

CS7643: Deep Learning
Fall 2020
HW3 Solutions

Tianxue Hu

October 7, 2020

Problem 1

$$W = \begin{bmatrix} w_{(0,0)} & w_{(0,1)} \\ w_{(1,0)} & w_{(1,1)} \end{bmatrix}$$

Stride=3 zero padding size = 1

$$X_{pad} = \begin{bmatrix} 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & X_{(0,0)} & X_{(0,1)} & X_{(0,2)} & 0 \\ 0 & X_{(1,0)} & X_{(1,1)} & X_{(1,2)} & 0 \\ 0 & X_{(2,0)} & X_{(2,1)} & X_{(2,2)} & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}, \quad \text{! ! are filters}$$

$$W_{\text{Pad}} = \begin{bmatrix} w_{(1,0)} x_{(0,0)} & w_{(1,0)} x_{(0,1)} \\ w_{(0,1)} x_{(1,0)} & w_{(0,0)} x_{(1,1)} \end{bmatrix}$$

$$Y = \begin{bmatrix} w_{(1,1)}X_{(0,0)}, w_{(1,2)}X_{(0,1)}, w_{(0,1)}X_{(1,0)}, w_{(0,2)}X_{(1,1)} \end{bmatrix}^T$$

$$= A \times$$

$$= \begin{bmatrix} W_{1,0} & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & W_{1,0} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & W_{0,1} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & W_{0,0} \end{bmatrix} \begin{array}{l} x(0,0) \\ x(0,1) \\ x(0,-1) \\ x(0,-2) \\ x(1,0) \\ x(1,1) \\ x(1,-1) \\ x(1,-2) \\ x(2,0) \\ x(2,1) \\ x(2,-1) \\ x(2,-2) \end{array}$$

Problem 2

$$w = \begin{bmatrix} w_{(0,0)} & w_{(0,1)} \\ w_{(1,0)} & w_{(1,1)} \end{bmatrix}$$

Stride=2 no padding

Affine transformation

$$\begin{array}{c} X_{(0,0)} \quad w_{(0,0)} \quad w_{(0,1)} \quad w_{(1,0)} \quad w_{(1,1)} \quad X_{(0,1)} \\ \hline X_{(1,0)} \quad w_{(0,0)} \quad w_{(0,1)} \quad w_{(1,0)} \quad w_{(1,1)} \quad X_{(1,1)} \\ \hline w_{(0,0)} \quad w_{(0,1)} \quad w_{(1,0)} \quad w_{(1,1)} \end{array}$$

$$Y = [X_{(0,0)} \ w_{(0,0)}, \ X_{(0,1)} \ w_{(0,1)}, \ X_{(0,0)} \ w_{(1,0)}, \ X_{(0,1)} \ w_{(1,1)}]$$

$$X_{(0,0)} \ w_{(1,0)}, \ X_{(0,1)} \ w_{(1,1)}, \ X_{(0,0)} \ w_{(1,0)}, \ X_{(0,1)} \ w_{(1,1)}$$

$$X_{(1,0)} \ w_{(0,0)}, \ X_{(1,0)} \ w_{(0,1)}, \ X_{(1,1)} \ w_{(0,0)}, \ X_{(1,1)} \ w_{(0,1)}$$

$$X_{(1,0)} \ w_{(1,0)}, \ X_{(1,0)} \ w_{(1,1)}, \ X_{(1,1)} \ w_{(1,0)}, \ X_{(1,1)} \ w_{(1,1)}]$$

$$= A \cdot X$$

$$= \left[\begin{array}{cccc|c} w_{(0,0)} & 0 & 0 & 0 & X_{(0,0)} \\ w_{(0,1)} & 0 & 0 & 0 & X_{(0,1)} \\ 0 & w_{(0,0)} & 0 & 0 & X_{(1,0)} \\ 0 & w_{(0,1)} & 0 & 0 & X_{(1,1)} \\ \hline w_{(1,0)} & 0 & 0 & 0 & \\ w_{(1,1)} & 0 & 0 & 0 & \\ 0 & w_{(1,0)} & 0 & 0 & \\ 0 & w_{(1,1)} & 0 & 0 & \\ \hline 0 & 0 & w_{(0,0)} & 0 & \\ 0 & 0 & w_{(0,1)} & 0 & \\ 0 & 0 & 0 & w_{(1,0)} & \\ 0 & 0 & 0 & w_{(1,1)} & \\ \hline 0 & 0 & 0 & w_{(1,0)} & \end{array} \right]$$

