

ДИФFUЗНЫЕ АЛГОРИТМЫ ОБУЧЕНИЯ НЕЙРОННЫХ СЕТЕЙ ПРЯМОГО РАСПРОСТРАНЕНИЯ

Ключевые слова: *нейронные сети прямого распространения, алгоритмы обучения, расширенный фильтр Калмана.*

ВВЕДЕНИЕ

Нейронные сети прямого распространения (НСПР) широко используются в разных приложениях, связанных с такими задачами, как прогнозирование временных рядов, классификация данных, идентификация и управление нелинейными объектами. Обзоры по методам их решения приведены в [1–3].

Обучение НСПР может рассматриваться как задача минимизации среднеквадратической ошибки относительно неизвестных параметров, входящих в ее описание (весов и смещений) при заданном обучающем множестве. Алгоритм обратного распространения (АОР) разработан в [4] и успешно применяется для их обучения. Вместе с тем известна присущая АОР медленная скорость сходимости, что существенно затрудняет или делает практически невозможным использование алгоритма в сложных задачах.

В многочисленных публикациях, посвященных НСПР, предложены различные алгоритмы обучения, превосходящие АОР по скорости сходимости и получаемой точности аппроксимации. В [5, 6] применяются методы Левенберга–Маквардта и квазиньютоновские, использующие информацию о матрице производных второго порядка критерия качества. В [7–9] алгоритмы обучения основываются на расширенном фильтре Калмана (РФК), использующем аппроксимационную ковариационную матрицу ошибки оценивания. Для НСПР с линейной функцией активации (ФА) в выходном слое в [10–12] представлены алгоритмы, учитывающие сепарабельную структуру сети [13]. В соответствии с VP-алгоритмом (variable projection) [10] исходная задача оптимизации преобразуется к эквивалентной относительно нелинейно входящих параметров (весов и смещений скрытого слоя). При этом уменьшается как размерность задачи, так и ее обусловленность, что позволяет сократить количество итераций для получения решения. Вместе с тем такой подход может использоваться только в пакетном режиме, и даже при аддитивно входящих ошибках измерений выходов НСПР в преобразованный критерий они входят нелинейно. Кроме того, существенно усложняется процедура определения частных производных критерия по параметрам. В [11, 12] предлагается ELM (extreme learning machine) алгоритм, с помощью которого обучаются только линейно входящие параметры (веса выходного слоя), а нелинейные выбираются случайно, без учета обучающей выборки, что сокращает время обучения, но может приводить к невысокой точности аппроксимации.

В данной работе предлагаются и исследуются новые алгоритмы обучения двухслойных НСПР с нелинейными ФА в скрытом слое и линейными — в выходном. Известно, что такие нейронные сети могут обладать универсальными аппроксимирующими свойствами (для сигмоидных функций доказательство приведено в [14]). Алгоритмы основываются на РФК и учете сепарабельной структуры НСПР, как в VP- и ELM-алгоритмах, но при этом одновременно обучаются все нейроны. Более точно, линейно входящие веса интерпретируются как диффузные — случайные величины, имеющие нулевое математическое ожида-

ние и матрицу ковариации, пропорциональную произвольно большому параметру λ [15]. Находятся асимптотические представления РФК при $\lambda \rightarrow \infty$, которые мы называем диффузными алгоритмами обучения (ДАО). Показано, что они, в отличие от их прототипа РФК с большим, но конечным λ , обладают свойством робастности по отношению к накоплению ошибок округления, и из ДАО при определенных упрощающих предположениях следует ELM-алгоритм. Приведен численный пример, показывающий, что ДАО могут превосходить ELM-алгоритм по точности аппроксимации.

1. ОСНОВНЫЕ ДОПУЩЕНИЯ И СООТНОШЕНИЯ ДЛЯ РФК

Рассмотрим двухслойную НСПР с нелинейными ФА в скрытом слое и линейными в выходном:

$$y_{it} = \sum_{k=1}^m w_{ik} \sigma \left(\sum_{j=1}^{n-1} a_{kj} z_{jt} + b_k \right), \quad i=1, \dots, r, \quad t=1, \dots, N. \quad (1)$$

Здесь z_{jt} , $j=1, \dots, n$, — входы, y_{it} , $i=1, \dots, r$, — выходы, a_{kj} , b_k , $k=1, \dots, m$, $j=1, \dots, n-1$, — веса и смещения скрытого слоя, w_{ik} , $i=1, 2, \dots, r$, $k=1, 2, \dots, m$, — веса выходного слоя, $\sigma(x) \in C^1(R)$ — ФА, $C^1(R)$ — пространство дифференцируемых функций на действительной прямой, или, в более компактной форме, используются векторно-матричные обозначения

$$y_t = W(\sigma(v_1 z_t), \dots, \sigma(v_m z_t))^T = f(z_t, \beta, \alpha) \quad t=1, \dots, N, \quad (2)$$

где $y_t = (y_{1t}, \dots, y_{rt})^T \in R^r$, $z_t = (z_{1t}, \dots, z_{n-1t}, 1)^T \in R^n$, $W = (w_1^T, \dots, w_r^T)^T \in R^{r \times m}$, $w_i = (w_{i1}, w_{i2}, \dots, w_{im})$, $i=1, 2, \dots, r$, $v_k = (v_{k1}, \dots, v_{kn})$, $k=1, \dots, m$, $\beta = (v_1^T, \dots, v_m^T)^T \in R^{mn}$, $\alpha = (w_1^T, \dots, w_r^T)^T \in R^{rm}$, R^l , $R^{l \times p}$ — пространства векторов и матриц размерностей l и $l \times p$ соответственно, $(\cdot)^T$ — операция транспонирования соответствующей матрицы.

Обучение нейронных сетей при последовательной обработке и обучающем множестве $\{z_t, y_t\}_{t=1}^N$ рассматривается как задача оценки состояния нелинейной динамической системы

$$\begin{aligned} \beta_{t+1} &= \beta_t, \quad \alpha_{t+1} = \alpha_t, \\ y_t &= f(z_t, \beta_t, \alpha_t) + \xi_t, \quad t=1, \dots, N, \end{aligned} \quad (3)$$

где $\xi_t \in R^r$ — случайный процесс с некоррелированными значениями, нулевым математическим ожиданием и матрицей ковариации $E[\xi_t \xi_t^T] = V_t$, характеризующий ошибки измерения выхода. Начальное состояние (3) удовлетворяет следующим условиям:

А1) векторы β_1, α_1 случайны и некоррелированы между собой и ξ_t , $t=1, \dots, N$;

А2) задана выборка весов и смещений скрытого слоя нейронных сетей $m\beta = (\bar{\beta}_1, \dots, \bar{\beta}_{mn})^T$ из непрерывного распределения с матрицей ковариации S_β , определяемая характером решаемой задачи и используемой ФА, при этом $\beta_1 = m\beta$;

А3) априорная информация относительно весов выходного слоя α отсутствует, и они интерпретируются как случайные величины, имеющие нулевое математическое ожидание и матрицу ковариации, пропорциональную большому параметру $\lambda > 0$, т.е.

$$E(\alpha) = 0, \quad E(\alpha\alpha^T) = \lambda I_{rm}, \quad (4)$$

где $I_{rm} \in R^{rm \times rm}$ — единичная матрица.

