

УДК 004.934.11

Федоров Є.Є.

Методика нейромережного аналізу форм слів

Постановка проблеми. На сьогодні актуальною є розробка моделей і алгоритмів інтелектуальних систем керування, що базуються на природно-мовних інтерфейсах. При витягу знань із текстів важливу роль відіграє аналіз слів і їхніх форм.

Аналіз досліджень. Аналіз останніх досягнень і публікацій [1-3], присвячених проблемі аналізу і синтезу форм слів, дозволяє зробити висновок, що ці моделі не враховують морфонологічні перетворення і не використовують кількісні оцінки мовних структур.

Рішення задачі. У даній статті пропонується методика нейромережного аналізу форм слів, для якої розробляються такі правила:

- а) формалізації і чисельного дослідження елементів аналізу форм слів;
- б) закріплення взаємозв'язків між елементами аналізу форм слів;
- в) аналізу форм слів.

Основний матеріал. Перед аналізом форм слів проводиться формалізація таких мовних конструкцій: частин мови H_i ; значень Z_m^k морфологічних ознак Z^k і їхніх наборів \bar{Z}_j ; букв A_i і фонем Φ_i ; слів \bar{C}_r^s ; основ слів $\bar{\Lambda}_q^s$; форм основ слів $\bar{\Lambda}_q^s$, $\bar{\Lambda}_q^s$, отриманих при лінійних і нелінійних морфонологічних перетвореннях; неосновних частин слів \bar{V}_u^s ; словозмінних афіксів (суфіксів \bar{B}_p^{s1} , флексій \bar{B}_p^{s2} , постфіксів \bar{B}_p^{s3}) і їхніх наборів \bar{B}_p^s ; буквених і фонемних послідовностей, що усикаються/нарощуються ($\bar{\Delta}_w^{s1}/\bar{\Delta}_w^{s2}$) і чергуються ($\bar{\Delta}_v^{s1}/\bar{\Delta}_v^{s2}$), і їхніх наборів $\bar{\Delta}_w^s, \bar{\Delta}_v^s$, що беруть участь у лінійних і нелінійних морфонологічних перетвореннях основ слів; форм слів $\bar{\Omega}_i^s$, при цьому $s=1$ – якщо буква, $s=2$ – якщо фонема.

Цим конструкціям привласнюються кількісні оцінки:

- а) ранги – $r(H_i), r(Z_m^k), r(A_i), r(\Phi_i)$;
- б) інформаційні міри (що обчислюються для векторів рангів морфологічних ознак, букв і фонем) – $M(\bar{Z}_j), M(\bar{C}_r^s), M(\bar{\Lambda}_q^s), M(\bar{\Lambda}_q^s), M(\bar{\Lambda}_q^s), M(\bar{V}_u^s), M(\bar{B}_p^{sk}), M(\bar{B}_p^s), M(\bar{\Delta}_w^s), M(\bar{\Delta}_v^{sk}), M(\bar{\Delta}_w^s), M(\bar{\Delta}_v^s), M(\bar{\Omega}_i^s)$.

Далі вводяться матриці бінарних відносин $\Gamma(H_i, \bar{\Omega}_i^s), \Gamma(H_i, \bar{B}_p^s), \Gamma(\bar{Z}_j, \bar{\Omega}_i^s), \Gamma(\bar{Z}_j, \bar{B}_p^s), \Gamma(H_i, \bar{Z}_j), \Gamma(H_i, \bar{C}_r^s), \Gamma(H_i, \bar{V}_u^s)$, що установлюють взаємозв'язок між формою слова $\bar{\Omega}_i^s$ і частиною мови H_i , набором словозмінних афіксів \bar{B}_p^s і частиною мови H_i , формою слова $\bar{\Omega}_i^s$ і набором морфологічних ознак \bar{Z}_j , набором словозмінних афіксів \bar{B}_p^s і набором морфологічних ознак \bar{Z}_j , набором морфологічних ознак \bar{Z}_j і частиною мови H_i , словом \bar{C}_r^s і частиною мови H_i ,

неосновною частиною слова \bar{V}_u^s і частиною мови H_i .