Problem 3

$$3. o=1, r=2, i=1, k=1 \dots$$

$$\text{conv layer } (o \times r^2, i, k, k) = (4, 1, 1, 1)$$

$$\text{conv layer } (o, r, k \times r, k \times r) = (1, 1, 2, 2)$$

want to prove affine transformation operation are identical for the two conv layers.

Affine transformation for $(4, 1, 1, 1)$

have $X = \begin{bmatrix} X_{(0,0)} & X_{(0,1)} \\ X_{(1,0)} & X_{(1,1)} \end{bmatrix}$

$$\text{Output channel } o' = 4$$

$$\text{Input channel } i' = 1$$

$$\text{filter row, col } k' = 1$$

so we have 4 input of 1×1 filter

let filter be

$$[w_1] [w_2] [w_3] [w_4]$$

$$\begin{array}{c}
 Y_i = \begin{bmatrix} w_1 & 0 & 0 \\ 0 & w_1 & 0 \\ 0 & 0 & w_1 \\ 0 & 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} X_{(0,0)} \\ X_{(0,1)} \\ X_{(1,0)} \\ X_{(1,1)} \end{bmatrix} = \begin{bmatrix} X_{00}w_1 \\ X_{01}w_1 \\ X_{10}w_1 \\ X_{11}w_1 \end{bmatrix} \\
 w_2 \quad 0 \quad 0 \quad 0 \\
 0 \quad w_2 \quad 0 \quad 0 \\
 0 \quad 0 \quad w_2 \quad 0 \\
 0 \quad 0 \quad 0 \quad w_2 \\
 w_3 \quad 0 \quad 0 \quad 0 \\
 0 \quad w_3 \quad 0 \quad 0 \\
 0 \quad 0 \quad w_3 \quad 0 \\
 0 \quad 0 \quad 0 \quad w_3 \\
 w_4 \quad 0 \quad 0 \quad 0 \\
 0 \quad w_4 \quad 0 \quad 0 \\
 0 \quad 0 \quad w_4 \quad 0 \\
 0 \quad 0 \quad 0 \quad w_4
 \end{array}$$

Affine transformation for $(1, 1, 2, 2)$

out put channel $o' = 1$

input channel $i' = 1$

filter row, col $k' = 2$

which gives 2×2 filter matrix

$$\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}$$

which is the same as problem 2. Given X

$$so Y_2 = \begin{bmatrix} w_1 & 0 & 0 & 0 \\ w_2 & 0 & 0 & 0 \\ 0 & w_1 & 0 & 0 \\ 0 & w_2 & 0 & 0 \\ w_3 & 0 & 0 & 0 \\ w_4 & 0 & 0 & 0 \\ 0 & w_3 & 0 & 0 \\ 0 & w_4 & 0 & 0 \\ 0 & 0 & w_1 & 0 \\ 0 & 0 & w_2 & 0 \\ 0 & 0 & 0 & w_1 \\ 0 & 0 & 0 & w_2 \\ 0 & 0 & w_3 & 0 \\ 0 & 0 & w_4 & 0 \\ 0 & 0 & 0 & w_3 \\ 0 & 0 & 0 & w_4 \end{bmatrix} \begin{bmatrix} X_{00} \\ X_{01} \\ X_{10} \\ X_{11} \end{bmatrix} = \begin{bmatrix} X_{00}w_1 \\ X_{00}w_2 \\ X_{01}w_1 \\ X_{01}w_2 \\ X_{10}w_3 \\ X_{10}w_4 \\ X_{11}w_1 \\ X_{11}w_2 \\ X_{11}w_3 \\ X_{11}w_4 \end{bmatrix}$$

So $Y_1 = Y_2$. proved that the affine transformation operation are identical.