Оценка состояния (3) $x_t = (\beta_t^T, \alpha_t^T)^T$ с помощью РФК удовлетворяет нелинейному разностному уравнению [16]

$$\hat{x}_{t+1} = \hat{x}_t + K_t(y_t - \tilde{f}(z_t, \hat{x}_t)), \quad \hat{x}_1 = (m^T, 0_{rm}^T)^T, \quad t = 1, \dots, N, \quad (5)$$

где $\tilde{f}(z_t, x_t) = f(z_t, \beta_t, \alpha_t)$, $0_{rm} \in R^{rm}$ — вектор с нулевыми элементами,

$$K_t = P_t C_t^T N_t^{-1}, \quad N_t = C_t P_t C_t^T + V_t, \quad (6)$$

$$P_{t+1} = P_t - P_t C_t^T N_t^{-1} C_t P_t, \quad P_1 = \text{block diag}(S_\beta, \lambda I_{rm}), \quad (7)$$

$$C_t = ((C_t^\beta)^T, (C_t^\alpha)^T)^T, \quad C_t^\beta = \left. \frac{\partial \tilde{f}(z_t, x_t)}{\partial \beta_t} \right|_{x_t = \hat{x}_t} \in R^{r \times mn}, \quad (8)$$

$$C_t^\alpha = \left. \frac{\partial \tilde{f}(z_t, x_t)}{\partial \alpha_t} \right|_{x_t = \hat{x}_t} \in R^{r \times rm},$$

матрица C_t определяется по выражениям

$$\frac{\partial y_{it}}{\partial v_k} = w_{ik} \sigma^{(1)}(v_k z_t) z_t, \quad i = 1, \dots, r, \quad k = 1, \dots, m, \quad (9)$$

$$\frac{\partial y_{it}}{\partial w_l} = (\sigma(v_1 z_t), \sigma(v_2 z_t), \dots, \sigma(v_m z_t))^T, \quad i = 1, \dots, r, \quad (10)$$

$\sigma^{(1)}(x)$ — производная $\sigma(x)$.

При пакетной обработке и обучающем множестве $\{z_t, y_t\}_{t=1}^N$ обучение рассматривается как задача оценки состояния нелинейной динамической системы

$$\beta_{t+1} = \beta_t, \quad \alpha_{t+1} = \alpha_t, \quad Y_t = F(Z, \beta_t, \alpha_t) + \xi_t, \quad t = 1, \dots, M, \quad (11)$$

где $F(Z, \beta_t, \alpha_t) = (f^T(z_1, \beta_t, \alpha_t), \dots, f^T(z_N, \beta_t, \alpha_t))^T$, $Z = (z_1^T, \dots, z_N^T)^T$, M — количество итераций. Предполагается, что начальное состояние (11), как и при последовательном обучении, удовлетворяет условиям А1)–А3).

Оценка состояния (11) с помощью РФК удовлетворяет нелинейному разностному уравнению

$$\hat{x}_{t+1} = \hat{x}_t + K_t(Y - \tilde{F}(Z, \hat{x}_t)), \quad \hat{x}_1 = (\bar{\beta}^T, 0_{rm}^T)^T, \quad t = 1, \dots, M, \quad (12)$$

где $\tilde{F}(Z, \hat{x}_t) = F(Z, \beta_t, \alpha_t)$, $Y = (y_1^T, \dots, y_N^T)^T$,

$$K_t = P_t C_t^T N_t^{-1}, \quad N_t = C_t P_t C_t^T + V_t, \quad (13)$$

$$P_{t+1} = P_t - P_t C_t^T N_t^{-1} C_t P_t, \quad P_1 = \text{block diag}(S_\beta, \lambda I_{rm}), \quad (14)$$

$$C_t = ((C_t^\beta)^T, (C_t^\alpha)^T)^T, \quad C_t^\beta = \left. \frac{\partial \tilde{F}(Z, \hat{x}_t)}{\partial \beta_t} \right|_{x_t = \hat{x}_t} \in R^{Nr \times mn}, \quad (15)$$

$$C_t^\alpha = \left. \frac{\partial \tilde{F}(Z, \hat{x}_t)}{\partial \alpha_t} \right|_{x_t = \hat{x}_t} \in R^{Nr \times rm},$$

матрица C_t определяется по выражениям (9), (10).

2. АСИМПТОТИКА РФК

Рассмотрим несколько вспомогательных утверждений.

Лемма 1. При $t=1, \dots, N$ справедливы представления

$$P_t = \tilde{S}_t + \tilde{R}_t M_t^{-1} \tilde{R}_t^T, \quad (16)$$

$$K_t = (\tilde{S}_t + \tilde{R}_{t+1} M_{t+1}^{-1} \tilde{R}_t^T) C_t^T N_{1t}^{-1}, \quad (17)$$

где

$$\tilde{S}_t = \begin{pmatrix} S_t & 0_{mn \times rm} \\ 0_{rm \times mn} & 0_{rm \times rm} \end{pmatrix}, \quad \tilde{R}_t = \begin{pmatrix} R_t \\ I_{rm} \end{pmatrix}, \quad (18)$$

$$S_{t+1} = S_t - S_t (C_t^\beta)^T N_{1t}^{-1} C_t^\beta S_t, \quad N_{1t} = C_t^\beta S_t (C_t^\beta)^T + V_t, \quad S_1 = S_\beta, \quad (19)$$

$$R_{t+1} = (I - S_t (C_t^\beta)^T N_{1t}^{-1} C_t^\beta) R_t - S_t (C_t^\beta)^T N_{1t}^{-1} C_t^\alpha, \quad R_1 = 0_{mn \times rm}, \quad (20)$$

$$M_{t+1} = M_t + (C_t^\beta R_t + C_t^\alpha)^T N_{1t}^{-1} (C_t^\beta R_t + C_t^\alpha), \quad M_1 = \frac{I_{rm}}{\lambda}, \quad \lambda > 0. \quad (21)$$

Доказательство. Пусть \tilde{P}_t и \bar{P}_t — два произвольных решения (7). Тогда [17] $Q_t = \tilde{P}_t - \bar{P}_t$, $t=1, \dots, N$, где

$$Q_{t+1} = A_t Q_t A_t^T - A_t Q_t C_t^T N_{2t}^{-1} C_t Q_t A_t, \quad Q_1 = \tilde{P}_1 - \bar{P}_1, \quad (22)$$

$$A_t = I - \bar{P}_t C_t^T N_{3t}^{-1} C_t, \quad N_{2t} = C_t \tilde{P}_t C_t^T + V_t, \quad N_{3t} = C_t \bar{P}_t C_t^T + V_t.$$

Положим $\tilde{P}_1 = P_1$, $\bar{P}_1 = \text{block diag}(S_\beta, 0)$. Тогда $\tilde{P}_t = P_t$, $\bar{P}_t = \tilde{S}_t$ и $P_t = \tilde{S}_t + Q_t$. Покажем, что $Q_t = \tilde{R}_t M_t^{-1} \tilde{R}_t^T$. Так как

$$\tilde{R}_{t+1} = A_t \tilde{R}_t, \quad \tilde{R}_1 = (e_p, \dots, e_q), \quad (23)$$

где $p = mn+1$, $q = mn+rm$, $e_i \in R^q$ — i -единичный вектор, то

$$\tilde{R}_{t+1} M_{t+1}^{-1} \tilde{R}_{t+1}^T = \tilde{R}_{t+1} M_t^{-1} \tilde{R}_{t+1}^T - \tilde{R}_{t+1} M_t^{-1} \tilde{R}_t^T C_t^T N_{2t}^{-1} C_t \tilde{R}_t M_t^{-1} \tilde{R}_{t+1}^T.$$

Это равенство выполняется, при условии, что M_t^{-1} удовлетворяет разностному уравнению

$$M_{t+1}^{-1} = M_t^{-1} - M_t^{-1} \tilde{R}_t^T C_t^T N_{2t}^{-1} C_t \tilde{R}_t M_t^{-1}. \quad (24)$$

Преобразовывая его с помощью тождества [15]

$$(P^{-1} + H^T R^{-1} H)^{-1} = P - P H^T (H P H^T + R)^{-1} H P,$$

с $P = M_t^{-1}$, $H = C_t \tilde{R}_t$, $R = N_{2t}$, получим

$$M_{t+1}^{-1} = (M_t + \tilde{R}_t^T C_t^T N_{2t}^{-1} C_t \tilde{R}_t)^{-1}, \quad (25)$$

отсюда следуют (16).