При аналізі форми слова відповідно до частини мови і набором морфологічних ознак виробляється визначення слова (у початковій формі) і набору словозмінних афіксів.

При використанні рангів й інформаційних мір мовних конструкцій, а також матриць бінарних відносин здійснюється нейромережний аналіз форм слів.

На рис.1 наведена п'ятишарова повнозв'язна з прямими і зворотними зв'язками нейронна мережа аналізу форм слів.

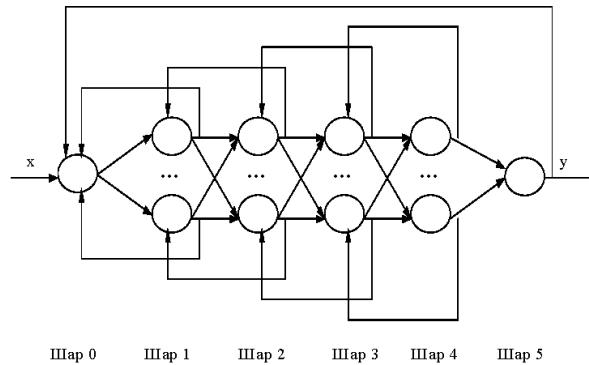


Рис. 1. Нейронна мережа аналізу форм слів

Нульовий (вхідний) шар містить один нейрон і здійснює зіставлення інформаційної міри форми слова x_1 з інформаційною мірою форми слова $M(\bar{\Omega}_l^s)$ з бази даних у вигляді

$$\rho(x_1, M(\bar{\Omega}_l^s)) = |x_1 - M(\bar{\Omega}_l^s)|, 77$$

зіставлення рангу частини мови слова і його форм x_2 з рангом частини мови $r(H_i)$ з бази даних у вигляді:

$$\rho(x_2, r(H_i)) = |x_2 - r(H_i)|,$$

зіставлення інформаційної міри набору значень морфологічних ознак форми слова x_3 з інформаційною мірою набору значень $M(\bar{Z}_{1_j}^s)$ з бази даних у вигляді:

$$\rho(x_3, M(\bar{Z}_{1_j}^s)) = |x_3 - M(\bar{Z}_{1_j}^s)|.$$

Якщо $\exists M(\bar{\Omega}_l^s) \rho(x_1, M(\bar{\Omega}_l^s)) = 0$ й $\exists r(H_i) \rho(x_2, r(H_i)) = 0$ і $\exists M(\bar{Z}_{1_j}^s) \rho(x_3, M(\bar{Z}_{1_j}^s)) = 0$, то встановлюється відповідність $x_1 \leftrightarrow M(\bar{\Omega}_l^s)$, $x_2 \leftrightarrow r(H_i)$, $x_3 \leftrightarrow M(\bar{Z}_{1_j}^s)$ і по прямому зв'язку здійснюється перехід до першого шару.

Перший шар складається з $\eta(\bar{B}_1^s)$ нейронів, кожний з яких відповідає вектору словозмінних афіксів (суфікс, флексія, постфікс) $\mu(\bar{B}_1^s) = (M(\bar{\sigma}_{p1}^s), M(\bar{\sigma}_{p2}^s), M(\bar{\sigma}_{p3}^s))$, $M(\bar{\sigma}_{pk}^s) \in M(\bar{B}^{sk})$ з інформаційною мірою $M(\bar{B}_1^s)$. На цьому шарі виконується виділення інформаційної міри форм основ (з нарощенням/усіканням і чергуванням букв/фонем) шляхом усікання інформаційної міри форми слова $M(\bar{\Omega}_l^s)$, що надійшла із вхідного шару, за допомогою інформаційних мір суфікса $M(\bar{\sigma}_{p1}^s)$, флексії $M(\bar{\sigma}_{p2}^s)$ і постфікса $M(\bar{\sigma}_{p3}^s)$ в такому вигляді:

$$M(\bar{\zeta}^s) = M(\bar{\Omega}_l^s) \circ (M(\bar{\sigma}_{p1}^s) \oslash M(\bar{\sigma}_{p2}^s) \oslash M(\bar{\sigma}_{p3}^s)).$$

