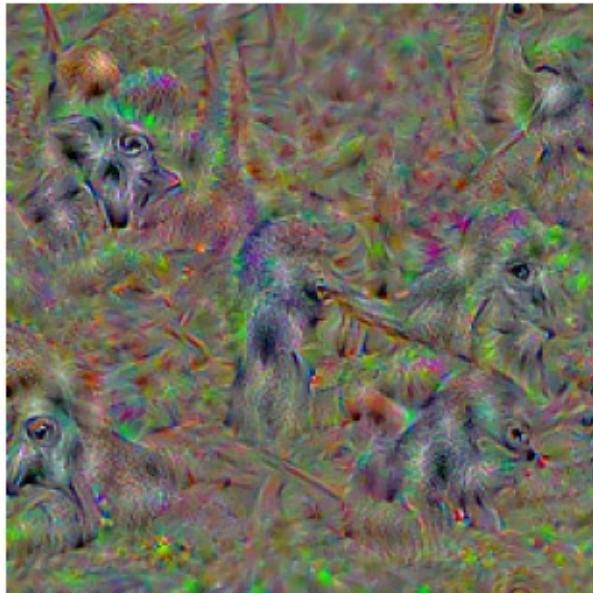
## Fooling Image Insights (1 point)

What insights do you get from fooling images:

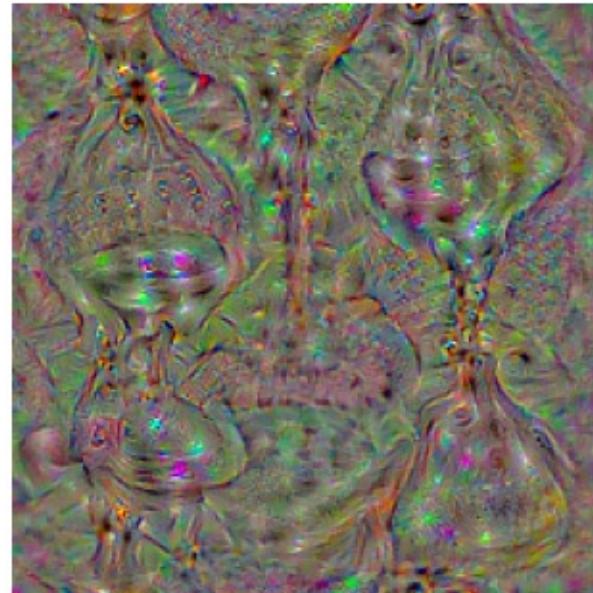
- Even though the fooling image would fool the network into classifying it as a stingray, it still looks the same as hay to the human. This tells us that networks can be easily tricked since it only takes very small differences in the pixels to confuse a nn.

## Class Visualization (3 points)

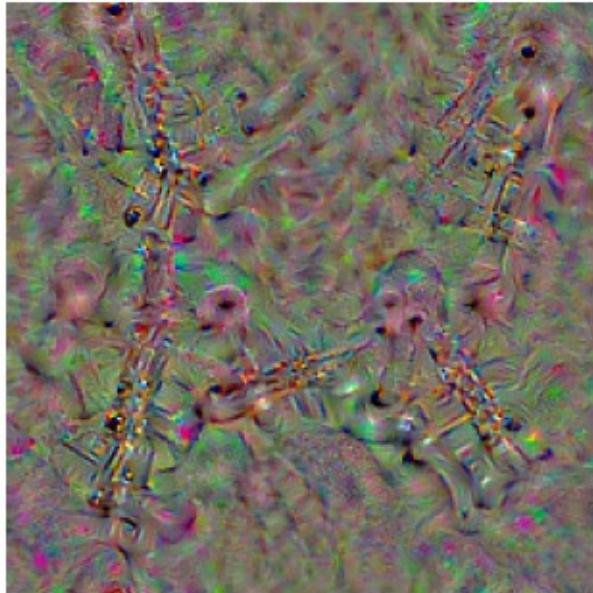
gorilla, Gorilla gorilla  
Iteration 100 / 100