Докажем (17). Сначала покажем, что

$$K_t = P_t C_t^T N_t^{-1} = \tilde{S}_t C_t^T N_{3t}^{-1} + A_t Q_t C_t^T N_{2t}^{-1}.$$

Так как $P_t = \tilde{S}_t + Q_t$, то должно быть $(\tilde{S}_t + Q_t) C_t^T = \tilde{S}_t C_t^T N_{3t}^{-1} N_{2t} + A_t Q_t C_t^T$.

Преобразовывая правую часть этого выражения, устанавливаем, что

$$\begin{aligned} \tilde{S}_t C_t^T N_{3t}^{-1} N_{2t} + A_t Q_t C_t^T &= \tilde{S}_t C_t^T (I_m + N_{3t}^{-1} C_t Q_t C_t^T) + \\ &+ (I - \tilde{S}_t C_t^T N_{3t}^{-1} C_t) Q_t C_t^T = (\tilde{S}_t + Q_t) C_t^T. \end{aligned}$$

Используя (23), (25), находим

$$\begin{aligned}
A_t Q_t C_t^T N_{2t}^{-1} &= A_t \tilde{R}_t M_t^{-1} \tilde{R}_t^T C_t^T N_{2t}^{-1} = \\
&= \tilde{R}_{t+1} M_{t+1}^{-1} (M_t + \tilde{R}_t^T C_t^T N_{1t}^{-1} C_t \tilde{R}_t) M_t^{-1} \tilde{R}_t^T C_t^T N_{2t}^{-1} = \\
&= \tilde{R}_{t+1} M_{t+1}^{-1} (\tilde{R}_t^T C_t^T + \tilde{R}_t^T C_t^T N_{1t}^{-1} C_t \tilde{R}_t M_t^{-1} \tilde{R}_t^T C_t^T) N_{2t}^{-1} = \\
&= \tilde{R}_{t+1} M_{t+1}^{-1} \tilde{R}_t^T C_t^T N_{1t}^{-1} (N_{1t} + C_t \tilde{R}_t M_t^{-1} \tilde{R}_t^T C_t^T) N_{2t}^{-1} = \tilde{R}_{t+1} M_{t+1}^{-1} \tilde{R}_t^T C_t^T N_{1t}^{-1},
\end{aligned}$$

отсюда следует (17).

Лемма доказана.

Лемма 2. Для любых $(m \times n)$ -матриц F_t , $t = 1, \dots, N$, справедливы равенства

$$(I_n - \Xi_t \Xi_t^+) F_t^T = 0, \quad t = 1, \dots, N, \quad (26)$$

где $\Xi_t = \sum_{s=1}^t (F_s)^T F_s$, $(\cdot)^+$ — псевдообратная матрица соответствующей матрицы.

Доказательство. Без ограничения общности будем полагать, что $F_t \neq 0$. Обозначим $\tilde{F}_t = (F_1^T, \dots, F_t^T)$. Пусть $l_{1,t}, \dots, l_{k(t),t}$ — любые линейно независимые столбцы матрицы \tilde{F}_t такие, что последний $l_{k(t),t}$ из них совпадает с последним ненулевым столбцом матрицы \tilde{F}_t . Использование скелетного разложения для \tilde{F}_t дает $\tilde{F}_t = L_t \Gamma_t$, где $L_t = (l_{1,t}, \dots, l_{k(t),t})$, $\Gamma_t = (\Gamma_{1t}, \dots, \Gamma_{tt})$ — $(n \times k(t))$ -, $(k(t) \times mt)$ -матрицы ранга $k(t)$, Γ_{it} , $i = 1, \dots, t$ — $(k(t) \times m)$ -матрицы. Покажем вначале, что $I_n - \Xi_t \Xi_t^+ = I - L_t (L_t^T L_t)^{-1} L_t^T$. Имеем $\Xi_t = \tilde{F}_t \tilde{F}_t^T = L_t \tilde{\Gamma}_t L_t^T$, где $\tilde{\Gamma}_t = \Gamma_t \Gamma_t^T$. Так как $\tilde{\Gamma}_t > 0$ — матрица Грамма, построенная по линейно независимым строкам матрицы Γ_t , $\text{rank}(L_t) = \text{rank}(\tilde{\Gamma}_t L_t^T)$, и L_t — матрица полного ранга по столбцам, то $\Xi_t^+ = (L_t \tilde{\Gamma}_t L_t^T)^+ = (L_t^T)^+ \tilde{\Gamma}_t^{-1} L_t^+$, $L_t^+ = (L_t^T L_t)^{-1} L_t^T$. Таким образом,

$$\begin{aligned}
I_n - \Xi_t \Xi_t^+ &= I_n - L_t \tilde{\Gamma}_t L_t^T (L_t \tilde{\Gamma}_t L_t^T)^+ = I_n - L_t \tilde{\Gamma}_t L_t^T (L_t^T)^+ \tilde{\Gamma}_t^{-1} L_t^+ = \\
&= I_n - L_t \tilde{\Gamma}_t L_t^T L_t (L_t^T L_t)^{-1} \tilde{\Gamma}_t^{-1} (L_t^T L_t)^{-1} L_t^T = I_n - L_t (L_t^T L_t)^{-1} L_t^T.
\end{aligned}$$

Поскольку $F_t^T = L_t \Gamma_t$ для некоторой матрицы $\bar{\Gamma}_t$, то

$$(I_n - \Xi_t \Xi_t^+) F_t^T = (I_n - L_t (L_t^T L_t)^{-1} L_t^T) L_t \bar{\Gamma}_t = 0.$$

Далее будет использоваться обозначение $O(\lambda^{-k})$ для функций $\Phi(\lambda, y_t, \hat{x}_t)$, удовлетворяющих условию $P(\|\Phi(\lambda, y_t, \hat{x}_t)\| / \lambda^{-k} < \delta_1) > 1 - \delta_2$ при $\lambda \rightarrow \infty$, $k \geq 0$, некоторых постоянных $\delta_1, \delta_2 > 0$ и t , принадлежащих ограниченному множеству, где $\|\cdot\|$ — евклидова норма соответствующей матрицы.