Отриманий результат зіставляється з інформаційною мірою форми основи (з нарощенням/усіканням і чергуванням букв/фонем) $M(\bar{\Lambda}_q^S)$ з бази даних у вигляді:

$$\rho(M(\zeta^S), M(\bar{\Lambda}_q^S)) = |M(\zeta^S) - M(\bar{\Lambda}_q^S)|,$$

Якщо $\forall M(\bar{\Lambda}_q^S) \quad \rho(M(\zeta^S), M(\bar{\Lambda}_q^S)) > 0$ чи

$\neg(\Gamma(H_i, \bar{\Omega}_i^S) = 1 \wedge \Gamma(H_i, \bar{B}_i^S) = 1 \wedge \Gamma(\bar{Z}_i^S, \bar{\Omega}_i^S) = 1 \wedge \Gamma(\bar{Z}_i^S, \bar{B}_i^S) = 1 \wedge \Gamma(H_i, \bar{Z}_i^S) = 1)$, то по зворотному зв'язку здійснюється повернення до вхідного шару.

Якщо $\exists M(\bar{\Lambda}_q^S) \quad \rho(M(\zeta^S), M(\bar{\Lambda}_q^S)) = 0$ і

$\Gamma(H_i, \bar{\Omega}_i^S) = 1 \wedge \Gamma(H_i, \bar{B}_i^S) = 1 \wedge \Gamma(\bar{Z}_i^S, \bar{\Omega}_i^S) = 1 \wedge \Gamma(\bar{Z}_i^S, \bar{B}_i^S) = 1 \wedge \Gamma(H_i, \bar{Z}_i^S) = 1$, то встановлюється відповідність $M(\zeta^S) \leftrightarrow M(\bar{\Lambda}_q^S)$ і по прямого зв'язку здійснюється перехід до другого шару.

Другий шар складається з $\eta(\bar{\Delta}_1^S)$ нейронів, кожний з яких відповідає вектору буквених/фонемних послідовностей $\mu(\bar{\Delta}_1^S) = (M(\tilde{\delta}_{v1}^S), M(\tilde{\delta}_{v2}^S))$, $M(\tilde{\delta}_{vk}^S) \in M(\bar{\Lambda}^{sk})$, що чергуються, з інформаційною мірою $M(\bar{\Delta}_1^S)$. На цьому шарі здійснюється формування форми основи слова, що породжує (з нарощенням/усіканням букв/фонем, без чергування) шляхом перетворення інформаційних мір форми слова (з чергуванням і нарощенням/усіканням букв/фонем $M(\bar{\Lambda}_q^S)$), отриманої на першому шарі за допомогою інформаційних мір що $M(\tilde{\delta}_{v1}^S)$ чергуються $M(\tilde{\delta}_{v2}^S)$, буквених/фонемних послідовностей у такому вигляді:

а) якщо $M(\tilde{\delta}_{v2}^S) = 0 \wedge M(\tilde{\delta}_{v1}^S) \neq 0$, тобто в одержуваній формі основи слова з'являється гласна між приголосними, що йдуть наприкінці основи, використовується правило