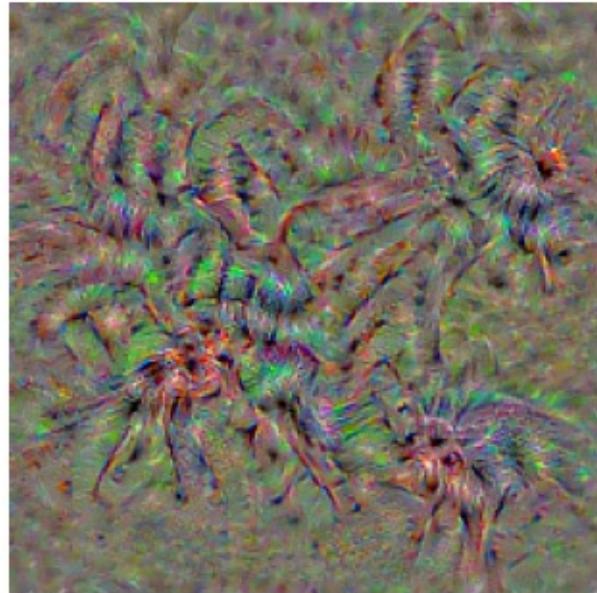
hourglass  
Iteration 100 / 100



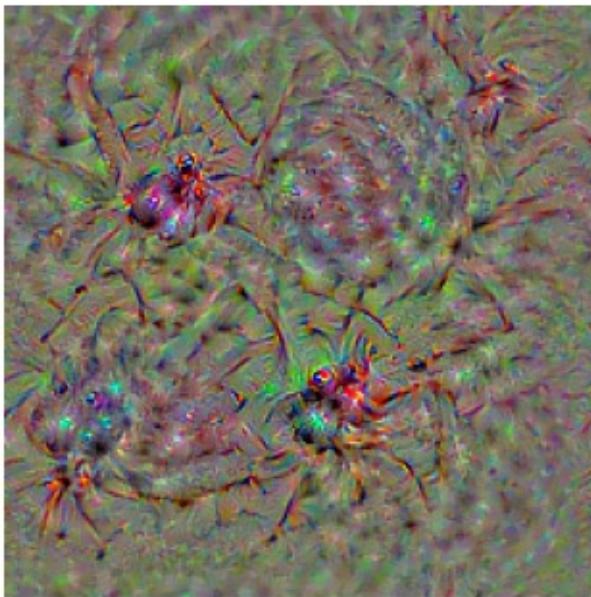
oboe, hautboy, hautbois  
Iteration 100 / 100



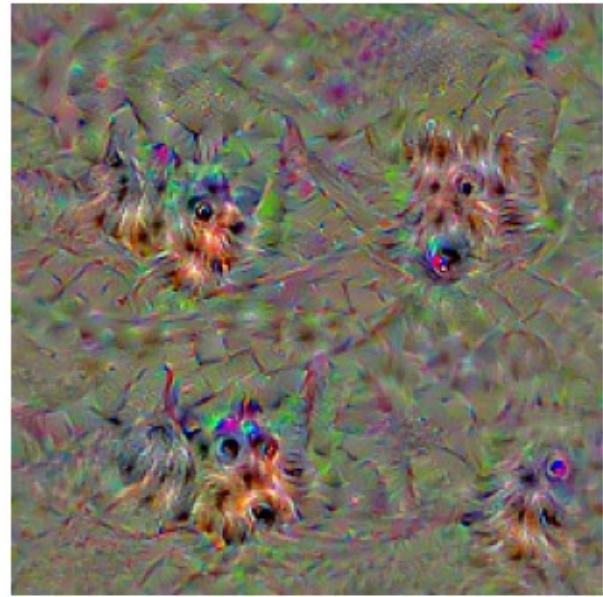
tarantula  
Iteration 100 / 100



tick  
Iteration 100 / 100



Yorkshire terrier  
Iteration 100 / 100



# Style Transfer

Content Source Img.

Composition VII + Tubingen (1 point)



Style Source Img.

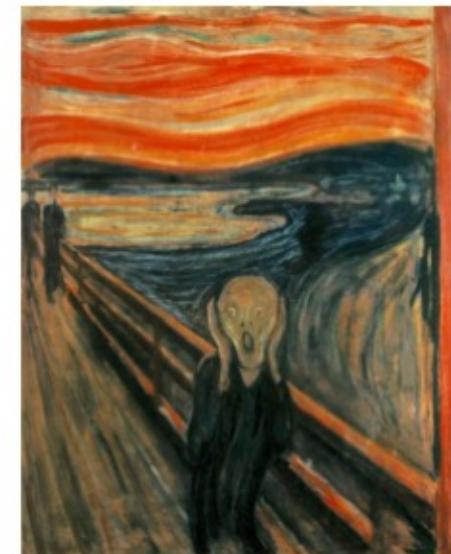


Content Source Img.