Лемма 3. Пусть матрица Ω_t определена выражением

$$\Omega_t = \frac{\Omega_1}{\lambda} + \sum_{s=1}^{t-1} F(y_s, \hat{x}_s)^T F(y_s, \hat{x}_s), \quad t = 2, \dots, N, \quad (27)$$

где $\Omega_1 > 0$, $m \times n$ -матрица $F(y_t, \hat{x}_t)$ удовлетворяют условию $F(y_t, \hat{x}_t) = O(1)$ при $\lambda \rightarrow \infty$ и $t = 2, \dots, N$. Тогда

$$\Omega_t^{-1} = \Omega_1^{-1} (I_n - \Xi_t \Xi_t^+) \lambda + \Xi_t^+ + O(\lambda^{-1}), \quad t = 2, \dots, N, \quad (28)$$

при $\lambda \rightarrow \infty$, где $\Xi_t = \sum_{s=1}^{t-1} F(y_s, \hat{x}_s)^T F(y_s, \hat{x}_s)$.

Доказательство. Имеем $\Omega_t = \Omega_1(I_n + \lambda\Omega_a^{-1}\Xi_t) / \lambda$. Отсюда с помощью формулы обращения возмущенных матриц [18] получим

$$(I_n + \lambda V)^{-1} = (I_n - VV^+) + V^+ (V^+ + \lambda I_n)^{-1}, \quad (29)$$

где V — произвольная симметрическая матрица, находим

$$\begin{aligned} \Omega_t^{-1} &= (I_n + \lambda\Omega_1^{-1}\Xi_t)^{-1}\Omega_1^{-1}\lambda = \\ &= [(I_n - \Omega_1^{-1}\Xi_t\Xi_t^+\Omega_1)\Omega_1^{-1} + \Xi_t^+\Omega_1(\Xi_t^+\Omega_1 + \lambda I_n)^{-1}\Omega_1^{-1}]\lambda. \end{aligned}$$

Представим матрицу $\Omega_1^{-1}\Xi_t$ в виде $\Omega_1^{-1}\Xi_t = T_t^T D_t T_t$, где T_t и D_t — соответственно ортогональная и диагональная матрицы. Тогда

$$\begin{aligned} (\Xi_t^+\Omega_1 + \lambda I_n)^{-1}\lambda &= T_t(D_t^+ + \lambda I_n)^{-1}T_t^T\lambda = \\ &= I_n - \Xi_t^+\Omega_1\lambda^{-1} + O(\lambda^{-2}), \quad \lambda \rightarrow \infty. \end{aligned} \quad (30)$$

Отсюда следует (28).

Теорема 1. 1. При ограниченном N справедливы асимптотические представления

$$P_t = \tilde{R}_t(I_{rm} - W_t W_t^+) \tilde{R}_t^T \lambda + \tilde{S}_t + \tilde{R}_t W_t^+ \tilde{R}_t^T + O(\lambda^{-1}), \quad t=2, \dots, N, \quad (31)$$

$$K_t = K_t^{dif} + O(\lambda^{-1}), \quad C_t \neq 0, \quad t=1, \dots, N, \quad \lambda \rightarrow \infty, \quad (32)$$

где

$$W_{t+1} = W_t + (C_t^\beta R_t + C_t^\alpha)^T N_{1t}^{-1} (C_t^\beta R_t + C_t^\alpha), \quad W_1 = 0_{mn}, \quad (33)$$

$$K_t^{dif} = (\tilde{S}_t + \tilde{R}_{t+1} W_{t+1}^+ \tilde{R}_t^T) C_t^T N_{1t}^{-1}. \quad (34)$$

2. Если дополнительно $\sigma(x) \in C^2(R)$, то

$$\hat{x}_t = \hat{x}_t^{dif} + O(\lambda^{-1}), \quad t=2, \dots, N, \quad \lambda \rightarrow \infty, \quad (35)$$

где

$$\hat{x}_{t+1}^{dif} = \hat{x}_t^{dif} + K_t^{dif} (y_t - \tilde{f}(z_t, \hat{x}_t^{dif})), \quad \hat{x}_1^{dif} = (\bar{\beta}^T, 0_{rm}^T)^T. \quad (36)$$

Доказательство. Покажем, что $\hat{x}_t = O(1)$ при $\lambda \rightarrow \infty$, $t=2, \dots, N$. Отсюда и леммы 3, если положить в (28) $F(y_t, \hat{x}_t) = N_{1t}^{-1/2} (C_t^\beta R_t + C_t^\alpha)$, будем иметь

$$M_{t+1}^{-1} = (I_n - \Xi_t \Xi_t^+) \lambda + \Xi_t^+ + O(\lambda^{-1}), \quad t=1, \dots, N, \quad \lambda \rightarrow \infty. \quad (37)$$

Подставляя (37) в (16), получаем (31). Воспользуемся индукцией. Так как C_1 не зависит от λ , то (31) справедливо при $t=1$ и $\hat{x}_2 = O(1)$, $\lambda \rightarrow \infty$. Предположим, что $\hat{x}_i = O(1)$ при $\lambda \rightarrow \infty$ и каждом $i \in Z = \{3, \dots, t-1\}$. В силу ограниченности Z и $\hat{x}_2 = O(1)$, $\lambda \rightarrow \infty$ эта оценка будет выполняться равномерно на $\{1, \dots, t-1\}$, что влечет $C_i = O(1)$ на Z и $\hat{x}_t = O(1)$ при $\lambda \rightarrow \infty$, $t=2, \dots, N$.

В результате подстановки (37) в (17) и в силу леммы 2 получаем (32).

Воспользуемся индукцией для доказательства (35). Обозначим $e_t = \hat{x}_t - \hat{x}_t^{dif}$.

Из (5) и (36) следует

$$\begin{aligned} e_{t+1} &= e_t + (K_t(\hat{x}_t) - K_t^{dif}(\hat{x}_t - e_t))y_t - K_t(\hat{x}_t)\tilde{f}(z_t, \hat{x}_t) + \\ &+ K_t^{dif}(\hat{x}_t - e_t)\tilde{f}(z_t, \hat{x}_t - e_t), \quad t=1, \dots, N. \end{aligned} \quad (38)$$

Поскольку $e_1 = 0$, $\tilde{f}(z_t, \hat{x}_1) = 0$, $K_1 = K_1^{dif} + O(\lambda^{-1})$ (без ограничения общности полагаем, что $C_1 \neq 0$), то из (38) следует оценка $e_2 = O(\lambda^{-1})$, $\lambda \rightarrow \infty$. Пусть t произвольно и $e_t = O(\lambda^{-1})$, $\lambda \rightarrow \infty$. Если $C_t = 0$, то $e_{t+1} = e_t = O(\lambda^{-1})$. При

$C_t \neq 0$ воспользуемся оценкой $K_t = K_t^{dif} + O(\lambda^{-1})$ и теоремой о среднем

$$K_t^{dif}(\hat{x}_t - e_t) = K_t^{dif}(\hat{x}_t) + O(\lambda^{-1}), \quad \tilde{f}(z_t, \hat{x}_t - e_t) = \tilde{f}(z_t, \hat{x}_t) + O(\lambda^{-1}), \quad \lambda \rightarrow \infty.$$

Подстановка этих выражений в (38) дает $e_{t+1} = O(\lambda^{-1})$, $\lambda \rightarrow \infty$. Отсюда в силу ограниченности множества $t=1, \dots, N$ следует оценка (35).

Соотношения (36), (33), (34), (19), (20) определяют ДАО.