$$\begin{aligned} & |\bar{\Lambda}_q^S| = \text{len} \wedge (r(\bar{\Lambda}_{qj}^S) \in \{r(T_i^S)\} \wedge r(\bar{\Lambda}_{q, j+1}^S) \in \{r(T_i^S)\}) \wedge \\ & \left((j+1 = \text{len}) \vee \left(\bigwedge_{k=j+1}^{\text{len}-1} \neg(r(\bar{\Lambda}_{qk}^S) \in \{r(T_i^S)\}) \wedge r(\bar{\Lambda}_{q, k+1}^S) \in \{r(T_i^S)\} \right) \right) \rightarrow \\ & \left(\bigwedge_{p=1}^j r(\zeta_p^S) = r(\bar{\Lambda}_{qp}^S) \right) \wedge r(\zeta_{j+1}^S) = r(\tilde{\delta}_{v1}^S) \wedge \left(\bigwedge_{b=j+2}^{\text{len}} r(\zeta_b^S) = r(\bar{\Lambda}_{q, b-1}^S) \right), j \in 1, |\bar{\Lambda}_q^S| - 1, \end{aligned}$$

де $\{r(T_i^1)\}$ – множина приголосних,

б) якщо $M(\tilde{\delta}_{v2}^S) \neq 0 \wedge M(\tilde{\delta}_{v1}^S) = 0$, тобто в одержуваній формі основи слова зникає гласна між приголосними, що йдуть наприкінці основи, використовується правило

$$\begin{aligned} & |\bar{\Lambda}_q^S| = \text{len} \wedge (r(\bar{\Lambda}_{q, j-1}^S) \in \{r(T_i^S)\} \wedge r(\bar{\Lambda}_{qj}^S) = r(\tilde{\delta}_{v2}^S) \wedge r(\bar{\Lambda}_{q, j+1}^S) \in \{r(T_i^S)\}) \wedge \\ & ((j+1 = \text{len}) \vee \\ & \left(\bigwedge_{k=j+1}^{\text{len}-1} \neg(r(\bar{\Lambda}_{q, k-1}^S) \in \{r(T_i^S)\}) \wedge r(\bar{\Lambda}_{qk}^S) = r(\tilde{\delta}_{v2}^S) \wedge r(\bar{\Lambda}_{q, k+1}^S) \in \{r(T_i^S)\} \right)) \rightarrow \\ & \left(\bigwedge_{p=1}^j r(\zeta_p^S) = r(\bar{\Lambda}_{qp}^S) \right) \wedge \left(\bigwedge_{b=j+1}^{\text{len}-1} r(\zeta_b^S) = r(\bar{\Lambda}_{q, b+1}^S) \right), j \in 2, |\bar{\Lambda}_q^S| - 1; \end{aligned}$$

в) якщо $M(\tilde{\delta}_{v2}^S) \neq 0 \wedge M(\tilde{\delta}_{v1}^S) \neq 0$, тобто в одержуваній формі основи слова заміняються послідовності букв/фонем, що йдуть наприкінці основи, використовується правило

$$|\tilde{\delta}_{v2}^S| = \text{len}1 \wedge |\tilde{\delta}_{v1}^S| = \text{len}2 \wedge |\bar{\Lambda}_q^S| = \text{len}3 \wedge \left(\bigwedge_{m=1}^{\text{len}1} r(\bar{\Lambda}_{q, j+m}^S) = r(\tilde{\delta}_{v2m}^S) \right) \wedge$$

$$\left((j + len1 = len3) \vee \bigwedge_{k=j+1}^{len3} \neg \left(\bigwedge_{n=1}^{len1} (r(\bar{\Lambda}_{q,k+n}) = r(\bar{\delta}_{v2n}^s)) \right) \right) \rightarrow$$

$$\left(j = 0 \vee \left(\bigwedge_{p=1}^j r(\zeta_p^s) = r(\bar{\Lambda}_{qp}^s) \right) \right) \wedge \left(\bigwedge_{z=1}^{len2} r(\zeta_{(j+1)+z}^s) = r(\bar{\delta}_{v1z}^s) \right) \wedge$$

$$\left(\bigwedge_{b=1}^{len3-(j+len1)} r(\zeta_{(j+len2+1)+b}^s) = r(\bar{\Lambda}_{q,(j+len1+1)+b}^s) \right), \quad j \in 0, |\bar{\Lambda}_q^s| - 1.$$