Problem 4

4.

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

AND: $x_1 \quad x_2 \quad f_{\text{AND}}(x)$

0	0	0
0	1	0
1	0	0
1	1	1

One possible solution:

$$w_{\text{AND}} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \quad b_{\text{AND}} = -1$$

prove table

$$f(x) = w^T x + b = [0.5 \ 0.5] \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 1 = -1 < 0 \Rightarrow 0$$

$$\vdots \quad \quad = [0.5 \ 0.5] \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 1 = -0.5 < 0 \Rightarrow 0$$

$$\vdots \quad \quad = [0.5 \ 0.5] \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 1 = -0.5 < 0 \Rightarrow 0$$

$$= [0.5 \ 0.5] \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 1 = 0 = 0 \Rightarrow 1$$

OR : $x_1 \quad x_2 \quad f_{OR}(x)$

0	0	0
0	1	1
1	0	1
1	1	1

one possible sol'n:

$$w_{OR} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \quad b = -0.1$$

prove table

$$f(x) = w^T x + b = [0.5 \ 0.5] \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 = -0.1 < 0 \Rightarrow 0$$

$$[0.5 \ 0.5] \begin{bmatrix} 0 \\ 1 \end{bmatrix} - 0.1 = 0.4 > 0 \Rightarrow 1$$

$$[0.5 \ 0.5] \begin{bmatrix} 1 \\ 0 \end{bmatrix} - 0.1 = 0.4 > 0 \Rightarrow 1$$

$$[0.5 \ 0.5] \begin{bmatrix} 1 \\ 1 \end{bmatrix} - 0.1 = 0.9 > 0 \Rightarrow 1.$$

Problem 5

	x_1	x_2	$f_{\text{XOR}}(x)$
1	0	0	0
2	0	1	1
3	1	0	1
4	1	1	0

From table

$$f_{\text{XOR}}([0]) > f_{\text{XOR}}([0])$$

$$w^T [0] + b > w^T [0] + b$$

$$w_2 + b > b$$

$$w_2 > 0 \quad (1)$$

$$f_{\text{XOR}}([1]) > f_{\text{XOR}}([1])$$

$$w^T [1] + b > w^T [1] + b$$

$$w_2 + b > w_1 + w_2 + b$$

$$w_1 < 0 \quad (2)$$

$$f_{\text{XOR}}([0]) > f_{\text{XOR}}([1])$$

$$w^T [0] + b > w^T [1] + b$$

$$w_1 + b > w_1 + w_2 + b$$

$$w_2 < 0 \quad (3)$$

Then (1) and (3) are conflict.

Thus XOR cannot be represented using linear model



Problem 6

$$6. h(x) = w^{(3)} \max\{0, w^{(2)} \max\{0, w^{(1)}x + b^{(1)}\} + b^{(2)}\} + b^{(3)}$$

$$w^{(1)} = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix} \quad b^{(1)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

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

$$w^{(3)} = \begin{bmatrix} 1 & 1 \end{bmatrix} \quad b^{(3)} = [1]$$

$$x = 1$$

$$h(x) = [1 \ 1] \max\{0, [1 \ 1] \max\{0, [0.5 \ 0.5] [1] + [0 \ 1]\} + [0 \ 0]\} + [1]$$

$$\begin{array}{c} \downarrow \\ \boxed{\begin{bmatrix} 0.5 \\ 1.5 \end{bmatrix}} \end{array}$$