Следствие 1. Диффузная составляющая $P_t^{dif} = \tilde{R}_t(I_{n-q} - W_t W_t^+) \tilde{R}_t^T \lambda$ — компонента в разложении P_t , пропорциональная большому параметру, равна нулю, начиная с $t_{tr} = \min_t \{t : W_t > 0, t=2, \dots, N\}$.

Следствие 2. Матрица K_t не зависит от диффузной составляющей и в отличие от матрицы P_t как функция λ равномерно ограничена по норме при $t=1, \dots, N$ и $\lambda \rightarrow \infty$.

Следствие 3. Ошибки численной реализации могут приводить к расходимости РФК при больших λ . Действительно, пусть δW_t^+ — ошибка, связанная с вычислением псевдообратной матрицы W_t^+ . Тогда использование (17) и леммы 3 дает

$$K_t = (\tilde{S}_t + \tilde{R}_{t+1} M_{t+1}^{-1} \tilde{R}_t^T) C_t^T N_{1t}^{-1} = \\ = [\tilde{S}_t + \tilde{R}_{t+1} ((I_{rm} - W_t(W_t^+ + \delta W_t^+))\lambda + O(1)) \tilde{R}_t^T] C_t^T N_{1t}^{-1}, \quad t=1, \dots, N, \quad \lambda \rightarrow \infty. \quad (39)$$

При $\delta W_t^+ \neq 0$ матрица K_t становится зависимой от диффузной составляющей. Поскольку операция псевдообращения не является непрерывной, то вызванное ею изменение K_t может быть существенным и приводить к расходимости. Более того, поскольку матрица K_t пропорциональна большому параметру λ , то даже если и выполнено условие непрерывности операции псевдообращения $\text{rank}(W_t) = \text{rank}(W_t + \delta W_t)$, для произвольных, достаточно малых по норме матриц δW_t , диффузная составляющая по-прежнему может приводить к потере точности. Кроме того, этот эффект может усиливаться при поступлении измерений высокой точности и малой по норме матрице S_t . ДАО при численной реализации не имеют указанных особенностей. Об этом свидетельствует отсутствие в его конструкции диффузных компонент — величин, пропорциональных большому параметру, причем характеристики предельного алгоритма не зависят от неизвестной априорной информации о начальном векторе состояния.

Замечание 1. Теорема 1 и следствия 1–3 справедливы для алгоритма обучения в пакетном режиме, описываемого соотношениями (12)–(15).

3. СВЯЗЬ С ELM-АЛГОРИТМОМ

Полагая в ДАО $S_\beta = 0$, получаем

$$\hat{\alpha}_{t+1} = \hat{\alpha}_t + K_t^{\alpha, dif} (y_t - C_t^\alpha \hat{\alpha}_t), \quad \hat{\alpha}_1 = 0_{rm}, \quad \hat{\beta}_t = m_\beta^T, \quad (40)$$

$$K_t^{\alpha, dif} = W_{t+1}^+ (C_t^\alpha)^T V_t^{-1}, \quad W_{t+1} = W_t + (C_t^\alpha)^T V_t^{-1} C_t^\alpha, \quad W_a = 0_{m \times m}, \quad (41)$$

т.е., как и в ELM-алгоритме, оцениваются только веса выходного слоя, а веса и смещения скрытого слоя в процессе обучения не изменяются. Покажем, что оценки, получаемые с использованием ДАО и ELM, совпадают.

Рассмотрим задачу минимизации взвешенной суммы квадратов

$$J_t(\alpha) = \sum_{i=1}^t (y_i - C_i^\alpha \alpha)^T V_i^{-1} (y_i - C_i^\alpha \alpha)^T \quad (42)$$

при фиксированном $\beta = \bar{\beta}$. Пусть $\text{rank}(\tilde{C}_t) = rm$, где $\tilde{C}_t = ((C_1^\alpha)^T, \dots, (C_t^\alpha)^T)$.

Тогда

$$\alpha_t^{\min} = \left(\sum_{i=1}^t (C_i^\alpha)^\top V_i^{-1} C_i^\alpha \right)^{-1} \sum_{i=1}^t (C_i^\alpha)^\top V_i^{-1} y_i. \quad (43)$$

Теорема 2. Оценка ДАО удовлетворяет условию $\hat{\alpha}_{t+1} = \alpha_t^{\min}$ при любом $\hat{\alpha}_1 = \bar{\alpha}$ и $t_{tr} \geq t$.

Доказательство. Имеем

$$\hat{\alpha}_{t+1} = \Phi_{t+1,1} \bar{\alpha} + \sum_{s=1}^t \Phi_{t+1,s+1} K_s^{\alpha, dif} y_s, \quad (44)$$

где переходная матрица $\Phi_{t,s}$ удовлетворяет уравнению

$$\Phi_{t+1,s} = (I_{rm} - K_t^{\alpha, dif} C_t^\alpha) \Phi_{t,s} = A_t \Phi_{t,s}, \quad \Phi_{s,s} = I_{rm}, \quad t = s, s+1, \dots$$

Покажем, что

$$\Phi_{t,s} = W_t^{-1} W_s, \quad t \geq tr, \quad 1 \leq s \leq t, \quad (45)$$

но вначале подтвердим, что решения матричных уравнений

$$G_{t,s} = (I_{rm} - K_t^{\alpha, dif} C_t^\alpha)^\top G_{t+1,s}, \quad G_{s,s} = I_{rm},$$

$$Z_{t,s} = Z_{t+1,s} - (C_t^\alpha)^\top V_t^{-1} C_t^\alpha W_s^{-1}, \quad Z_{s,s} = I_{rm}$$

совпадают, где $s \geq R = \min \{s : W_s > 0, s = 2, \dots, N\}$, $t = s-1, \dots, 1$. Для этого долж-

но выполняться условие $(K_t^{\alpha, dif})^\top Z_{t+1,s} = V_t^{-1} C_t^\alpha W_s^{-1}$. Имеем

$$\begin{aligned} Z_{t,s} &= I_{rm} - \sum_{i=t}^{s-1} (C_i^\alpha)^\top V_i^{-1} C_i^\alpha W_s^{-1}, \\ (K_t^{\alpha, dif})^\top Z_{t+1,s} &= (W_{t+1}^{-1} (C_t^\alpha)^\top V_t^{-1})^\top \left(I_{rm} - \sum_{i=t+1}^{s-1} (C_i^\alpha)^\top V_i^{-1} C_i^\alpha W_s^{-1} \right) = \\ &= (W_{t+1}^{-1} (C_t^\alpha)^\top V_t^{-1})^\top \left(W_s - \sum_{i=t+1}^{s-1} (C_i^\alpha)^\top V_i^{-1} C_i^\alpha \right) W_s^{-1} = V_t^{-1} C_t^\alpha W_s^{-1}. \end{aligned}$$

Так как $G_{s,t} = A_s^\top A_{s+1}^\top \dots A_{t-1}^\top$, $\Phi_{t,s} = A_{t-1} A_{t-2} \dots A_s$, то выполняется (45). Подстановка (45) в (44) дает

$$\hat{\alpha}_{t+1} = W_{t+1}^{-1} \sum_{s=1}^t W_{s+1} W_{s+1}^+ (C_s^\alpha)^\top V_s^{-1} y_s. \quad (46)$$

Покажем, что

$$W_{s+1} W_{s+1}^+ (C_s^\alpha)^\top = (C_s^\alpha)^\top, \quad s = 1, \dots, t. \quad (47)$$

