#### האוניברסיטה העברית בירושלים

בית הספר להנדסה ולמדעי המחשב ע'ש רחל וסלים בנין

# $\mathsf{C++r}$ סדנאות תכנות בשפת C++רס (67315) סדנאות תכנות בשפת - C++

22:00 עד השעה של - 2022 במאי, ביעי, ה־25 במאי, יום רביעי, יום רביעי, ה

נושאי התרגיל: Introduction to C++, Classes, Operator Overloading, References, Rule of Three, Exceptions

אנא הקפידו לקרוא את כל התרגיל מתחילתו ועד סופו לפני שתגשו לממשו.

## רקע 1

בתרגיל זה נכתוב תוכנה לזיהוי ספרות הנכתבות בכתב יד. התוכנה שלנו תקבל כקלט תמונה של ספרה בין 0 ל־9 ותחזיר כפלט את הספרה אשר זוהתה. נעשה זאת על ידי בניית מודל של רשת נוירונים. הרשת שנריץ תגיע לדיוק של כ־96 אחוזים בזיהוי ספרות.

#### 1.1 הקדמה

רשת נוירונים היא מודל בלמידת מכונה המבוסס על מבנה המוח האנושי: נוירון מקבל גירוי חשמלי מנוירונים אחרים <sup>-</sup> אם הגירוי הזה עובר סף מסוגלים מסוים הוא שולח בעצמו אות אל נוירונים אחרים. המוח מורכב ממספר רב של נוירונים המקושרים זה לזה ברשת מורכבת, ויחד הם מסוגלים לבצע את הפעולות הנדרשות ממנו. רשת נוירונים מלאכותית(Artificial neural network) פועלת באופן דומה. ברשתות אלו נעשה שימוש בזיהוי ומיקום של עצמים בתמונה, הבנת שפה אנושית וניתוחה, יצירת טקסט ועוד. מוצרים רבים בחיינו משתמשים ברשתות נוירונים: עוזרים קוליים(Amazon Alexa, Apple Siri), השלמה אוטומטית לתוכן המייל ב־Gmail, זיהוי מחלות בתמונות סריקה רפואית ועוד.

שימו לב: למידת מכונה בכלל, ורשתות נוירונים בפרט, הינם נושאים רחבים ומורכבים ולכן לא יכללו בתוכן תרגיל זה. לצורך מימוש התרגיל אין צורך להבין איך ולמה עובדת רשת הנוירונים. הרקע התיאורטי הנדרש יוצג בסעיף 1.2, פרטי הרשת שנבנה יובאו בסעיף 2.2, והמחלקות למימוש יובאו בסעיף 3.

מומלץ לצפות בסרטון הבא המפרט על המבנה של רשת נוירונים וכיצד ניתן לממש אותה באמצעות אלגברה לינארית: https://www.youtube.com/watch?v=aircAruvnKk

#### רקע תיאורטי 1.2

#### Fully Connected רשת נוירונים 1.2.1

- רשת בנויה משכבות (סעיף 1.2.2).
- הקלט של כל שכבה הוא וקטור, והפלט הוא וקטור אחר.
  - הפלט של כל שכבה הוא הקלט של השכבה הבאה.
- קלט הרשת הוא וקטור המייצג את האובייקט שהרשת תעבד. ברשת שלנו, הוא מייצג תמונה של ספרה (סעיף2.2.2).
- פלט הרשת הוא וקטור המייצג את המסקנה של הרשת. ברשת שלנו, הוא מייצג את הספרה שהרשת זיהתה (סעיף 2.2.3).

#### 1.2.2 שכבה ברשת

כל שכבה ברשת מקבלת וקטור  $y \in \mathbb{R}^n$  ומחזירה וקטור ומחזירה וקטור באמצעות הפעולה המתמטית כל שכבה ברשת מקבלת וקטור באהי

$$y = f(W \cdot x + b)$$

:כאשר

- .(Weights) מטריצה שאיבריה נקראים המשקולות של מטריצה  $W \in M_{n \times m}(\mathbb{R})$ 
  - של השכבה. (Bias) איז שנקרא השנקרא וקטור שנקרא  $b \in \mathbb{R}^n$
  - $(1.2.3 \; ext{even})$  פונקציית האקטיבציה של השכבה פונקציית  $f: \mathbb{R}^n o \mathbb{R}^n \; ullet$

. כלומר, בהינתן וקטור קלט $y=f(W\cdot x+b)\in\mathbb{R}^n$  הוקטור הא הפלט של השכבה. כלומר, בהינתן האינתן וקטור הא השכבה.

