

## Прийняття рішення при розпізнаванні мови на основі нейронних мереж

О.М. КАРПОВ, Г.В. ЗІРНЄЄВА

Дніпропетровський національний університет

В статті представлена архітектура нейронної мережі для прийняття рішення при розпізнаванні мови, алгоритм сегментації мовного сигналу, алгоритм навчання та проаналізована ефективність побудованої нейронної мережі.

В статье представлена архитектура нейронной сети для принятия решения при распознавании речи, алгоритм сегментации речевого сигнала, алгоритм обучения и проанализирована эффективность построенной нейронной сети.

In article the architecture of a neural network for making-decision is presented at recognition of speech, algorithm of segmentation of a speech signal, algorithm of training and efficiency of the constructed neural network is analyzed.

**Постановка проблеми.** Проблема розпізнавання мови є дуже актуальною в наш час і мало вирішеною. Особливо актуальною є проблема побудови системи штучного інтелекту за допомогою нейронних мереж. Однією з важливих наукових задач, є створення нейронної мережі, такої архітектури, щоб забезпечити найбільш високий відсоток розпізнавання мови.

**Аналіз останніх досліджень і публікацій.** Уперше теорія нейронних мереж як новий напрямок була позначена в роботі Маккаллока і Питтса. Великий внесок у розпізнаванні мовного сигналу за допомогою нейронних мереж внесли наступні учені: П.Е. Овчинников, Ю.А.Семин - навчання перцептрона без сегментації слів з навчальної вибірки в задачі розпізнавання звуків, В. Л. Арлазаров - використання острівного нейромережевого аналізу мовного сигналу в кореляції з виділенням стійких ознак і застосуванні фонологічних і інших "інженерних" знань, О.И.Федяев розробив нейромережевий метод аналізу фонетичної структури мовного сигналу. А також розпізнавання мовних сигналів за допомогою нейронних мереж займалися наступні учені Д.Разумихин, О.М. Карпов, П.А. Лалетин та ін.

Проблема розпізнавання мови як одне з складових штучного інтелекту давно залучала дослідників, і на сьогоднішній день хоч і досягнуті визначені успіхи, вона залишається відкритою. Не цілком вирішеними задачами в даній області, є задача сегментації мовного сигналу і безпосередньо розпізнавання мовного сигналу. У розпізнаванні мови за допомогою нейронних мереж запропоновані архітектури не дають високого відсотка розпізнавання вхідних сигналів.

**Постановка задачі.** Побудувати нейронну мережу для задачі ухвалення рішення при розпізнаванні мови. Побудувати таку архітектуру зв'язків синоптичних ваг, та такий спосіб їх з'єднання, який забезпечував найкраще розпізнавання мовного сигналу. Створити програмну реалізацію нейронної мережі, алгоритму навчання і безпосередньо розпізнавання мови. Використовувати мережну архітектуру для тимчасової обробки. Проаналізувати ефективність створеної нейронної мережі. Дані про навчений нейронної мережі зберігати у форматі xml.

**Основний матеріал.** Для реалізації системи прийняття рішень при розпізнаванні мови на основі нейронної мережі необхідними етапами є обробка вхідного сигналу.

На першому етапі сигнал, що надійшов, нормалізуємо, відтиснемо латентні періоди, переводимо в спектральне представлення за допомогою швидкого перетворення Фур'є, далі перетворимо в спектрально-смугове представлення і забираємо шуми за допомогою фільтрації ФНЧ.

На другому етапі проводимо розбивку вхідного сигналу на фонемі за допомогою методу верифікації по сукупності параметрів [1]. Метод верифікації по сукупності параметрів полягає в наступному:

Припустимо, кожна еталонна послідовність  $Y_k$  вміщує  $[y_{k0}(\omega), \dots, y_{ki}(\omega), \dots, y_{km}(\omega)]$  елементів, об'єднаних в  $\mu_k$  груп послідовності, кожна група в еталонній послідовності містить  $\gamma_j$  фонем,  $j=1-\mu_k$ , вхідна послідовність  $X$  складається з  $x_0(\omega), \dots, x_i(\omega), \dots, x_r(\omega)$  елементів, об'єднаних в  $\theta$  груп послідовності, кожна група вхідної послідовності містить  $\sigma_p$  фонем,  $p=1-\theta$ . Задача сегментації визначити межі між групами чи фонемами. Сегментація мови верифікацією на приналежність елементів  $[y_{k0}(\omega), \dots, y_{ki}(\omega), \dots, y_{km}(\omega)]$ ;  $[x_0(\omega), \dots, x_i(\omega), \dots, x_r(\omega)]$  до деякого класу с близькими параметрами відбувається по правилу

$$d = \text{extrem} \{x_u\} \# \{x_i\}. \quad (1)$$

Межа визначається як

$$d_i = \max \{x_u\} \# \{x_{u+\Delta}\}, \quad (2)$$

$$u = v + \delta; v = 0, \Delta, 2\Delta, \dots, r - 2\Delta; \delta = 1 + \Delta$$