Используя скелетное разложения $\tilde{C}_{s+1} = L_{s+1} \Gamma_{s+1}$, получаем

$$W_{s+1} = \tilde{C}_{s+1} \tilde{C}_{s+1}^\top = L_{s+1}^\top \tilde{\Gamma}_{s+1} L_{s+1}, \quad W_{s+1}^+ = (L_{s+1}^\top)^+ \tilde{\Gamma}_{s+1}^{-1} L_{s+1}^+, \quad C_{s+1} = L_{s+1} \bar{\Gamma}_{s+1},$$

где $\tilde{\Gamma}_{s+1} = \Gamma_{s+1} \Gamma_{s+1}^\top$, $\bar{\Gamma}_{s+1}$ — некоторая матрица. Поскольку L_{s+1} — матрица полного ранга по столбцам, то $L_{s+1}^+ = (L_{s+1}^\top L_{s+1})^{-1} L_{s+1}^\top$ и $W_{s+1} W_{s+1}^+ = L_{s+1} (L_{s+1}^\top L_{s+1})^{-1} L_{s+1}^\top$. Отсюда следует (47), подстановка которого в (46) дает $\hat{\alpha}_{t+1} = \alpha_t^{\min}$.

Замечание 2. Соотношения (40), (41) представляют собой новую, рекуррентную версию (диффузную) ELM-алгоритма.

4. ДВУХЭТАПНЫЕ АЛГОРИТМЫ ОБУЧЕНИЯ

Получим упрощенные ДАО, используя двухэтапную процедуру оценивания [19]. Фиксируя последовательно в (3) $\hat{\beta}_t$ и $\hat{\alpha}_t$ и используя РФК, находим:

$$\hat{\beta}_{t+1} = \hat{\beta}_t + K_t^\beta (y_t - f(z_t, \hat{\beta}_t, \hat{\alpha}_t)), \hat{\beta}_1 = m_\beta^T, \quad (48)$$

$$\hat{\alpha}_{t+1} = \hat{\alpha}_t + K_t^\alpha (y_t - f(z_t, \hat{\beta}_t, \hat{\alpha}_t)), \hat{\alpha}_1 = 0_{rm}, \quad (49)$$

$$K_t^\beta = S_t (C_t^\beta)^T N_t^{-1}, K_t^\alpha = P_t (C_t^\alpha)^T N_t^{-1}, \quad (50)$$

$$S_{t+1} = S_t - S_t (C_t^\beta)^T N_t^{-1} C_t^\beta S_t, N_{1t} = C_t^\beta S_t (C_t^\beta)^T + V_t, S_1 = S_\beta, \quad (51)$$

$$P_{t+1} = P_t - P_t (C_t^\alpha)^T N_t^{-1} C_t^\alpha P_t, N_{2t} = C_t^\alpha P_t (C_t^\alpha)^T + V_t, P_1 = \lambda I_{rm}. \quad (52)$$

Теорема 3. 1. При ограниченном N справедливы асимптотические представления

$$P_t = (I_r - W_t W_t^+) \lambda + W_t^+ + O(\lambda^{-1}), t = 2, \dots, N, \quad (53)$$

$$K_t^\alpha = K_t^{\alpha, dif} + O(\lambda^{-1}), C_t^\alpha \neq 0, t = 1, \dots, N, \lambda \rightarrow \infty, \quad (54)$$

где

$$W_{t+1} = W_t + (C_t^\alpha)^T V_t^{-1} C_t^\alpha, W_1 = 0_{mn}, \quad (55)$$

$$K_t^{\alpha, dif} = W_{t+1}^+ (C_t^\alpha)^T V_t^{-1}. \quad (56)$$

2. Если дополнительно $\sigma(x) \in C^2(R)$, то

$$\hat{\beta}_t = \hat{\beta}_t^{dif} + O(\lambda^{-1}), \hat{\alpha}_t = \hat{\alpha}_t^{dif} + O(\lambda^{-1}), t = 1, \dots, N, \lambda \rightarrow \infty, \quad (57)$$

где

$$\hat{\beta}_{t+1}^{dif} = \hat{\beta}_t^{dif} + K_t^\beta \delta_t, \hat{\alpha}_{t+1}^{dif} = \hat{\alpha}_t^{dif} + K_t^{\alpha, dif} \delta_t, \quad (58)$$

$$\delta_t = y_t - f(z_t, \hat{\beta}_t^{dif}, \hat{\alpha}_t^{dif}), \hat{\beta}_1^{dif} = \bar{m}^\beta, \hat{\alpha}_1^{dif} = 0_{rm}. \quad (59)$$

Доказательство. Доказательство того, что $\hat{x}_t = (\hat{\alpha}_t^T, \hat{\beta}_t^T)^T = O(1)$ при $\lambda \rightarrow \infty, t = 2, \dots, N$, проводится так же, как и в теореме 1, поэтому его опускаем. Поскольку

$$P_{t+1}^{-1} = P_t^{-1} + (C_t^\alpha)^T V_t^{-1} C_t^\alpha, P_1^{-1} = I_{rm} / \lambda, t = 1, \dots, N,$$

то из леммы 3, если положить в (27) $F(y_t, \hat{x}_t) = V_t^{-1/2} C_t^\alpha$, будет следовать (53).

Имеем

$$\begin{aligned} K_t^\alpha &= P_t (C_t^\alpha)^T N_{2t}^{-1} = P_{t+1} P_{t+1}^{-1} P_t (C_t^\alpha)^T N_{2t}^{-1} = \\ &= P_{t+1} (P_t^{-1} + (C_t^\alpha)^T V_t^{-1} C_t^\alpha) P_t (C_t^\alpha)^T N_{2t}^{-1} = P_{t+1} (C_t^\alpha)^T V_t^{-1}. \end{aligned} \quad (60)$$

В результате подстановки (53) в (60) и в силу леммы 2 получаем (54).

Пусть $e_t^\beta = \hat{\beta}_t - \hat{\beta}_t^{dif}, e_t^\alpha = \hat{\alpha}_t - \hat{\alpha}_t^{dif}$. Тогда

$$\begin{aligned} e_{t+1}^\beta &= e_t^\beta + (K_t^\beta (\hat{\beta}_t, \hat{\alpha}_t) - K_t^\beta (\hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha)) y_t - K_t^\beta (\hat{\beta}_t, \hat{\alpha}_t) f(z_t, \hat{\beta}_t, \hat{\alpha}_t) + \\ &+ K_t^\beta (\hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha) f(z_t, \hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha), \end{aligned} \quad (61)$$

$$\begin{aligned} e_{t+1}^\alpha &= e_t^\alpha + (K_t^\alpha (\hat{\beta}_t, \hat{\alpha}_t) - K_t^{\alpha, dif} (\hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha)) y_t - K_t^\alpha (\hat{\beta}_t, \hat{\alpha}_t) f(z_t, \hat{\beta}_t, \hat{\alpha}_t) + \\ &+ K_t^{\alpha, dif} (\hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha) f(z_t, \hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha), t = 1, \dots, N. \end{aligned} \quad (62)$$