#### 1.2.3 פונקציית אקטיבציה

פונקציות בתרגיל בתרגיל בתרגיל ממש שתי פונקציות של שכבה ברשת הנוירונים. פונקציה זו אינה לינארית. בתרגיל נממש שתי פונקציות  $f:\mathbb{R}^n o \mathbb{R}^n$  אקטיבציה שונות:

ReLU פונקציית

$$\forall x \in \mathbb{R}$$
  $ReLU(x) = \begin{cases} x & x \ge 0 \\ 0 & else \end{cases}$ 

. כאשר הפונקציה פועלת על וקטור את מבצעת את מבצעת איז  $x\in\mathbb{R}^n$ היא פועלת על הפונקציה הפונקציה איז מבצעת את האיז איז וקטור

Softmax פונקציית

$$\forall x \in \mathbb{R}^n \qquad Softmax(x) = \frac{1}{\sum_{k=1}^n e^{x_k}} \begin{bmatrix} e^{x_1} \\ e^{x_2} \\ e^{x_3} \\ \vdots \\ e^{x_n} \end{bmatrix} \in \mathbb{R}^n$$

 $x \in \mathbb{R}^n$  של -  $x_l$ 

.cmath המיובאת std::exp הפונקציה בפונקציה להשתמש ביתן לצורך החישוב מ־exp(t) המיובאת -  $e^t$ 

הפונקציה מקבלת וקטור  $x\in\mathbb{R}^n$  וממירה אותו לוקטור התפלגות (וקטור שאיבריו הם מספרים אי־שליליים שסכומם 1) באופן שתואם את הפלט הסופי של הרשת שלנו.

#### מימוש הרשת

## ודיוק הרשת float: שימוש ב־2.1

.float (32-bit) איברי המטריצה שנממש יהיו מטיפוס

מאופן המימוש של float במעבד, פעולות האריתמטיקה אינן בהכרח אסוציאטיביות, כלומר לא בהכרח יתקיים:

$$(a+b) + c = a + (b+c)$$

לכן, כדי להימנע משגיאות נומריות בעת ביצוע כפל מטריצות, אנא ממשו את סדר הפעולות לפי ההגדרה המתמטית שלמדתם בלינארית 1:

$$(A \cdot B)_{ij} = \sum_{k=1}^{n} a_{ik} b_{kj}$$

הרשת שנריץ מגיעה ל**כ־96 אחוזי דיוק**. לכן, הרשת עלולה לטעות בחלק מהתמונות שתזינו לה - זוהי התנהגות תקינה של התוכנה. גם אם אחוזי ההצלחה של הרשת שלכם נמוכים במעט מ־96 אחוזים, ייתכן כי הדבר נובע משגיאות נומריות, ולא צפויה הורדה של נקודות במקרה זה. עם זאת, אם אחוזי ההצלחה של הרשת נמוכים משמעותית מרף זה, ייתכן שהדבר נובע מטעות במימוש.

#### 2.2 תיאור הרשת

#### 2.2.1 שכבות הרשת

• הרשת מורכבת מ־4 שכבות:

פונקציית אקטיבציה	Bias - היסט	משקולות - Weights	שכבה
Relu	$b_1 \in \mathbb{R}^{128}$	$W_1 \in M_{128 \times 784}$	(כניסה) 1
Relu	$b_2 \in \mathbb{R}^{64}$	$W_2 \in M_{64 \times 128}$	2
Relu	$b_3 \in \mathbb{R}^{20}$	$W_3 \in M_{20 \times 64}$	3
Softmax	$b_4 \in \mathbb{R}^{10}$	$W_4 \in M_{10 \times 20}$	4 (מוצא)