$$\downarrow$$

$$\boxed{\begin{bmatrix} 2 \\ 2 \end{bmatrix}}$$

$$= [1 \ 1] \boxed{\begin{bmatrix} 2 \\ 2 \end{bmatrix}} + 1 = 5$$

$$h(x) = Wx + b = 3x + 2 \quad \text{where } w=3 \quad b=2$$

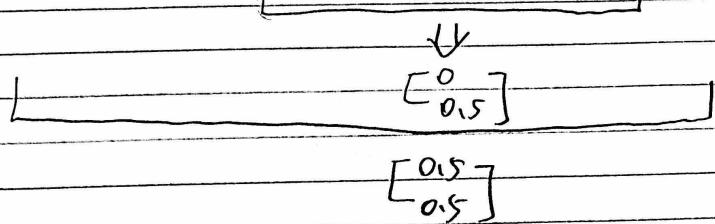
$$\frac{dh}{dx} = 3 \quad h(x) = 3x + 2$$

Problem 7

7.

$$x = -1$$

$$h(x) = \lceil -1 \rceil \max\{0, \lceil -1 \rceil \max\{0, \lceil \frac{0.5}{0.5} \rceil \lceil -1 \rceil + \lceil 0 \rceil \} + \lceil 0 \rceil \} + 1$$



$$= \lceil -1 \rceil \lceil \frac{0.5}{0.5} \rceil + 1 = 2$$

A possibility could be:

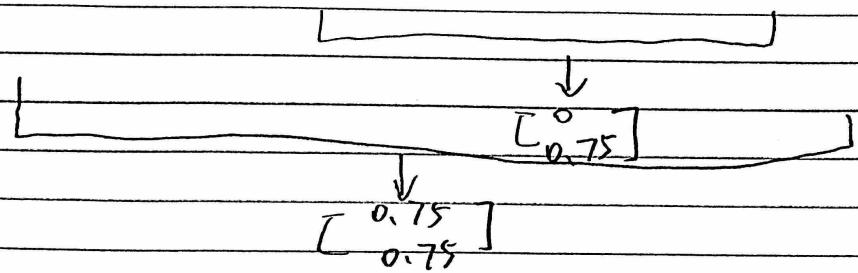
$$h(x) = w \cdot x + b \quad \text{where } w=1, b=3$$

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

Problem 8

$$8. \quad x = -0.5$$

$$h(x) = [1] \max\{0, [1] \max\{0, \underbrace{[0.5] \}_{0.5} [-0.5] + [0] \}_{1} \}_{0} + [0] \}_{0}$$



$$= 2.5$$

A possibility would be

$$h(x) = x + 3 \quad \text{when } w=1 \quad b=3$$

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

Problem 9

9.

$$w_{ij}^{(1)} = \begin{cases} 2 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases} \quad w = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$b_i^{(1)} = -1 \quad b = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$$

when
 $d=1 \quad O = (0,1)^T$ input space $R = \{-1, 0, 1\}$
 $= 2$

$R_1 = (-1, 0) \quad R_2 = (0, 1)$ input regions

when $d=2$, we choose either R_1 or R_2 for corresponding dimensions. Then there are $2^2 = 4$ input regions

Similarly, for d dimension, there are 2^d input regions.

Problem 10

10. g is n^d regions of $(0,1)^d$ onto $(0,1)^d$
 f is n^d regions of $(0,1)^d$ onto $(0,1)^d$

$f \circ g(1)$ means for each n^d regions has n^d regions.
Thus, there are $(0,1)^d = n^d$ input regions
for $f \circ g(1)$

Problem 11

ii.

$$h_1(x) = |w_1x + b_1|$$

$$h_2(x) = |w_2h_1 + b_2| = |w_2(w_1x + b_1) + b_2|$$

$$h_3(x) = |w_3(w_2(w_1x + b_1) + b_2) + b_3|$$

:

$$f(x) = h_L(x) = |w_Lh_{L-1} + b_L|$$

where $x \in (0,1)^d$

we proved in 9 that each layer has 2^d input regions for d dimensions.

we prove by induction

Base: for h_1 , there are 2^d regions $2^d = 2^{1 \times d} = 2^{2^d}$ when $L=1$.