Так как $e_1^\beta = 0, e_1^\alpha = 0, f(u_t, \hat{\beta}_1, \hat{\alpha}_1) = 0, K_1^\alpha = K_1^{\alpha, dif} + O(\lambda^{-1})$ (без ограничения общности полагаем, что $C_t^\alpha \neq 0$), то из (61), (62) следует $e_2^\beta = 0, e_2^\alpha = O(\lambda^{-1})$. $\lambda \rightarrow \infty$. Пусть t произвольно и $e_t = O(\lambda^{-1}), \lambda \rightarrow \infty$. Если $C_t^\alpha = 0$, то $e_{t+1} = e_t = O(\lambda^{-1})$. При $C_t^\alpha \neq 0$ используем оценку $K_t^\alpha = K_t^{\alpha, dif} + O(\lambda^{-1})$ и теорему о среднем

$$f(z_t, \hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha) = f(z_t, \hat{\beta}_t, \hat{\alpha}_t) + O(\lambda^{-1}),$$

$$K_t^{\alpha, dif}(\hat{\beta}_t - e_t^\beta, \hat{\alpha}_t - e_t^\alpha) = K_t^{\alpha, dif}(\hat{\beta}_t, \hat{\alpha}_t) + O(\lambda^{-1}), \lambda \rightarrow \infty.$$

Подстановка этих выражений в (61), (62) дает $e_{t+1} = O(\lambda^{-1})$, $\lambda \rightarrow \infty$. Отсюда в силу ограниченности множества $t=1, \dots, N$ следуют оценки (57).

Замечание 3. Теорема 3 остается справедливой для двухэтапного ДАО в пакетном режиме, соотношения для которого нетрудно выписать, используя (11).

5. АСИМПТОТИЧЕСКИ ТОЧНЫЕ ИЗМЕРЕНИЯ

Рассмотрим поведение ДАО при малых шумах измерения, матрица интенсивности которых V_t входит в алгоритмы через обратную матрицу $N_{1t}^{-1} = (C_t^\beta S_t (C_t^\beta)^\top + V_t)^{-1}$. С учетом этого важно понимать, как будут вести себя ДАО при $V_t \rightarrow 0$.

Теорема 4. При ограниченном N и $V_t = I_{mn} / \lambda$, $\lambda \rightarrow \infty$ и $S_1 > 0$

$$\tilde{K}_t = \lim_{\lambda \rightarrow \infty} K_t^{dif} = \begin{pmatrix} Q_{t+1} (C_t^\beta)^\top \\ 0_{r \times rm} \end{pmatrix} + \begin{pmatrix} G_{t+1} \\ I_{rm} \end{pmatrix} L_{t+1}^+ \begin{pmatrix} G_t \\ I_{rm} \end{pmatrix} \begin{pmatrix} (C_t^\beta)^\top \\ (C_t^\alpha)^\top \end{pmatrix} \bar{H}_t, \quad (63)$$

где

$$Q_{t+1} = Q_t + (C_t^\beta)^\top C_t^\beta, \quad Q_1 = 0_{mn \times mn}, \quad (64)$$

$$G_{t+1} = (I_{mn} - Q_{t+1}^+ (C_t^\beta)^\top C_t^\beta) G_t - Q_{t+1} C_t^\alpha, \quad G_1 = 0_{mn \times r}, \quad (65)$$

$$L_{t+1} = L_t + (C_t^\beta G_t + C_t^\alpha)^\top \bar{H}_t (C_t^\beta G_t + C_t^\alpha), \quad L_1 = 0_{rm \times rm}, \quad (66)$$

$$\bar{H}_t = (I_r - \tilde{H}_t \tilde{H}_t^+) (I_r + C_t^\beta Q_t^+ (C_t^\beta)^\top)^{-1}, \quad (67)$$

$$\tilde{H}_t = (I_r + C_t^\beta Q_t^+ (C_t^\beta)^\top)^{-1} C_t^\beta H_t (C_t^\beta)^\top,$$

$$H_t = S_1 (I_{mn} - Q_t Q_t^+). \quad (68)$$

Доказательство. Так как

$$S_t^{-1} = \left(\frac{S_1^{-1}}{\lambda} + \sum_{s=a}^{t-1} (C_s^\beta)^\top C_s^\beta \right) \lambda, \quad t = 2, \dots, N,$$

то из леммы 3, если положить в (27) $\Omega_1 = S_1^{-1}$, $F(y_t, \hat{x}_t) = C_t^\beta$, будет следовать

$$S_t = S_1 (I_{mn} - Q_t Q_t^+) + Q_t^+ / \lambda + O(1/\lambda^2), \quad t = 2, \dots, N. \quad (69)$$

Поскольку $K_t^\beta = S_t (C_t^\beta)^\top N_{1t}^{-1} = S_{t+1} (C_t^\beta)^\top V_t^{-1}$, то из (69) и леммы 2 имеем

$$\lim_{\lambda \rightarrow \infty} S_t (C_t^\beta)^\top N_{1t}^{-1} = Q_{t+1}^+ (C_t^\beta)^\top, \quad (70)$$

но так как $R_{t+1} = (I_{mn} - K_t^\beta C_t^\beta) R_t - K_t^\beta C_t^\alpha$, $R_1 = 0_{mn}$, то отсюда вытекает (65).

Преобразуем N_{1t} с помощью (69):

$$\begin{aligned} N_{1t} &= C_t^\beta S_t (C_t^\beta)^\top + V_t = C_t^\beta H_t (C_t^\beta)^\top + (I_{mn} + C_t^\beta Q_t^+ (C_t^\beta)^\top) / \lambda + O(1/\lambda^2) = \\ &= (I_{mn} + C_t^\beta Q_t^+ (C_t^\beta)^\top) (\tilde{H}_t + I_{mn} / \lambda + O(1/\lambda^2)). \end{aligned}$$

Использование леммы 3 дает

$$N_{lt}^{-1} = (\tilde{H}_t + I_r / \lambda)^{-1} (I + C_t^\beta Q_t^+ (C_t^\beta)^T)^{-1} + O(1/\lambda^2) = \\ = \lambda(I_r - \tilde{H}_t \tilde{H}_t^+) (I + C_t^\beta Q_t^+ (C_t^\beta)^T)^{-1} + O(1) = \lambda \bar{H}_t + O(1).$$

Имеем

$$W_{t+1} = W_t + \lambda(C_t^\beta R_t + C_t^\alpha)^T (\bar{H}_t + O(1/\lambda))(C_t^\beta R_t + C_t^\alpha), \quad W_a = 0_{rm \times rm}.$$

Отсюда следует, что

$$\lim_{\lambda \rightarrow \infty} W_t / \lambda = \sum_{s=a}^{t-1} (C_s^\beta G_s + C_s^\alpha)^T \bar{H}_s (C_s^\beta G_s + C_s^\alpha), \\ \lim_{\lambda \rightarrow \infty} \begin{pmatrix} R_{t+1} \\ I_{rm} \end{pmatrix} W_{t+1}^+ \begin{pmatrix} R_t \\ I_{rm} \end{pmatrix}^T \begin{pmatrix} (C_t^\beta)^T \\ (C_t^\alpha)^T \end{pmatrix} N_{lt}^{-1} = \begin{pmatrix} G_{t+1} \\ I_{rm} \end{pmatrix} L_{t+1}^+ \begin{pmatrix} G_t \\ I_{rm} \end{pmatrix}^T \begin{pmatrix} (C_t^\beta)^T \\ (C_t^\alpha)^T \end{pmatrix} \bar{H}_t.$$

Теорема 5. При ограниченном N и $V_t = I_{mn} / \lambda$, $\lambda \rightarrow \infty$ и $S_1 > 0$ справедливо

$$\bar{K}_t^\beta = \lim_{\lambda \rightarrow \infty} K_t^\beta = Q_{t+1} (C_t^\beta)^T, \quad \bar{K}_t^\alpha = \lim_{\lambda \rightarrow \infty} K_t^{dif, \alpha} = L_{t+1} (C_t^\alpha)^T,$$

где

$$Q_{t+1} = Q_t + (C_t^\beta)^T C_t^\beta, \quad Q_1 = 0_{mn \times mn}, \quad L_{t+1} = L_t + (C_t^\alpha)^T \bar{H}_t C_t^\alpha, \quad L_1 = 0_{rm \times rm}.$$

Доказательство повторяет вывод формулы (70) в теореме 4.