Scream + Tubingen (1 point)



Style Source Img.



Content Source Img.

Starry Night + Tubingen (1 point)



Style Source Img.



- Be sure to append the jupyter notebook i.e. \$root/test\_style\_transfer.ipynb for testing sections of the style transfer implementation to this report.

In [1]:

```
import numpy as np
import torch
from torch.autograd import Variable
import torchvision
import PIL

from style_modules import ContentLoss, StyleLoss, TotalVariationLoss
from style_utils import features_from_img, extract_features, rel_error, style_transfer
from image_utils import preprocess
```

In [2]:

```
def content_loss_test(correct, content_loss):
    content_image = 'styles_images/tubingen.jpg'
    image_size = 192
    content_layer = 3
    content_weight = 6e-2

    c_feats, content_img_var = features_from_img(content_image, image_size, cnn, dtype)

    bad_img = Variable(torch.zeros(*content_img_var.data.size()))
    feats = extract_features(bad_img, cnn)

    student_output = content_loss(content_weight, c_feats[content_layer], feats[content_layer]).data.numpy()
    error = rel_error(correct, student_output)
    print('Content Loss Maximum error is {:.3f}'.format(error))

def gram_matrix_test(correct, style_loss):
    style_image = 'styles_images/starry_night.jpg'
    style_size = 192
    feats, _ = features_from_img(style_image, style_size, cnn, dtype)
    student_output = style_loss.gram_matrix(feats[5].clone()).data.numpy()
    error = rel_error(correct, student_output)
    print('Gram Matrix Maximum error is {:.3f}'.format(error))

def style_loss_test(correct, style_loss):
    content_image = 'styles_images/tubingen.jpg'
    style_image = 'styles_images/starry_night.jpg'
    image_size = 192
    style_size = 192
    style_layers = [1, 4, 6, 7]
    style_weights = [300000, 1000, 15, 3]
```

```

c_feats, _ = features_from_img(content_image, image_size, cnn, dtype)
feats, _ = features_from_img(style_image, style_size, cnn, dtype)
style_targets = []
for idx in style_layers:
    style_targets.append(style_loss.gram_matrix(feats[idx].clone()))

student_output = style_loss(c_feats, style_layers, style_targets, style_weights).data.numpy()
error = rel_error(correct, student_output)
print('Style Loss Error is {:.3f}'.format(error))

def tv_loss_test(correct, tv_loss):
    content_image = 'styles_images/tubingen.jpg'
    image_size = 192
    tv_weight = 2e-2

    content_img = preprocess(PIL.Image.open(content_image), size=image_size)
    content_img_var = Variable(content_img.type(dtype))

    student_output = tv_loss(content_img_var, tv_weight).data.numpy()
    error = rel_error(correct, student_output)
    print('TV Loss Error is {:.3f}'.format(error))

```

In [3]:

```

dtype = torch.FloatTensor
# Uncomment out the following line if you're on a machine with a GPU set up for PyTorch!
# dtype = torch.cuda.FloatTensor

cnn = torchvision.models.squeezenet1_1(pretrained=True).features
cnn.type(dtype)

# Fix the weights of the pretrained network
for param in cnn.parameters():
    param.requires_grad = False

answers = np.load('data/style-transfer-checks.npz') # for local test

```

## Test content loss (1pt)

In [4]:

```

content_loss = ContentLoss()
content_loss_test(answers['cl_out'], content_loss) # should print out 0.0

```

Content Loss Maximum error is 0.000

## Test style loss (2pt)

In [5]:

```
style_loss = StyleLoss()  
gram_matrix_test(answers['gm_out'], style_loss) # should print out 0.0  
style_loss_test(answers['sl_out'], style_loss) # should print out 0.0
```

Gram Matrix Maximum error is 0.000  
Style Loss Error is 0.000

## Test total variation loss (1pt)

In [6]:

```
tv_loss = TotalVariationLoss()  
tv_loss_test(answers['tv_out'], tv_loss) # should print out 0.0
```

TV Loss Error is 0.000