при умові

$$T_{\text{seg}} > 0,2 \text{сек}, \quad (3)$$

де  $\Delta$  – шаг верифікації( в даній роботі  $\Delta=3;2;2$ ); # – операція зіставлення (бітового в хемінгеном або десяткового в Евклідовому просторах).

$$d_{1u} = \sum_{v=0}^{r-2\Delta} \sum_{u=v}^{v+\Delta} \sum_{j=1}^n [x_u(\omega_j) - x_{u+\Delta}(\omega_j)]^2$$

В евклідовому просторі

$$d_{2u} = \sum_{v=0}^{r-2\Delta} \sum_{u=v}^{v+\Delta} \sum_{\delta=0}^{\Delta-1} \sum_{j=1}^n [x_u(\omega_j) - x_{u+\Delta+\delta}(\omega_j)]^2$$

Если  $y_{ki}(\omega)$  и  $x_i(\omega)$  – битовые последовательности по  $\omega$ , то хеммінгово расстояние

$$d_{1u}^h = \sum_{v=0}^{r-2\Delta} \sum_{u=v}^{v+\Delta} \sum_{j=1}^n \text{unc} [x_u(\omega_j) \wedge x_{u+\Delta}(\omega_j)]^2$$

$$d_{2u}^h = \sum_{v=0}^{r-2\Delta} \sum_{u=v}^{v+\Delta} \sum_{\delta=0}^{\Delta-1} \sum_{j=1}^n \text{unc} [x_u(\omega_j) \wedge x_{u+\Delta+\delta}(\omega_j)]^2$$

де  $\text{unc}()$  – функція несовпадений бит операції “исключающее ИЛИ – ^”, а  $x_u(\omega_j)$  – это  $S(\Omega_r, t_j)$ .

Верифікацію і, відповідно, сегментацію можна здійснювати на всій довжині  $T_{\text{фр}}$  слова чи фрази. На рис.1. представлений результат сегментації методом верифікації по сукупності параметрів.

Далі необхідно спроектувати архітектуру нейронної мережі. Задача застосування нейронних мереж для ухвалення рішення при розпізнаванні слів мови має свої особливості в порівнянні з іншими методами останнього кроку, а саме, Марківськими процесами, динамічного програмування, принципу максимуму Понтрягіна, локальних екстремумів, для яких вибір розпізнаного образу вважається як найкраще наближення пред'явленої й еталонної реалізації. Звичайно в цьому випадку чи обчислюється мінімальна відстань чи максимальна подоба.

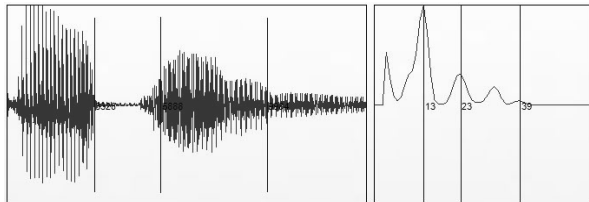


Рис. 1. Сегментація методом верифікації по сукупності параметрів. Вхідний сигнал - слово «один».

Для нейронних мереж у принципі адаптація при навчанні теж реалізується як найкраще наближення деякого опису, що відповідає одиничному стану деякого нейрона, призначеному вчителем для даного образу.

Однак у випадку мовного сигналу  $s(t)$ , що є функцією часу, параметричний опис також розподілений у часі. Нехай, для сигналу  $s(t)$  знайдено спектрально-часове представлення  $S(\omega_k, T_l)$ , де  $\omega_k$  - частота  $k$ -ї гармоніки перетворення Фур'є,  $T_l$  -  $l$ -й інтервал аналізу при  $T_l = 12,5 \text{ Мсек}$ ,  $k = \overline{1, 128}$ ,  $l = \overline{1, L}$ ,  $L = T_0 / T_l$ ,  $T_0$  - тривалість мовного висловлення.

В [2] розглядається деяка просторово-тимчасова модель нейрона, заснована на фокусованому нейронному фільтрі і FIR-фільтрі (фільтр із кінцевою імпульсною характеристикою).

Вихід  $v_j(n)$  має вид

$$v_j(n) = \sum_{i=1}^{m_0} \sum_{l=0}^p w_{ji}(l)x_i(n-l) + b_j,$$

де  $w_{ji}$  – вага  $l$ -го вторинного синапса, що належить  $i$ -му первинному синапсу;  $x_i(n)$  – вхідний сигнал, застосований до  $i$ -му первинного синапсу в момент часу  $n$ ;  $b_j$  – зсув, застосований до нейрона,  $v_j(n)$  – індуковане локальне поле чи нейрона аргумент функції активації  $\varphi$ .