• כלומר, על מנת לממש את הרשת עלינו לשרשר את רצף הפעולות הבא:

$$r_1 = Relu(W_1 \cdot x + b_1)$$

$$r_2 = Relu(W_2 \cdot r_1 + b_2)$$

$$r_3 = Relu(W_3 \cdot r_2 + b_3)$$

$$r_4 = Softmax(W_4 \cdot r_3 + b_4)$$

- .(2.2.2 סעיף x וקטור הקלט לרשת (סעיף ).
- .i+1הפלט של השכבה ה־, שהוא גם הקלט לשכבה ה־  $r_i$
- .(2.2.3 סעיף) הפלט של הפלט הפרטית, שהוא הרביעית, שהוא הפלט של הרשת  $r_4$

### 2.2.2 וקטור הקלט

- כל מספר במטריצה הוא ערך בין 0 ל־1, כלומר (Grayscale), של פיקסלים בגווני אפור של 28  $\times$  28 של בגודל A בגודל מטריצה בתור מטריצה  $A\in M_{28\times28}\left([0,1]\right)$ 
  - לנוחיותכם, בקבצי העזר נמצא הקובץ plot img.py אשר מקבל כקלט נתיב לתמונה ומציג אותה בחלון חדש.
    - . וקטור  $28 \cdot 28 = 748$  עם עמודה אחת) וקטור (מטריצה עם שיישלח לרשת יהיה וקטור (מטריצה עם אחת)

#### 2.2.3 וקטור הפלט

- . (10 באורך שסכומם באי־שליליים שסכומם באורך ווען התפלגות הפלט מהשכבה האחרונה הוא וקטור התפלגות (וקטור שאיבריו הם מספרים אי־שליליים שסכומם 1
  - כל אינדקס בווקטור מייצג ספרה בין 0 ל־9.
  - הערך של האינדקס מייצג את הסיכוי שזוהי הספרה בתמונה, לפי הרשת.
- התשובה שתיתן הרשת היא האינדקס עם הערך המקסימלי, כלומר הספרה הסבירה ביותר, וההסתברות של אותה ספרה
  - במקרה של שיוויון, נחזיר את האינדקס הנמוך מבין השניים
    - לדוגמא:

Valu	C	)	0.003	0.08	0	0	0	0	0.9	0.007	0.01
Inde	; C	)	1	2	3	4	5	6	7	8	9

90% בהסתברות בתמונה היא 7 בהסתברות הינתן וקטור הפלט הזה, תשובת הרשת תהיה שהספרה בתמונה היא

#### 2.3 מהלך ריצת התוכנית

 $\frac{\mathsf{wight}}{\mathsf{main.cpp}}$  שעליכם לממש על מנת מומש עבורכם בקובץ  $\frac{\mathsf{main.cpp}}{\mathsf{main.cpp}}$  שעליכם לממש על מנת לב: חלק זה ממומש עבורכם בקובץ לכם עם קבצי התרגיל להשתמש.

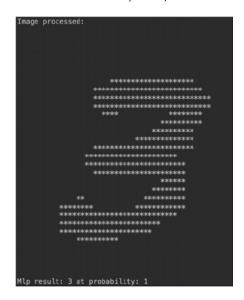
התוכנה תקבל בשורת הפקודה נתיבים לקבצי המשקולות וההיסטים כקבצים בינאריים. נריץ את התוכנה עם המשקולות וההיסטים כך:

\$./mlpnetwork w1 w2 w3 w4 b1 b2 b3 b4

- iרמיב לקובץ המשקולות של השכבה ה־ $w_i$ 
  - iמתיב לקובץ ההיסט של השכבה ה־ $b_i$

כאשר התוכנה רצה היא ממתינה לקלט מהמשתמש. הקלט יהיה נתיב לקובץ של תמונה המכילה ספרה. התוכנה:

- 1. תפתח את הקובץ ותטען את התמונה למטריצה
  - 2. תכניס את המטריצה אל הרשת כקלט
- 3. כאשר התקבלה תוצאה, התוכנה תדפיס את התמונה, את הספרה שהרשת זיהתה ובאיזו הסתברות. לדוגמא (מתוך פתרון בית הספר):



4. התוכנה תמתין לקלט חדש

.0 כאשר נזין לתוכנה q - חתוכנה תצא עם קוד

## 3 המחלקות למימוש

הינכם נדרשים לממש את המחלקות הבאות בלבד. אין להגיש מחלקות נוספות.

