

РЕАЛІЗАЦІЯ СИГМОЇДАЛЬНИХ ФУНКЦІЙ АКТИВАЦІЇ НА ПЛІС ДЛЯ НЕЙРОННИХ МЕРЕЖ

І. Цмоць¹, В. Рабик², І. Ігнатєв³

¹Національний університет “Львівська Політехніка”,
вул. Бандери, 12, 79013 Львів, Україна
ivan.tsmots@gmail.com

²Львівський національний університет імені Івана Франка,
вул. Ген. Тарнавського, 107, 79017 Львів, Україна
rabykv@ukr.net

³Тернопільський національний економічний університет,
вул. Львівська, 11, 46009, Тернопіль, Україна
ignatyevkii@gmail.com

Розроблено та модифіковано методи апроксимації функції активації нейронних мереж (сигмоїдальна функція, тангенс гіперболічний), використовуючи кусково-лінійну апроксимацію і апроксимацію поліномом другого порядку. Виконано оцінку точності апроксимації цими методами сигмоїдальної функції та її похідної. Для розглянутих методів апроксимації розроблено структурні схеми, які реалізовані мовою VHDL з використанням елементів бібліотек Quartus II на FPGA. Наведено апаратні ресурси, необхідні для реалізації.

Ключові слова: FPGA, VHDL, функція активації, сигмоїдальна функція, кусково-лінійна апроксимація, середня абсолютна похибка, максимальна абсолютна похибка, апаратні ресурси FPGA.

Вступ. В останні десятиліття зростає зацікавленість до апаратної реалізації штучних нейронних мереж (ШНМ). Про це свідчить кількість публікацій на цю тематику. Це пов'язано насамперед зі стрімким розвитком елементної бази, яка використовується при реалізації цифрових ШНМ (надвеликі інтегральні схеми - VLIS). Однією з актуальних проблем, які залишаються при цьому, є збільшення швидкості роботи нейронних мереж. Одним зі способів, що дозволяє прискорити їх роботу за рахунок паралелізму, є їх реалізація на FPGA.

Швидкість роботи штучних нейронів залежить від функції активації. Реалізація сигмоїдальної функції, тангенса гіперболічного на FPGA вимагає значних апаратних ресурсів [1, 2, 3]. В роботі [1] приведено огляд основних методів, які використовуються при реалізації сигмоїдальної функції та тангенса гіперболічного.

Для цифрової реалізації нелінійних функцій активації використовують різні методи апроксимації: табличний, розклад в ряд Тейлора, кусково-лінійну апроксимацію [1, 2]. При розкладі в ряд Тейлора потрібно виконувати багато перемножень, тому цей метод не є прийнятним для реалізації на FPGA. Табличний метод передбачає створення

таблиці можливих значень цільової функції з врахуванням обмеженої її розрядності. Але створення окремої локальної таблиці для кожного нейрона вимагатиме значних апаратних ресурсів FPGA. Використання ж однієї таблиці для різних нейронів призведе до великих часових затримок, так як поширення сигналів в нейронах одного шару виконується паралельно. Тому найбільш часто при реалізації функції активації сигмоїдального типу використовується кусково - лінійна апроксимація (PWL). Кусково-лінійна апроксимація сигмоїдальної функції полягає в заміні нелінійної функції на кожному з вибраних її інтервалів прямою лінією. Така апроксимація використовується в багатьох роботах [2, 3].

Вибір методу апроксимації функцій активації сигмоїдального типу і їх апаратна реалізація є головними аспектами від яких залежить точність і швидкодія алгоритму. При низькій точності апроксимації отримаємо малу швидкодію, а зменшення похибки апроксимації призводить до збільшення апаратних ресурсів і зменшення швидкості оброблення даних.

Іншим важливим моментом, який потрібно враховувати, є диференційованість апроксимованої функції активації [2, 3], так як методи навчання ШНМ включають як функцію активації, так і її похідну.

Функції активації нейронів сигмоїдального типу.

Сигмоїдальні функції активації - найчастіше використовується для нейронних мереж прямого поширення сигналу. Вони є монотонно зростаючими, неперервними і диференційованими. Для опису і аналізу функцій активації сигмоїдального типу в роботі [1] використано загальний клас функцій:

$$f(x, k, b, T, c) = k + \frac{c}{1 + be^{Tx}}, \quad \forall x \in R, \quad (1)$$

де $k \in R$, $b \in R^+$, $T, c \in R \setminus \{0\}$, R - множина дійсних чисел $(-\infty, +\infty)$, R^+ - множина дійсних додатних чисел $(0, +\infty)$, $R \setminus \{0\}$ - множина дійсних чисел за винятком точки 0 $(-\infty, 0)$ та $(0, +\infty)$. Сигмоїдальні нелінійності, які відносяться до цього класу:

- класична сигмоїдальна функція ($k=0$, $c=b=1$ та $T=-1$);

$$f(x) = \frac{1}{1 + e^{-x}}, \quad (2)$$

- тангенс гіперболічний ($k=1$, $c=-2$, $b=1$, $T=2$):

$$h(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}. \quad (3)$$

Необхідно відмітити, що обчислення сигмоїдальної функції активації (2) достатньо виконувати тільки для додатних аргументів x . Для від'ємних значень x її можна знайти з виразу:

$$f(-x) = 1 - f(x). \quad (4)$$

Дійсно, виконавши прості перетворення для виразу (2), отримаємо:

$$f(-x) = \frac{1}{1 + e^x} = \frac{1}{1 + \frac{1}{e^{-x}}} = \frac{e^{-x}}{1 + e^{-x}} = \frac{1 + e^{-x} - 1}{1 + e^{-x}} = 1 - f(x),$$

Нелінійну функцію тангенса гіперболічного (3) також можна обчислити через класичну сигмоїдальну функцію (2).

$$h(x) = 1 + \frac{e^{-x}}{e+x} - 1 = \frac{2e^{-x}}{e+x} - 1 = \frac{2}{1+e^x} - 1 = 2f(2x) - 1, \quad (5)$$

$$h(-x) = 2f(-2x) - 1 = 2[1 - f(2x)] - 1 = 1 - 2f(2x). \quad (6)$$

Для оцінки точності апроксимації використовується максимальна і середня похибки [2]. Середня абсолютна ε_{ave} , максимальна абсолютна ε_{max} похибки функції $f(x)$, яка апроксимується функцією $\bar{f}(x)$ в інтервалі (x_{min}, x_{max}) визначаються як:

$$\varepsilon_{ave} = \frac{\sum_{i=0}^{N_p-1} |\bar{f}(x_i) - f(x_i)|}{N_p}, \quad (7)$$

$$\varepsilon_{max} = \max |\bar{f}(x_i) - f(x_i)|, \quad i = 0, \dots, N_p, \quad (8)$$

де N_p - кількість точок, на які розбивається інтервал (x_{min}, x_{max}) .

Ефективність апроксимації порівнюється по досягнутій точності, швидкості та апаратним ресурсам.

Арифметичні операції в нейронах виконуються над дійсними числами. При реалізації арифметичних операцій на FPGA над дійсними числами з фіксованою та плаваючою крапками зростає час їх виконання та апаратні ресурси необхідні для їх реалізації. Тому при реалізації компонент штучних нейронів на мові VHDL дійсні числа перетворювалися до цілих шляхом їх домноження на 2^{10} і відсіканням дробової частини. Формат представлення дійсних чисел має розмірність 16 біт: 1 біт – знаковий і 15 біт – для зберігання отриманого цілого числа. Від’ємні числа представляються в доповнюючому коді. Максимальна ціла частина дробового числа не повинна перевищувати 15. Необхідно зауважити, що вхідні дані в нейронних мережах нормуються. Для представлення дійсного додатного числа 2.1625 у форматі цілих чисел домножуємо його на $2^{10} = 1024$ і відсікаємо дробову частину: $2.1625 * 1024 = 2214.4 \approx 2214$. В шістнадцятковій системі числення це число рівне: $2214_{10} = 0x08A6$. Це ж саме від’ємне число в шістнадцятковій системі числення $-2214_{10} = 0xF75A$.

Методи апроксимації сигмоїдальних функцій.

Розглянемо декілька алгоритмів кусково-лінійної апроксимації, які відрізняються кількістю точок, розміщенням початкової і кінцевої точок на апроксимуючих лініях та критеріями їх вибору.

1. Кусково-лінійна апроксимація нелінійної функції [2, 3]. В цьому методі апроксимація сигмоїдальної функції (2) виконується виразом [2]:

$$f(x) = \begin{cases} 1, & x \geq 5.0 \\ 0.03125 * x + 0.84375, & 2.375 \leq x < 5.0 \\ 0.125 * x + 0.625, & 1.0 \leq x < 2.375 \\ 0.25 * x + 0.5, & 0 \leq x < 1.0 \end{cases} \quad (9)$$

Обчислення повинні виконуватися тільки для додатних значень x . Для від'ємних вхідних даних x сигмоїдальна функція обчислюється з допомогою виразу (4). Перетворимо вираз (9) в цілочисельний, помноживши його доданки без змінної x на 2^{10} :

$$f(x) = \begin{cases} 1024, & x \geq 5120 \\ 2^{-5} * x + 864, & 2432 \leq x < 5120 \\ 2^{-3} * x + 640, & 1024 \leq x < 2432 \\ 2^{-2} * x + 512, & 0 \leq x < 1024 \end{cases} \quad (10)$$

Для реалізації сигмоїдальної функції відповідно до виразу (10) розроблена структурна схема пристрою, яка наведена на рис. 1, де Rg – регістр, LessThan – пристрої порівняння, Add – суматори, Buf – буферний пристрій з третім станом.