Таким чином, формується вектор рангів основи

$$\mu(\zeta^s) = (r(\zeta_1^s), \dots, r(\zeta_t^s), \dots, r(\zeta_M^s)).$$

Інформаційна міра цього вектора обчислюється у виді

$$M(\zeta^s) = \|\mu(\zeta^s)\| = \sum_{t=1}^{|\zeta^s|} \left(r(\zeta_t^s) * 2^{(|\zeta^s| - t)n} \right)$$

де n – кількість біт, що приходяться на одну букву/фонему.

Отриманий результат зіставляється з інформаційною мірою форми основи слова $M(\bar{\Lambda}_q^s)$ з бази даних у вигляді:

$$\rho(\zeta^s, M(\bar{\Lambda}_q^s)) = |\zeta^s - M(\bar{\Lambda}_q^s)|$$

Якщо $\forall M(\bar{\Lambda}_q^s) \rho(\zeta^s, M(\bar{\Lambda}_q^s)) > 0$, то по зворотному зв'язку здійснюється повернення до першого шару.

Якщо $\exists M(\bar{\Lambda}_q^s) \rho(\zeta^s, M(\bar{\Lambda}_q^s)) = 0$, то встановлюється відповідність $\zeta^s \leftrightarrow M(\bar{\Lambda}_q^s)$ і по прямому зв'язку здійснюється перехід до третього шару.

Третій шар складається з $\eta(\bar{\Delta}_1^s)$ нейронів, кожний з яких відповідає вектору буквених/фонемних послідовностей $\mu(\bar{\Delta}_1^s) = (M(\bar{\delta}_{w1}^s), M(\bar{\delta}_{w2}^s))$, $M(\bar{\delta}_{wk}^s) \in M(\bar{\Delta}^{sk})$, що нарощуються/усікаються, з інформаційною мірою $M(\bar{\Delta}_w^s)$. На цьому шарі здійснюється формування основи слова, що породжує шляхом перетворення інформаційної міри форми слова (з нарощенням/усіканням букв/фонем, без чергування $M(\bar{\Delta}_q^s)$), отриманої на другому шарі, за допомогою інформаційних мір, що $M(\bar{\delta}_{w1}^s)$ усікаються і $M(\bar{\delta}_{w2}^s)$ нарощуваних буквених/фонемних послідовностей у вигляді:

$$M(\zeta^s) = (M(\bar{\Delta}_q^s) \circ M(\bar{\delta}_{w2}^s)) \diamond M(\bar{\delta}_{w1}^s).$$

Отриманий результат зіставляється з інформаційною мірою форми основи $M(\bar{\Lambda}_q^s)$ з бази даних у вигляді:

$$\rho(\zeta^s, M(\bar{\Lambda}_q^s)) = |\zeta^s - M(\bar{\Lambda}_q^s)|.$$

Якщо $\forall M(\bar{\Lambda}_q^s) \rho(M(\zeta^s), M(\bar{\Lambda}_q^s)) > 0$, то по зворотному зв'язку здійснюється повернення до другого шару.

Якщо $\exists M(\bar{\Lambda}_q^s) \rho(M(\zeta^s), M(\bar{\Lambda}_q^s)) = 0$, то встановлюється відповідність $M(\zeta^s) \leftrightarrow M(\bar{\Lambda}_q^s)$ і по прямому зв'язку здійснюється перехід до четвертого шару.