Вихідна функція нейрона має вид  $y_j(n) = \varphi(v_j(n)) = 1/(1 + \exp(-av_j(n)))$  для сигмоїдальної функції активації. З погляду розпізнавання мовних сигналів необхідно аналізувати функцію  $S(\omega_k, T_l)$  розподілену в часі. В часі функція представлена деякою

сукупністю сегментів (фонем), кожен сегмент – сукупністю відліків – спектральних зрізів. Кожен спектральний зріз – вектор спектральних параметрів (128 чи 110 коефіцієнтів) чи спектрально-смугових (9 коефіцієнтів).

У даній роботі будується нейронна мережа для спектрально-полосових параметрів. Схема підключення синапсів може бути прямої й інверсної. Будемо вважати схему з'єднання нейронів у FIR-фільтрі в роботі [2] прямою. Інверсна схема припускає, що кожен спектральний зріз є варіантом образу і нейрон повинний навчатися на сукупність спектральних зрізів, що на інтервалі сегмента утворюють образ. Нейронна мережа, таким чином, повинна бути багатошаровою: Перший шар настроювання на спектральний зріз з вихідним значенням  $v_{1j}(n)$ , де  $n$  – номер спектрального зрізу,  $j$  номер сегмента (фонем) у слові; другий шар – розподілена в часі  $n$  послідовність  $v_{1j}(n)$  на сегменті  $j$ . Вихід другого шару  $-v_{2j}(j^*n, j^*n - p_j)$ , де  $(j^*n, j^*n - p_j)$  відповідає перебору значень  $v_{1j}(n)$  від  $(j^*n - p_j)$  до  $(j^*n)$  для кожного сегмента  $j = \overline{1, m}$ ,  $m$  – кількість сегментів.

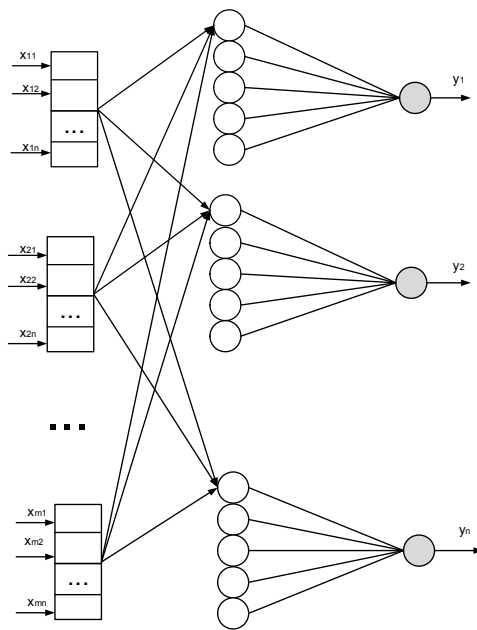


Рис. 2. Структура нейронної мережі

На рис. 2. представлена архітектура нейронної мережі, де  $n$  – кількість фонем в слові,  $m=9$  – кількість смуг.

Для навчання нейронної мережі необхідно вибрати алгоритм навчання. Існує велике число алгоритмів навчання, орієнтованих на рішення різних задач. Серед них виділяється алгоритм зворотного поширення помилки, що є одним з найбільш успішних сучасних алгоритмів. Його основна ідея полягає в тім, що зміна ваг синапсів відбувається з урахуванням локального градієнта функції помилки. Різниця між реальними і правильними відповідями нейронної мережі, обумовленими на вихідному шарі, поширюється в зворотному напрямку – назустріч потоку сигналів. У підсумку кожен нейрон здатний визначити внесок кожної своєї ваги в сумарну помилку мережі.

Найпростіше правило навчання відповідає методу найшвидшого спуску, тобто зміни синаптичних ваг пропорційно їх внеску в загальну помилку.

Навчання нейронної мережі здійснюється послідовно за допомогою алгоритму зворотного поширення [3]. Суть алгоритму зворотного поширення полягає в наступному:

Сигнал похибки вихідного нейрона  $j$  на ітерації  $n$  (відповідає  $n$  – прикладу навчання) визначається співвідношенням:

$$e_j(n) = d_j(n) - y_j(n) \quad (4)$$

Поточне значення загальної енергії помилки обчислюється по формулі:

$$E(n) = \frac{1}{2} \sum_{j \in C} e_j^2(n) \quad (5)$$

де множина  $C$  включає всі нейрони вихідного шару мережі.

Енергія середньоквадратической помилки:

$$E_{av}(n) = \frac{1}{N} \sum_{n=1}^N E(n) \quad (6)$$

Метою процесу навчання є настроювання вільних параметрів мережі з метою мінімізації величини  $E_{av}(n)$ .