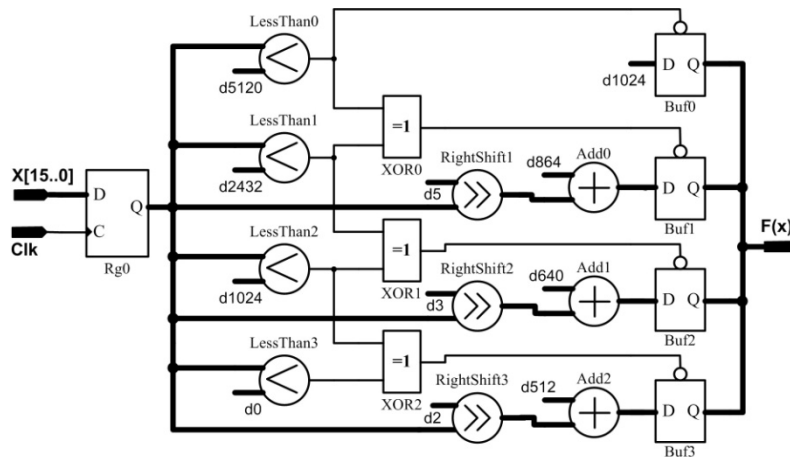


Рис. 1. Структурна схема пристрою, що реалізує сигмоїдальну функцію згідно виразу (10)

Для обчислення сигмоїдальної функції в цьому пристрої необхідно три суматори Add, чотири пристрої порівняння LessThan, чотири шинних формувачі Buf. Множення на величини 2^{-2} , 2^{-3} і 2^{-5} (зсув вправо) здійснюється шляхом відповідного з'єднання входів та виходів відповідних компонентів схеми. На схемі це позначено з допомогою пристроїв зсуву RightShift. Вибір результатів обчислення сигмоїдальної функції, які отримуються на виходах Add1, Add2 і Add3, здійснюється за результатами порівняння вхідних даних x з числами 5120, 2432, 1024 та 0. На одні входи (A_i) пристроїв порівняння LessThan поступають вхідні дані (x), які порівнюються з даними на інших входах (B_i). На виходах LessThan формується результат порівняння:

$$Out_{LessThan} = \begin{cases} 0, & x < B_i \\ 1, & x \geq B_i \end{cases} \quad (11)$$

Інформація з виходів пристроїв порівняння LessThan надходить на входи логічних елементів виключне АБО (XOR). Сигнали з LessThan1 та виходів логічних елементів XOR поступають на входи шинних формувачів, які передають дані на вихід, якщо на вході дозволу сигнал логічної '1' та переводять вихід в третій стан (Z), якщо на цьому вході – сигнал логічного '0'.

Час обчислення сигмоїдальної функції для цієї схеми визначається за формулою:

$$t_1 = t_{Rg} + t_{LessThen} + t_{XOR} + t_{Add} + t_{Buf}, \quad (12)$$

де t_{Rg} , $t_{LessThen}$, t_{XOR} , t_{Add} , t_{Buf} – час затримки проходження сигналу відповідно через регістр, пристрої порівняння, логічні елементи XOR, суматори та шинні формувачі.

Моделювання сигмоїдальної функції (2) та її кусково-лінійної апроксимації (9) зображено на рис. 2а, а їх похідних - на рис. 2б.

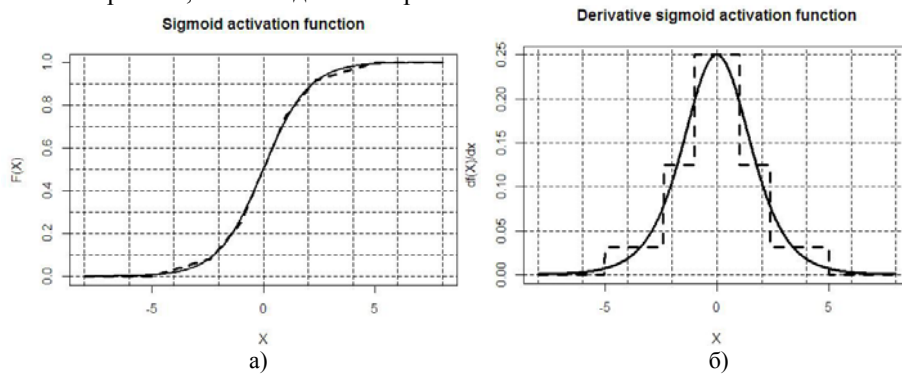


Рис. 2. а) - вигляд сигмоїдальної функції та її кусково-лінійної апроксимації; б) – вигляд похідних сигмоїдальної функції та її кусково-лінійної апроксимації

Розглянемо діапазон зміни вхідного сигналу (-8, 8) і розіб'ємо його на $N_p=1000$ інтервалів. В цьому діапазоні середня і максимальна похибки кусково-лінійної апроксимації сигмоїдальної функції рівні $\epsilon_{ave} = 0.00587$, $\epsilon_{max} = 0.0185$, та її похідної $d\epsilon_{ave} = 0.01412$, $d\epsilon_{max} = 0.07088$.

Вигляд абсолютної похибки між сигмоїдальною функцією та її кусково-лінійною апроксимацією (суцільна крива), між похідними сигмоїдальної функції і її кусково-лінійною апроксимацією (пунктирна крива) зображений на рис. 3.