6. ЧИСЛЕННЫЙ ПРИМЕР

Проиллюстрируем эффективность полученных в работе результатов на одном из стандартных примеров, используемых при тестировании алгоритмов [11]. Рассмотрим задачу аппроксимации функции

$$y(x) = \begin{cases} \sin(x)/x, & x \neq 0 \\ 1, & x = 0. \end{cases}$$

Обучающее и тестирующее множества (x_i, y_i) включают по 2000 точек, а значения x_i равномерно распределены на интервале $[-10, 10]$. К обучающим значениям y_i добавляется равномерно распределенный на интервале $[-0.2, 0.2]$ шум, а тестирующее множество предполагается незашумленным. Используется сигмоидная ФА и сравнивается точность аппроксимации ELM и ДАО с последовательной обработкой. В качестве меры точности используется оценка 90-й перцентили ошибки аппроксимации по 500 выборкам. Веса и смещения скрытого слоя предполагаются равномерно распределенными случайными величинами на интервалах $[-1, 1]$, $[0, 1]$ соответственно, $V_t = 0.16/12$. Результаты моделирования приведены в табл. 1, для алгоритмов, следующих из теорем 1, 3 (ДАО1 и ДАО2 соответственно). Видно, что ДАО превосходят по точности ELM как на обучающем, так и на тестирующем множествах. ДАО1 и ДАО2 дают практически неразличимые результаты по точности, но при этом ДАО2 превосходит ДАО1 по быстродействию на 32 %. На рис. 1 приведены графики аппроксимирующей зависимости 1, 2 (ELM и ДАО1 соответственно) для пяти нейронов скрытого слоя и одного использованного обучающего множества.

Таблица 1

Количество нейронов скрытого слоя	Точность алгоритмов					
	при обучении			при тестировании		
	ELM	ДАО1	ДАО2	ELM	ДАО1	ДАО2
4	0.315	0.169	0.174	0.295	0.125	0.13
5	0.272	0.159	0.159	0.246	0.108	0.11
6	0.205	0.148	0.147	0.17	0.092	0.092

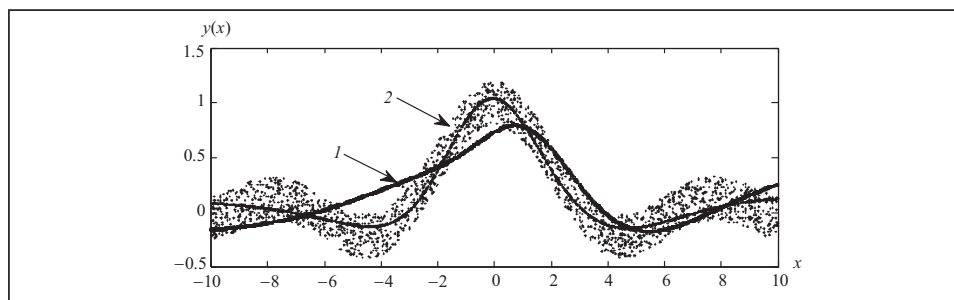


Рис. 1

ЗАКЛЮЧЕНИЕ

В настоящей работе рассмотрена задача обучения НСПР. Для ее решения предложены алгоритмы, основанные на асимптотическом анализе поведения РФК и сепарабельной структуре сети. Получение условий сходимости разработанных алгоритмов — одно из направлений дальнейших исследований.

СПИСОК ЛИТЕРАТУРЫ

1. Palit A., Popovic D. Computational intelligence in time series forecasting: Theory and engineering applications. — New York: Springer, 2005. — 372 p.
2. Devroye L., Györfi L., Lugosi G. Probabilistic theory of pattern recognition. — New York: Springer, 1996. — 636 p.
3. Neural systems for control / D.L Elliott, O.M. Omidvar (Eds). — Boston: Academ. Press, 1997. — 362 p.
4. Rumelhart D.E, Hinton G.E., Williams R.J. Learning internal representation by error propagation // Parallel Distributed Proces. — Cambridge, MA: MIT, 1986. — P. 318–362.
5. Battiti R. First and second order methods for learning: Between steepest descent and Newton's method // Neural Comput. — 1992. — 4, N 2. — P. 141–166.
6. Hagan M.T., Menhaj M. Training multilayer networks with the Marquardt algorithm // IEEE Transact. on Neural Networks. — 1994. — 5, N 6. — P. 989–993.
7. Singhal S., Wu I. Training multilayer perceptrons with the extended Kalman filter // Advances in Neural Inform. Proces. Systems. — 1989. — 1. — P. 133–140.
8. Iiguni Y., Sakai H., Tokumaru H. A real time learning. Algorithm for a multilayered neural network based on the extended. Kalman filter // Signal Proces., IEEE Transact. — 1992. — 40, N 4. — P. 959–966.
9. Bertsekas D.P. Incremental least squares methods and the extended Kalman filter // SIAM J. Optim. — 1996. — 6. — P. 807–822.
10. Pereyra V., Scherer G., Wong F. Variable projections neural network training // Mathematics and Comput. in Simulat. — 2006. — 73. — P. 231–243.
11. Huang G.B., Zhu Q.Y., Siew C.K. Extreme learning machine: Theory and applications // Neurocomputing. — 2006. — 70, N 1–3. — P. 489–501.
12. Huang G.B., D.H. Wang D.H., Lan Y. Extreme learning machines: A survey // Intern. Journ. of Machine Learning and Cybernetics. — 2011. — 2, N 2. — P. 107–122.
13. Cybenko G. Aproximation by superpositions of a sigmoidal function // Mathematics of Control, Signals and Systems. — 1989. — 2. — P. 303–314.
14. Golub G.H., Pereyra V. Separable nonlinear least squares: The variable projection method and its applications // Inverse Problems. — 2003. — 19, N 2. — P. 1–26.
15. Ansley C.F., Kohn R. Estimation, filtering and smoothing in state space models with incompletely specified initial conditions // Ann. Statist. — 1985. — N 13. — P. 1286–1316.
16. Anderson B.D.O., Moore J.B. Optimal filtering. — New York: Prentice-Hall, 1979. — 19. — 368 p.
17. Wimmer H.R. Stabilizing and unmixed solutions of the discrete time algebraic Riccati equation // Proc. Workshop on the Riccati equation in Control, Systems, and Sygnals. — 1989. — Italy. — P. 95–98.
18. Алберт А. Регрессия, псевдоинверсия и рекуррентное оценивание. — М.: Наука, 1977. — 224 с.
19. Ljung L. System identification. Theory for the user. — New York: Prentice-Hall, 1999. — 603 p.

Поступила 06.06.2012