Корекція  $\Delta w_{ij}(n)$ , застосовувана до ваг  $w_{ij}(n)$ , визначається відповідно до дельта-правила:

$$\Delta w_{ij}(n) = -\eta \frac{\partial E(n)}{\partial w_{ij}(n)} \quad (7)$$

Після того як мінімізували величину  $E_{av}(n)$ , водимо ідентифікатор слова і записуємо ваги мережі у файл еталона.

Збереження даних по навченій нейронній мережі організуємо у форматі xml, що дозволяє створювати універсальний файл еталона, що може бути використаний у будь-якій системі по розпізнаванню мови. У такий спосіб файл еталона зі значеннями ваг нейронної

```
<networks>
<word>
  <id>ноль</id>
  <phonems>3</phonems>
  <layer>
    <line>
      <item>1,12</item>
      ...
      <item>0,49</item>
    </line>
    ...
  </layer>
</layer>
...
</word>
...
<word>
<word>
  <id>девять</id>
  <phonems>5</phonems>
  <layer>
    ...
  </layer>
  <layer>
    ...
  </layer>
</word>
</networks>
```

Рис. 3. Структура файла еталона

мережі буде стерпним, не залежати від платформи, на якій розробляється програмне забезпечення, не залежати від мови розробки і також використання файлу в такому форматі полегшує реконфігурацію програмного забезпечення. Приклад створеного файлу еталона представлений на рис.3.

Провівши навчання нейронної мережі на вибірці достатнього розміру, можна проводити розпізнавання.

Для розпізнавання вхідний сигнал обробляється по тій же принципі, що і при навчанні. На вхід нейронної мережі надходить спектрально-смугове представлення сигналу, з файлу еталона послідовно зчитуються вагові коефіцієнти для нейронної мережі й у залежності від близькості до 1 вхідних значень приймається рішення про результат розпізнавання.

На рис.4. представлено результат роботи програми по прийняттю рішень на основі нейронної мережі. У таблиці відображені відстані до кожного з еталонів, структура нейронної мережі для слова "шість" і результат розпізнавання .

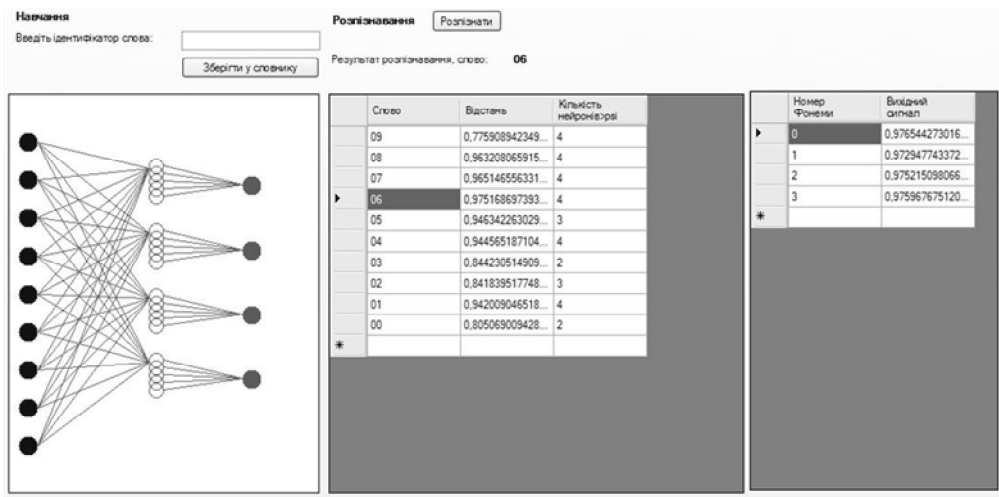


Рис. 4. Результат розпізнавання слова «шесть»

**Висновок**

Використання нейронної мережі запропонованої архітектури, дозволяє поліпшити якість розпізнавання

вхідних сигналів, завдяки вирішенню внутрішньої задачі апроксимації. А використання формату xml для збереження даних про нейронної мережі, дозволяє використовувати навчену нейрону мережу в будь-якому іншо-

му програмному забезпеченні та полегшує його реконфігурацію.

#### ЛІТЕРАТУРА

1. Карпов О.Н., Габович А.Г., Марченко Б.Г., Хорошко В.А., Щербак Л.Н. Компьютерные технологии распознавания речевых сигналов. Монография. –К.: ООО "Полиграф-Консалтинг", 2005 .–138 с. : ил. Парал. тит. англ.
2. Саймон Хайкин. Нейронные сети: полный курс, 2-е изд.: Пер. с англ. – М. : ООО «И.Д. Вильямс», 2006. – 1104 с.
3. Дж. Ф. Люгер Искусственный интеллект: Стратегии и методы решения сложных проблем, 4-е издание.: Пер. с англ. – М.: Издательский дом «Вильямс», 2003. – 804 с. : ил. Парал. тит. англ.