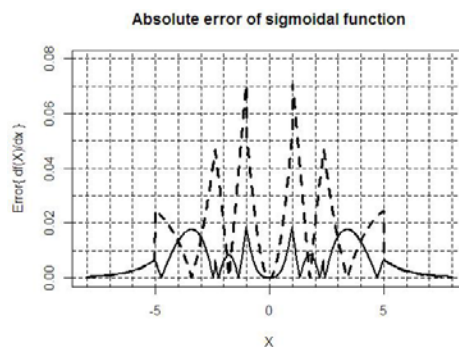


Рис. 3. Вигляд абсолютної похибки між сигмоїдальною функцією і її кусково-лінійною апроксимацією та між їх похідними

2. Апроксимація сигмоїдальної функції (2) кривою другого порядку. Для вхідних даних з діапазону $(0, x_{max})$ сигмоїдальна функція апроксимується поліномом [2, 3]:

$$\bar{f}(x) = c + bx + ax^2. \quad (13)$$

В рівнянні (13) невідомими є коефіцієнти a , b , c , які визначаються методом найменших квадратів. Система рівнянь для їх обчислення має вигляд:

$$\begin{bmatrix} \sum_{i=0}^{N_p-1} x_i^2 & \sum_{i=0}^{N_p-1} x_i^3 & \sum_{i=0}^{N_p-1} x_i^4 \\ \sum_{i=0}^{N_p-1} x_i & \sum_{i=0}^{N_p-1} x_i^2 & \sum_{i=0}^{N_p-1} x_i^3 \\ N_p & \sum_{i=0}^{N_p-1} x_i & \sum_{i=0}^{N_p-1} x_i^2 \end{bmatrix} * \begin{bmatrix} c \\ b \\ a \end{bmatrix} = \begin{bmatrix} \sum_{i=0}^{N_p-1} f(x_i) x_i^2 \\ \sum_{i=0}^{N_p-1} f(x_i) x_i \\ \sum_{i=0}^{N_p-1} f(x_i) \end{bmatrix}. \quad (14)$$

Розбиваємо діапазон $(0, 4)$ на $N_p=1000$ інтервалів. Використовуючи значення x_i та $f(x_i)$ в цих точках, формуємо систему рівнянь (14). Розв'язавши її отримаємо коефіцієнти a , b , c і вираз для апроксимації сигмоїдальної функції:

$$f(x) = \begin{cases} 1, & x \geq 4.0 \\ -0.03577 * x^2 + 0.25908 * x + 0.5038, & 0 \leq x < 4.0 \end{cases}. \quad (15)$$

Цілочисельний вираз (15) має вигляд:

$$f(x) = \begin{cases} 1024, & x \geq 4096 \\ -36 * 2^{-20} * x^2 + 265 * 2^{-10} * x + 515, & 0 \leq x < 4096 \end{cases}. \quad (16)$$

В діапазоні $(-8, 8)$ середня і максимальна похибки апроксимації сигмоїдальної функції поліномом другого порядку рівні $\varepsilon_{ave} = 0.00426$, $\varepsilon_{max} = 0.01798$ та її похідної - $d\varepsilon_{ave} = 0.00769$, $d\varepsilon_{max} = 0.04388$. Сигмоїдальна функція та апроксимація її поліномом другого порядку зображено на рис. 4а, а їх похідні - на рис. 4б. Суцільною лінією на рис. 4 зображена сигмоїдальна функція і її похідна, а пунктирною лінією – їх апроксимації.

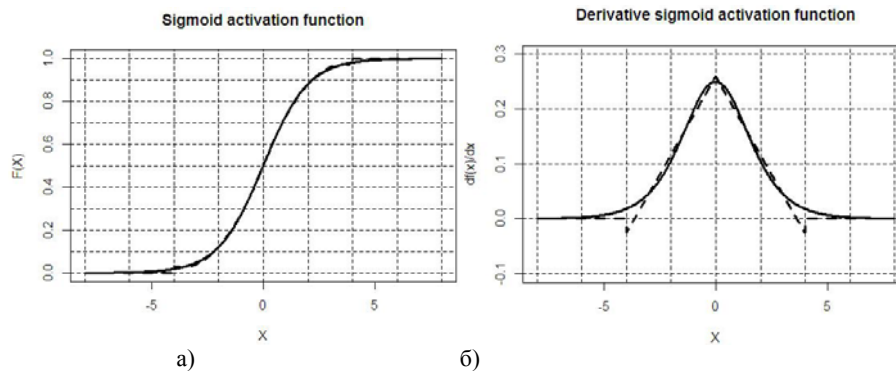


Рис. 4. а) - вигляд сигмоїдальної функції та її апроксимації поліномом другого порядку;
б) – вигляд похідних сигмоїдальної функції та її апроксимації поліномом другого порядку

На рис. 5 зображено абсолютні похибки між сигмоїдальною функцією та її апроксимацією поліномом другого порядку (суцільна лінія) і між похідною сигмоїдальної функції та похідною її апроксимації (пунктирна лінія).