Четвертий шар складається з $\eta(\bar{V}^s)$ нейронів, кожний з яких відповідає інформаційній мірі неосновної частини слова $M(\bar{V}_u^s)$. На цьому шарі виконується формування слова, що породжує, шляхом конкатенації інформаційної міри основи

слова, що $M(\bar{\Lambda}_q^S)$ породжує, отриманої на третьому шарі, за допомогою інформаційної міри неосновної частини $M(\bar{V}_u^S)$ слова в такому вигляді:

$$M(\zeta^S) = M(\bar{\Lambda}_q^S) \diamond M(\bar{V}_u^S).$$

Отриманий результат зіставляється з інформаційною мірою слова, що $M(\bar{C}_r^S)$ породжує, з бази даних у вигляді

$$\rho(M(\zeta^S), M(\bar{C}_r^S)) = |M(\zeta^S) - M(\bar{C}_r^S)|.$$

Якщо $\forall M(\bar{C}_r^S) \rho(M(\zeta^S), M(\bar{C}_r^S)) > 0$ чи $\neg(\Gamma(H_i, \bar{C}_r^S) = 1 \wedge \Gamma(H_i, \bar{V}_u^S) = 1)$, то по зворотному зв'язку здійснюється повернення до другого шару.

Якщо $\exists M(\bar{C}_r^S) \rho(M(\zeta^S), M(\bar{C}_r^S)) = 0$ і $\Gamma(H_i, \bar{C}_r^S) = 1 \wedge \Gamma(H_i, \bar{V}_u^S) = 1$, то встановлюється відповідність $M(\zeta^S) \leftrightarrow M(\bar{C}_r^S)$ і по прямому зв'язку здійснюється перехід до п'ятого шару.

П'ятий (вихідний) шар містить один нейрон і здійснює зіставлення сформованої на третьому шарі інформаційної міри слова $M(C_r^S)$, що породжує, отриманого на четвертому шарі, з інформаційною мірою очікуваного на виході слова y_1 , що породжує, у вигляді:

$$\rho(M(C_r^S), y_1) = |M(C_r^S) - y_1|.$$

Якщо $\rho(M(C_r^S), y_1) > 0$, то по зворотному зв'язку здійснюється повернення до вхідного шару.

Таким чином встановлюється асоціативний зв'язок $\bar{\Omega}_i^S \rightarrow \bar{B}1_p^S \rightarrow \bar{\Delta}1_v^S \rightarrow \bar{\Delta}1_w^S \rightarrow \bar{V}_u^S \rightarrow \bar{C}_r^S$

Наукова новизна. У даній роботі розроблені правила лінійних і нелінійних морфологічних перетворень. Для роботи з мовними конструкціями використовувалися кількісні оцінки – ранги й інформаційні міри.

Практична значимість. Основні положення роботи можуть бути реалізовані в інтелектуальній системі у вигляді алгоритмів, що забезпечують спілкування з користувачем природною мовою.

Література

1. Кисленко Ю.И. Информационный робот // Искусственный интеллект. – 2001. – № 1. – С. 61-68.
2. Cheblakov G.B., Dinenberg F.G., Levin D.Y., Popov I.G., Zagorulko Y.A. Approach to development of a system for speech interaction with an intelligent robot // Perspectives of System Informatics. – Berlin: Springer-Verlag. – 1999. – P. 517-529.
3. Литвинцева Л.В., Ульянов С.В., Танака Т., Охви Д., Ямафудзи К. Интеллектуальное управление мобильным роботом сервисного обслуживания // Труды VIII науч.-техн. конф. "Экстремальная робототехника". – Санкт-Петербург. – 1997. – С. 35-43.
4. Современный русский язык: Учеб. для филол. спец. высших учебных заведений / В.А.Белошапкова, Е.А.Брызгунова, Е.А.Земская и др.; Под ред. В.А.Белошапковой. – М.: Азбуковник, 1997. – 928 с.
5. Русская грамматика: В 2 т. – М.: Наука, 1980. – т. 1: Фонетика. Словообразование. Морфология. – 784 с.
6. Кривооубский О.А., Федоров Е.Е. Формальное представление русского языка и речи // Искусственный интеллект. – 2003. – № 4. – С. 402-410.