- בטבלאות המובאות לפניכם רשומות כלל הפונקציות והאופרטורים שעליכם לממש.
- אין להרחיב את ה־API המפורט, כלומר אין להוסיף פונקציות public אין להרחיב את ה-API.
- חישבו היטב על החתימה של כל פונקציה: מהו ערך ההחזרה שלה? האם היא משנה את האובייקט הנוכחי? כיצד נגדיר את הטיפוס של הארגומנטים שלה?
- שימו לב: <u>עליכם להחליט</u> איפה יש להשתמש ב־const, איפה המשתנים וערכי ההחזרה צריכים להיות by value והאם לממש (standalone, non-member function) או כפונקציה העומדת בפני עצמה (member function) או כפונקציה העומדת בפני שצמה
  - . std:: vector של איסור על שימוש ב־STL: בתרגיל זה אין להשתמש בספריית איסור על שימוש ב-STL: בתרגיל במבני נתונים כמו

## Matrix המחלקה 3.1

. float מחלקה זו מייצגת אובייקט של מטריצה (גם וקטור הינו מטריצה, בעלת עמודה אחת ו־n שורות). נזכיר כי איברי המטריצה יהיו מטיפוס

Description	Name	Comments				
	Constructors					
Constructor	Matrix(int rows, int cols)	Constructs Matrix of size rows×cols.				
		Inits all elements to 0.				
Default Constructor	Matrix()	Constructs Matrix of size 1×1.				
		Inits the single element to <b>0</b> .				
Copy Constructor	Matrix(Matrix m)	Constructs matrix from another Matrix m.				
Destructor	~Matrix()					

		Methods & Functions
Getter	get_rows()	returns the amount of rows as int.
Getter	get_cols()	returns the amount of cols as int.
	transpose()	Transforms a matrix into its <b>transpose</b> matrix,
		i.e $(A.transpose())_{ij} = A_{ji}$ .
		Supports function calling concatenation.
		e.g:
		Matrix $a(5,4), b(4,5);$
		$a.transpose();// a.get\_rows == 4, a.get\_cols == 5$
		b.transpose().transpose();// b is same as before
	vectorize()	Transforms a matrix into a column vector(section 3.1.2).
		Supports function calling concatenation.
		e.g.:
		Matrix $m(5,4)$ ;
		m.vectorize();
		$m.get\_cols() == 1$
		$m.get\_rows() == 20$
	$plain\_print()$	Prints matrix elements, no return value.
		Prints space after each element (including last
		element in row).
		Prints newline after each row (including last row).
	dot(Matrix m)	Returns a matrix which is the <b>elementwise</b>
		multiplication (Hadamard product) of this matrix and
		another matrix m:
		$\forall i, j : (A.dot(B))_{ij} = A_{ij} \cdot B_{ij}$
	$\operatorname{norm}()$	Returns the Frobenius norm of the given
		$\max_{i} A.norm() = \sqrt{\sum_{i} A_{ij}^2}$
		$\bigvee i,j$

	Operators			
+	Matrix addition	$\text{Matrix a, b;} \rightarrow \text{a + b}$		
=	Assignment	$Matrix a, b; \rightarrow a = b$		
*	Matrix multiplication	Matrix a, b; $\rightarrow$ a * b		
*	Scalar multiplication on the right	Matrix m; float $c; \rightarrow m * c$		
*	Scalar multiplication on the left	Matrix m; float $c; \rightarrow c * m$		
+=	Matrix addition accumulation	$\text{Matrix a, b;} \rightarrow \text{a += b}$		
()	Parenthesis indexing	For i,j indices, Matrix m:		
		m(i,j) will return the i,j element in the matrix		
		e.g.		
		Matrix $m(5,4)$ ;		
		m(1,3) = 10;		
		float $x = m(1,3); // x == 10$		
	Brackets indexing	For i index, Matrix m:		
		m[i] will return the i'th element (section 3.1.2)		
		e.g.		
		Matrix $m(5,4)$ ;		
		m[3] = 10;		
		float $x = m[3]; // x == 10$		
<<	Output stream	Pretty export of matrix as per section 3.1.1		
>>	Input stream	Fills matrix elements: has to read input stream fully, otherwise		
		it's an error (don't trust the user to validate it).		
		see section 3.1.3 for more details.		