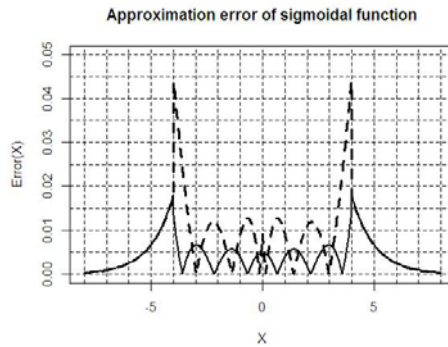


Рис. 5. Вигляд абсолютної похибки між сигмоїдальною функцією і її апроксимацією поліномом другого порядку та між їх похідними

Структурна схема пристрою, що реалізує сигмоїдальну функцію (16), наведена на рис. 6.

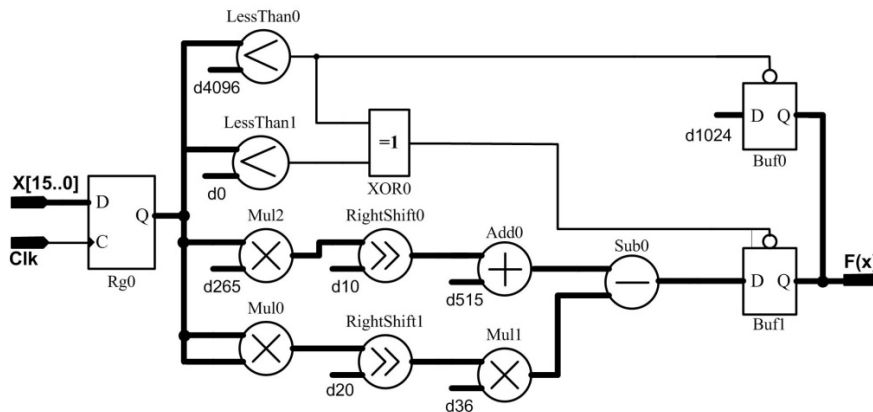


Рис. 6. Структурна схема пристрою, що реалізує сигмоїдальну функцію згідно виразу (16)

У цьому пристрої для реалізації сигмоїдальної функції (2) використовуються три перемножувачі Mul, суматор Add, пристрій віднімання Sub, дві схеми порівняння LessThan, два шинних формувачі Buf. Максимальний час обчислення сигмоїдальної функції в цьому пристрої визначається за формулою:

$$t_2 = t_{Rg} + 2t_{Mul} + t_{Sub} + t_{Buf}, \tag{17}$$

де t_{Rg} , t_{Mul} , t_{ADD} , t_{Buf} – час затримки проходження сигналу відповідно через регістр, перемножувач, пристрій віднімання та шинний формувач.

3. В роботах [2, 3] для апроксимації сигмоїдальної функції поліномом другого порядку використовується спрощений вираз:

$$f(x) = \begin{cases} 1, & x \geq 4.0 \\ -0.03125 * x^2 + 0.25 * x + 0.5, & 0 \leq x < 4.0 \end{cases} \quad (18)$$

реалізація якого вимагає тільки одного перемножувача та суматорів.

В цілочисельному форматі вираз (18) має вигляд:

$$f(x) = \begin{cases} 1024, & x \geq 4096 \\ -2^{-15} * x^2 + 2^{-2} * x + 512, & 0 \leq x < 4096 \end{cases} \quad (19)$$

Сигмоїдальна функція і апроксимація її виразом (18) в діапазоні $(-8, 8)$ зображені на рис. 7а, а їх похідні – на рис. 7б. На цих рисунках суцільною лінією зображена сигмоїдальна функція та її похідна, а пунктирною – апроксимація функції виразом (18) та її похідна.

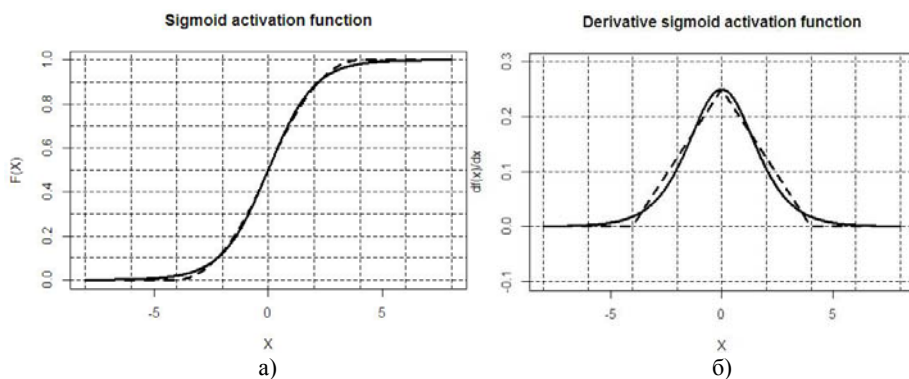


Рис. 7. Графіки а) – сигмоїдальної функції та її апроксимації; б) – похідних сигмоїдальної функції та її апроксимації

Середня і максимальна похибки апроксимації сигмоїдальної функції виразом (18) в діапазоні $(-8, 8)$ рівні $\varepsilon_{ave} = 0.00774$, $\varepsilon_{max} = 0.02160$, та їх похідних $d\varepsilon_{ave} = 0.00877$, $d\varepsilon_{max} = 0.02375$. Вигляд абсолютних похибок зображений на рис. 8 (суцільна лінія – похибка між сигмоїдальною функцією та апроксимуючим виразом (18), пунктирна лінія – похибка між їх похідними).