Inductive Case: assume for h_k there're $2^{k \cdot d}$ input regions, we want to prove for h_{k+1} there're $2^{(k+1) \cdot d}$ input regions.

$$\text{for } L_{k+1}, \dots, h_{k+1} = h_{k+1}(h_k(x))$$

which has $2^d \cdot 2^{kd}$ input regions

$$= 2^{(k+1)d}$$

inductive case proved.

proved.

Problem 12

12.

The optimization problem by least square method is

$$w = \arg \min_{w \in \mathbb{R}^n} \sum_{i=1}^n (y_i - w^T x_i)^2$$

to get optimized sol'n.

we take derivative wrt. to w and set to 0

$$\frac{\partial}{\partial w} \sum (y_i - w^T x_i)^2$$

$$= 2 \sum_i x_i (y_i - w^T x_i) = 0$$

$$\Rightarrow 2 \sum_i x_i (y_i - w^T x_i) = 0$$

$$2 \sum_i x_i y_i - 2 \sum_i w^T x_i x_i = 0$$

$$\sum_i x_i y_i = \sum_i w^T x_i^2$$

$$w_{gd} = \frac{\sum_i x_i y_i}{\sum_i x_i^2}$$

$$\text{as } x_w = y$$

$$w_{gd} = \frac{\sum_i \frac{y_i}{w_i} y_i}{\sum_i x_i^2} = \frac{1}{w_i} \cdot \frac{\sum_i y_i^2}{\sum_i x_i^2}$$

$$= \frac{1}{w_i} \sum_i w_i^2$$

$$= \arg \min_{w \in \mathbb{R}^n} \|w\|_2^2$$

Problem 13

As $\|x\|_2^2 = \sum_i x_i$, (25) is calculating the optimized gradient descent by finding at which w , the sum of w_i has the least value.

Problem 14

(23) tries to find the optimiaized w that have the least sum of differences between actual y and prediction $y_{hat} = w^T x$. The solution is not necessarily unique since $f(w)$ may falls into an local minimum and the result depends on w_0 .

Problem 15

Approximation error is the minimum generalization error achievable by a predictor in the hypothesis class Estimation error is minimize training error being the only estimate of predictor minimizing the generalization error.

Problem 16

Briefly summarize the key contributions, strengths and weaknesses of this paper.

Key contributions: This paper shows that double descent is a robust phenomenon that occurs in a variety of tasks, architectures, and optimization methods. They defined the effective model complexity (EMC) of a training procedure as the maximum number of samples on which it can achieve close to zero training error. And increasing training time will increase the EMC. The phenomenon indicates that training a deep network on a larger train set performs worse.

Strength: The paper investigate the double decent effect in different scenarios: with/without data augmentation, different optimization methods, different number of epochs, different noise settings, and trained on different kinds of networks. Their experiments making the conclusion reputable and universal.

Weakness: The paper proposed several EMC regimes: under-parameterized regime, over-parameterized regime, and critically parameterized regime. They indicates that in the critically parameterized regime the double descent phenomenon causes poor performance when there are more training data. However, they failed to find out how to define the critical interval if others want to apply it. The paper says "The width of the critical interval depends on both the distribution and the training procedure in ways we do not yet completely understand".

Problem 17

What is your personal takeaway from this paper? This could be expressed either in terms of relating the approaches adopted in this paper to your traditional understanding of learning parameterized models, or potential future directions of research in the area which the authors haven't addressed, or anything else that struck you as being noteworthy.

Despite the conventional wisdom of over-fitting and under-fitting usually happens in statistical modeling, deep learning researchers believe the larger the model the better the performance. This paper discovers the double descent phenomenon that the performance of deep learning models also has a "negative" relationship to the number of training data and model size (when EMC is in a critical interval). As mentioned in the weakness part, the researchers still cannot understand how to find the critical interval, I hope it will be defined in the future for effective model training.