## 3.1.1 הדפסת תמונה ממטריצה

בא: בפסאודו־קוד בפסאודו , << באמצעות אופרטור במטריצה במטריצה במטריצה התמונה את כדי להדפיס התמונה במטריצה במטריצה במטריצה בא

```
\begin{array}{lll} for & i = 1 \ to \ A. \, rows: \\ & for \ j = 1 \ to \ A. \, cols: \\ & & if \ A(i\,,j\,) > 0.1: \\ & & print "**" \ (double \ asterisk) \\ & & else: \\ & & print "" \ \ (double \ space) \\ & & print \ newline \end{array}
```

## אינדקס יחיד לזכרון דו־מימדי 3.1.2

מטריצה A הינה אובייקט דו־מימדי, ולכן אנו זקוקים לשני אינדקסים על מנת לגשת לאיבר בה. פעמים רבות, נוח יותר לגשת לכל איבר במטריצה באמצעות אינדקס יחיד. נבצע את המיפוי מזוג אינדקסים לאינדקס יחיד באופן הבא:

$$A(i, j) == A[i \cdot \text{rowsize} + j]$$

- אינדקס השורה i
- אינדקס העמודה j

(מספר העמודות) אורך - rowsize

לדוגמה, תהי  $M_{3 imes 4}$ , כלומר בעלת 3 שורות ו־4 עמודות. מתקיים:

$$A(2,1) == A[2 \cdot 4 + 1] == A[9]$$

וודאו שאתם מבינים מדוע מיפוי זה הינו חד־חד־ערכי ועל.

#### 3.1.3 קריאת תמונה מקובץ בינארי למטריצה

כדי לקרוא קובץ בינארי מתוך Input stream, עליכם להשתמש בפונקציה std::istream::read, שמקבלת ∗rhar ומספר בתים∖תווים לקרוא. תוכלו למצוא תיעוד מלא של הפונקציה בלינק הבא:

https://www.cplusplus.com/reference/istream/istream/read/

.char∗ במקרה הזה תוכלו לעשות casting מפורש ל

## Activation קבצי 3.2

בקבצים אלה נגדיר את פעולת פונקציות האקטיבציה. בקבצי ה־Activation עליכם לממש את פונקציות האקטיבציה וידצה softmax ורבא יהיטות בקבצים אלה נגדיר את פעולת פונקציות האקטיבציה. בקבצי ה-Activation::relu או Activation::relu. השייכות ל-namespace Activation כלומר, על מנת לקרוא לפונקציות צריך לרשום typedef של פוינטר לפונקצית אקטיבציה שימו לב שנשתמש בתרגיל זה ב־function pointers על מנת לגשת לפונקציות אלו. מומלץ להגדיר typedef של פוינטר לפונקצית אקטיבציה ולהשתמש בו בתוכנית.

## Dense המחלקה 3.3

מחלקה זו מייצגת שכבה ברשת הנוירונים.

Description	Name	Comments
Constructor	Dense(weights, bias, ActivationFunction)	Inits a new layer with given parameters.
		C'tor accepts 2 matrices and activation function
	Methods	
Getter	get_weights()	Returns the weights of this layer.
Getter	get_bias()	Returns the bias of this layer.
Getter	get_activation()	Returns the activation function of this layer.
	Operators	
()	Parenthesis	Applies the layer on input and returns output matrix
		Layers operate as per section 2.2.1
		e.g:
		Dense layer(w, b, act);
		Matrix output = layer(input);

## MlpNetwork המחלקה 3.4

מחלקה זו תשמש אותנו לסדר את השכבות השונות למבנה רשת ותאפשר הכנסה של קלט לרשת וקבלת הפלט המתאים. מחלקה זו מממשת ספציפית את הרשת המתוארת במסמך זה (סעיף 2.2.1).

שימו לב כי struct digit מומש עבורכם בקבצים שקיבלתם.

<u>נקודה למחשבה</u>: מה היה נדרש לממש במחלקה זו על מנת לתמוך ברשת עם מספר שכבות וגודל שכבות הניתן ב**זמן ריצה**?