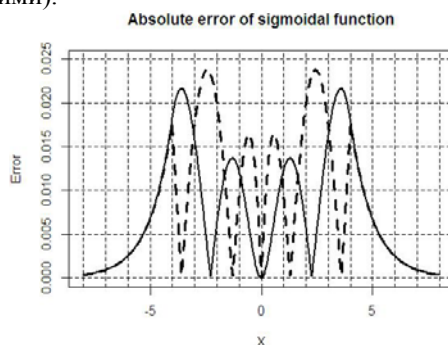


Рис. 8. Вигляд абсолютної похибки між сигмоїдальною функцією і її апроксимацією виразом (18) та між їх похідними

Для реалізації сигмоїдальної функції (19) розроблена структурна схема пристрою, яка наведена на рис. 9.

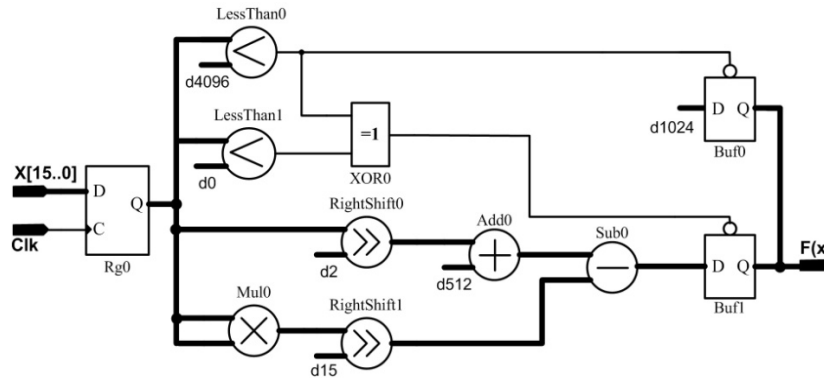


Рис. 9. Структурна схема пристрою, що реалізує сигмоїдальну функцію згідно виразу (19)

У даному пристрої для реалізації сигмоїдальної функції використовуються перемножувач, пристрій віднімання та суматор. У порівнянні з пристроєм на рис. 6, цей пристрій вимагає для своєї реалізації менших апаратних ресурсів. Максимальний час обчислення сигмоїдальної функції в даному пристрої визначається за формулою:

$$t_3 = t_{Rg} + t_{Mul} + t_{Add} + t_{Sub} + t_{Buf} \quad (20)$$

Реалізація сигмоїдальних функцій на FPGA.

Реалізація сигмоїдальних функцій виконувалася в середовищі розробки Quartus II для FPGA EP3C16F484C6 сімейства Cyclone III з використанням мови програмування апаратури VHDL та елементів бібліотек Quartus II.

Для цього використовувалися розглянуті методи апроксимації сигмоїдальної функції (вирази (10), (16), (19)) та структурні схеми (рис. 1, рис. 6, рис. 9), що реалізують ці вирази. Кожний з цих виразів дозволяє обчислити значення функції активації для додатних значень зваженої суми, що подається на її вхід. Для від'ємних значень суми при обчисленні функції активації використовується вираз (4).

Розглянемо реалізацію структурної схеми, наведеної на рис. 9. На рис. 10 зображено символ (FA_Sigm_N3) – зовнішній вигляд пристрою для апроксимації сигмоїдальної функції виразом (19). Входами пристрою є: Clk – вхід синхронізації, IN[15..0] – сума зважених входів нейрона розрядністю 16, а виходом – Out[15..0] – значення функції активації розрядністю 16. Обчислення сигмоїдальної функції виконується по передньому фронту імпульсів, які поступають на вхід Clk.

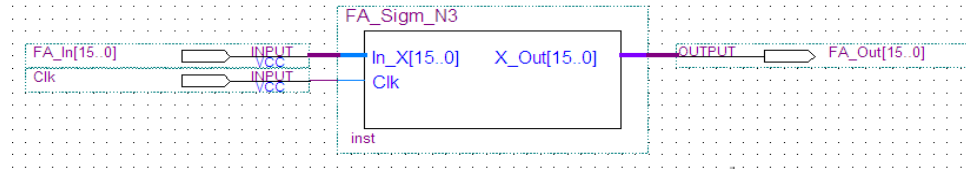


Рис. 10. Зовнішній вигляд символу FA_Sigm_N3

На рис. 11 наведений фрагмент часової діаграми роботи цього пристрою, який обчислює сигмоїдальну функцію. Часова діаграма отримана, використовуючи моделювання в часовій області в середовищі розробки Quartus II.

На входи модуля FA_Sigm_N3 подається послідовність імпульсів синхронізації з періодом 20 нс та послідовність входніх сигналів від -4192 до 4192 з кроком 32. На його виході отримуємо величину апроксимованої сигмоїдальної функції. На рис. 11а зображений фрагмент часової діаграми з входніми величинами в діапазоні від 544 до 1088, а на рис. 11б – в діапазоні від -1056 до -544. Так значенню входнього сигналу $FA_In=640$ відповідає величина сигналу $FA_Out=512+2^{-3}*640-2^{-15}*640*640=660$, що співпадає зі значенням на часовій діаграмі на рис. 11а. Для від'ємного значення $FA_In=-640$ згідно рис. 11б – $FA_Out=364$ ($1024-660=364$). Затримка вихідного сигналу по відношенню до імпульсів синхронізації (t_{CO}) для FA_Sigm_N3 складає 17...18 нс.

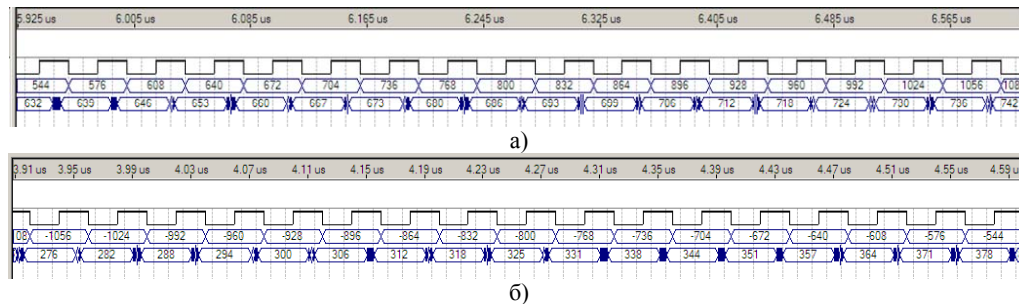


Рис. 11. Часові діаграми роботи модуля сигмоїдальної функції FA_Sigm_N3

Схема символу FA_Sigm_N3 зображена на рис. 12. На вхід регістру reg поступають імпульси синхронізації Clk та входні дані In_X, для яких обчислюється функція активації. Якщо входні дані від'ємні, то на виході символу Code, отримуємо їх доповнюючий код. Для додатних значень даних на вихід Code вони подаються без змін. Розпізнавання додатних і від'ємних входніх даних виконується з допомогою сигналу Sign_In. Якщо Sign_In='1', то входні дані додатні.

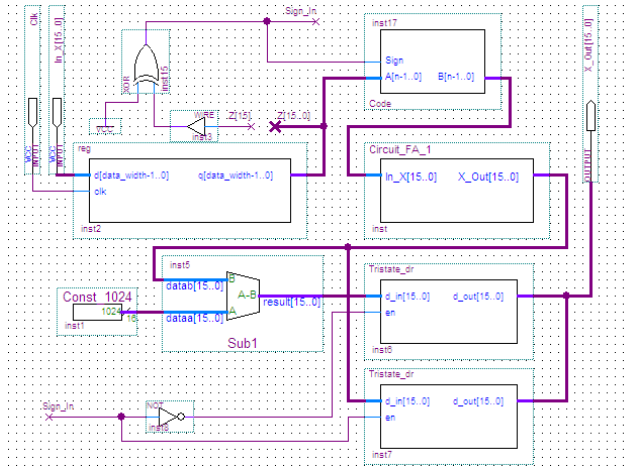


Рис. 12. Схема пристрою обчислення сигмоїдальної функції FA_Sigm_N3

З виходу символу Code вхідні дані поступають на пристрій обчислення сигмоїдальної функції (Circuit_FA_1) для додатних вхідних даних. З виходу Circuit_FA_1 дані поступають як на пристрій віднімання (Sub1), так і на буферний пристрій з третім станом (Tristate_dr). Для від’ємних вхідних значень (Sign_In=’0’) виконується віднімання від константи 1024 отриманої величини функції активації і отримана різниця поступає на інший буферний пристрій. На його виході отримуємо величину апроксимованої сигмоїдальної функції.

Апаратні ресурси, необхідні для реалізації сигмоїдальних функцій, порівнюються по кількості логічних елементів (ЛЕ) та виводів FPGA EP3C16F484C6 сімейства Cyclone III фірми Altera. Максимальна кількість ЛЕ для цієї FPGA [4] складає 15408, а кількість виводів – 347. Для реалізації сигмоїдальної функції активації (FA_Sigm_N3) використано 135 ЛЕ та 33 виводи FPGA.

Порівняння по точності апроксимації сигмоїдальної функції (2) виразами (9), (15), (18) та її похідної приведено в табл. 1. Похибка апроксимації сигмоїдальної функції цими виразами є меншою, ніж похибка апроксимації її похідних.