Description	Name	Comments
Constructor	MlpNetwork(weights[], biases[])	Accepts 2 arrays of matrices, size 4 each.
		one for weights and one for biases.
		constructs the network described (sec. 2.2)
	Operators	
()	Parenthesis	Applies the entire network on input.
		returns digit struct.
		MlpNetwork m();
		digit output = m(img);

#### טיפול בשגיאות

.exceptions בתרגיל זה נדרוש מכם להשתמש

- חשבו בעצמכם היכן יכולות להיות שגיאות (בדיקת הקלטים של הפונקציות, בדיקת ערכי החזרה של פעולות שונות, ועוד).
  - במקרה של שגיאה:
  - זרקו חריגה(exception) מתאימה בהתאם להוראות הבאות:
  - .std::length error אם התבצעה שגיאה הנוגעת לאורכים בעייתיים,זרקו
  - .std::out of range אם התבצעה שגיאה שקשורה לגישה למיקום לא חוקי, זרקו
    - .std::runtime error אם התבצעה שגיאה עקב קלט לא נכון מהמשתמש \*
- אניזרק באופן אוטומטי במקרה זה.  $std::bad\_alloc$  חריגה(החריגה לזרוק לזרוק אינכם בדרשים אינכרון נכשלה אינכם אינכם לזרוק אינכם לזרוק אינכם לזרוק (https://en.cppreference.com/w/cpp/memory/new/bad alloc לפירוט,קראו:
  - בתרגיל זה אינכם נדרשים לשחרר את הזכרון במקרה של שגיאה(כיוון שזאת אחריות משתמש הספרייה לשחרר את הזיכרון).

## 5 קימפול והרצה

בקבצי העזר לתרגיל הניתנים לכם, מצורף קובץ Makefile על מנת לקמפל את התוכנה. על התוכנה להתקמפל באמצעות הפקודה הבאה: make mlpnetwork

נריץ את התוכנית כמפורט בסעיף 2.3.

## Presubmit 5.1

קובץ ה presubmit זמין בנתיב

~labcc2/presubmit/ex4/srcs/presubmit.cpp

## 6 הקבצים להגשה

Matrix.h Matrix.cpp
Activation.h Activation.cpp
Dense.h Dense.cpp
MlpNetwork.h MlpNetwork.cpp

## 7 הערות וסיכום

#### 7.1 הנחיות כלליות

- קראו בקפידה את הוראות תרגיל זה ואת ההנחיות להגשת תרגילים שבאתר הקורס.
- קריות malloc: זכרו להשתמש ב-malloc מעבר לי-C++ זכרו להשתמש בפונקציות ובספריות של לי-C++ ולא למשל, נשתמש ב-math.h ולא ב-C במקום std::std::string
- איסור על שימוש ב־ $\mathrm{STL}$ : נזכיר בשנית כי בתרגיל זה נאסר השימוש בספריית איסור על שימוש ב־ $\mathrm{STL}$ : נזכיר בשנית כי בתרגיל זה נאסר השימוש בספריית std::vector
- ניהול זיכרון דינמי: כזכור, הקצאת זיכרון דינמית מחייבת את שחרור הזיכרון, למעט במקרים בהם ישנה שגיאה המחייבת סגירת התוכנית באופן מיידי עם קוד שגיאה (כלומר קוד יציאה 1). תוכלו להיעזר בתוכנה valgrind כדי לחפש דליפות זיכרון בתוכנית שכתבתם.
  - by reference הקפידו לא להעתיק by value משתנים כבדים, אלא להעבירם היכן שניתן reference שימוש ב-reference הקפידו לא

- שימוש ב־const: הקפידו להשתמש במילה השמורה const היכן שנדרש מכם בהגדרת הפונקציות והארגומנטים.
- ידאו כי התרגיל שלכם עובר את ה־Pre-submission Script ודאו כי התרגיל שלכם עובר את יציאות או אזהרות.
  - בונוס:
  - .22:00 בשעה 18/5 בשעה עד יום רביעי בונוס בינוס -
  - בונוס מראש(כרגיל). +1/+2/+3 הגשה הום/יומיים שלושה ימים מראש
- שימו לב: הפעם אתם לא צריכים להגיש טסטים. חרף זאת, אנו מעודדים אתכם לכתוב טסטים ולבדוק את עצמכם במהלך פתירת התרגיל.

## בהצלחה!