Табл. 1.

Порівняння по точності апроксимацій сигмоїдальної функції та її похідних

Вираз	ε_{ave}	ε_{max}	$d\varepsilon_{ave}$	$d\varepsilon_{max}$
(9)	0.00587	0.01850	0.01412	0.07088
(15)	0.00426	0.01798	0.00769	0.04388
(18)	0.00774	0.02160	0.00877	0.02375

Сигмоїдальні функції реалізовані на стенді DE0 [5]. Значення зваженої суми, що подається на вхід пристрою FA_Sigm_N3 задається з допомогою перемикачів, а величина функції активації виводиться на 4-ри розрядний семисегментний індикатор стенду.

1. *Beiu V., Peperstraete J.A., Vandewalle J., Lauwereins R.* Close Approximations of Sigmoid Functions by Sum of Steps for VLSI Implementation of Neural. [Електронний ресурс]. Режим доступу: <https://pdfs.semanticscholar.org/fdef/62a66787929bb80163219744d9eab041b203.pdf>
2. *Tommiska M.T.* Efficient digital implementation of the sigmoid function for reprogrammable logic. [Електронний ресурс]. Режим доступу: <https://pdfs.semanticscholar.org/9bf6/4bae3f8528cd5ac72a4ae869a74563ff6c26.pdf>
3. *Tisan Alin, Oniga Stefan, Mic Daniel, Buchman Attila.* Digital Implementation of the Sigmoid Function for FPGA Circuits. [Електронний ресурс]. Режим доступу: http://users.utcluj.ro/~ATN/papers/ATN_2_2009_4.pdf
4. Cyclone III Device Handbook. [Електронний ресурс]. Режим доступу: http://www.altera.com/literature/hb/сyc3/сyc3_ciii51001.pdf
5. DE0 User Manual. [Електронний ресурс]. Режим доступу: http://esca.korea.ac.kr/teaching/FPGA_boards/DE0/DE0_User_Manual.pdf

IMPLEMENTATION OF SIGMOID ACTIVATION FUNCTIONS ON FPGA FOR NEURAL NETWORKS

I. Tsmots¹, V. Rabyk², I. Ignatyev³

¹*Lviv Polytechnic National University,
12 Bandera St., Lviv, Ukraine, 79013
ivan.tsmots@gmail.com*

²*Ivan Franko National University of Lviv,
107 Tarnavskogo St., Lviv, Ukraine, 79017
rabykv@ukr.net*

³*Ternopil National Economic University
11 Lvivska St., Ternopil, Ukraine, 46009
ignatyevkii@gmail.com*

The paper presents a literature review of the main methods of neural networks activation function approximation (sigmoidal function, hyperbolic tangent). The attention is focused on the accuracy and speed of the methods of approximation of the activation function and the advantages and disadvantages of these methods in their hardware implementation on the FPGA.

The mathematical descriptions of the sigmoidal function, the hyperbolic tangent and their derivatives are considered. The ratios for calculating sigmoidal functions with negative arguments are given. In order to evaluate the accuracy of the approximation of the activation functions and their derivatives, the mean absolute and mean maximum errors were used. The conversion of valid input data into an integer format is achieved by its multiplication by 210 and cutting off the fractional part. In the implementation, the following format of representation of real numbers is used: 16 bits, 1 sign bit and 15 bits for storing the resulting integer number. Negative numbers are presented in two's complement.

Three methods of approximation of sigmoidal function are considered. The first one is a piecewise linear approximation. A nonlinear function with positive arguments is replaced by four straight lines. The second method is the approximation by the second-order polynomial. A system of equations for calculating a , b and c coefficients for polynomials is obtained and values of these coefficients are calculated. The third method is the approximation by the second-order polynomial with a simplified expression. A comparison is made for the accuracy of approximations by three methods of sigmoidal function and its derivatives. It is noted that the errors of the approximation of the sigmoidal function are smaller than the errors of approximation of its derivatives. For the considered approximation methods, integer expressions were obtained, according to which the structural schemes of the implementation of the sigmoidal function on the FPGA were developed. The main elements of these schemes are registers, comparators, multipliers, adders, subtraction devices and three-state buffers.

Digital hardware implementation of sigmoid activation functions was performed in VHDL language in the Quartus II development environment using elements from standard library. Devices with the functionality of sigmoidal function element are realized on FPGA EP3C16F484C6 of Cyclone III family. The simulation of the device calculating the sigmoidal function FA_Sigm_N3 in the time domain for positive and negative input data was carried out. The diagram of the device and its description and hardware resources necessary for the realization and calculation time are presented. The hardware realization of the sigmoid activation function on the DE0 stand was performed and the testing of its work was carried out.

Key words: FPGA, VHDL, activation function, sigmoid function, piecewise linear approximation, mean absolute error, maximum absolute error, hardware resources FPGA hardware costs.

*Стаття: надійшла до редакції 10.03.2019,
доопрацьована 14.03.2019,
прийнята до друку 15.03.2019.*