**ORIGINAL ARTICLE**

# Prediction of component concentrations in sodium aluminate liquor using stochastic configuration networks

**Wei Wang[1,2]** · **Dianhui Wang[2,3]**

## Abstract

Online measuring of component concentrations in sodium aluminate liquor is essential and important to Bayer alumina production process. They are the basis of closed-loop control and optimization and affect the final product quality. There are three main components in sodium aluminate liquor, termed caustic hydroxide, alumina and sodium carbonate (their concentrations are represented by $c_K$, $c_A$ and $c_C$, respectively). They are obtained off-line by titration analysis and suffered from larger time delays. To solve this problem, a hybrid model for $c_K$ and $c_A$ is proposed by combining a mechanism model and a stochastic configuration network (SCN) compensation model. An SCN-based model for $c_C$ is also proposed using the estimated values of $c_K$ and $c_A$ from the hybrid model. A real-world application conducted in Henan Branch of China Aluminum Co. Ltd demonstrates the effectiveness of the proposed modelling techniques. Experimental results show that our proposed method performs favourably in terms of the prediction accuracy, compared against the regress model, BP neural networks, RBF neural networks and random vector functional link model.

**Keywords** Stochastic configuration networks · Industrial data modelling · Component concentrations · Sodium aluminate liquor

## 1 Introduction

Sodium aluminate liquor is an intermediate product throughout the whole process of Bayer alumina production. The main component concentrations of sodium aluminate liquor $c_K$, $c_A$ and $c_C$ are the foundation for indicator control in the procedure of original ore pulp preparation, digestion, decomposition and evaporation. Due to the nonlinearity characteristics of sodium aluminate liquor, such as easy precipitation, large viscosity, high concentration and strong corrosiveness, online measuring of component concentrations is quite difficult and usually they are obtained through artificial sampling and laboratory titration analysis, which is not only complicated and costly, but also suffers from larger time delays as well. Therefore, online estimate of component concentrations in sodium aluminate liquor is significant to implement process optimization and control of the alumina industry.

The existing methods of estimating $c_K$, $c_A$ and $c_C$ can be classified into two categories. One is based on the principle of reagent titration [15, 16], including photometric titration, potentiometric titration, thermometric titration, microtitration, reagent automatic titration analysis, flow injection analysis and so on. However, most of them belong to off-line measurement, and the instrument pipeline is relatively thin and easy to scar and block up, resulting in degraded estimate accuracy. The automatic titration developed by American Matocha is highly automated with good precision for the chemical reaction process. But, such an instrument system has some limitations such as complex structure, low reliability, strict environmental requirements and high maintenance cost, which stop making a wide use in domestic alumina industry. Another category mainly includes the temperature–conductivity method and the

✉ Dianhui Wang
dh.wang@latrobe.edu.au

1   College of Information Engineering, Dalian Ocean University, Dalian 116023, China

2   State Key Laboratory of Synthetical Automation for Process Industry, Northeastern University, Shenyang 110819, China

3   Department of Computer Science and Information Technology, La Trobe University, Melbourne, VIC 3086, Australia

conductivity–density–ultrasonic method [2]. Among the methods of separating the reagents, the temperature–conductivity measurement is relatively simple with low cost, which is the first choice for online measurement of sodium aluminate liquor component concentrations in the alumina industry of China. In the literature, the design of sensor model of $c_K$, $c_A$ and $c_C$ using temperature and conductivity is mainly based on the orthogonal test and the least squares regression techniques [3]. From data modelling perspectives [7, 8], these methods do not take advantages of making full use of mechanistic knowledge and applying the data-driven error compensation model or model correction mechanism for domain problem-solving.

To address the problems mentioned above, we proposed a method of combining a mechanism model with an error compensation model based on neural networks for $c_K$ and $c_A$ in [22]. However, the learning process is time-consuming and easy to fall into a local minimum. Besides, it is quite challenging to properly set the initial weights and biases and determine the number of hidden nodes which are closely associated with both the learning and generalization performance. To resolve these problems, randomized methods for neural networks can be applied for developing fast learner models [14]. Randomized learning algorithms have much less computational cost through random assignment on the input weights and biases and evaluate the output weights by least squares methods. Compared to randomized RBF networks [1] and RVFL networks [6, 11], stochastic configuration networks (SCNs) share merit to ensure the universal approximation property of built randomized learner models. The essential and innovative contribution of the SCN framework is the way of assigning the random parameters with an inequality constraint and adaptively selecting the scope of the random parameters [18]. So an SCN-based learner model is employed in this paper to compensate the unknown parts which are not included in the mechanism model of $c_K$ and $c_A$. Another component concentration $c_C$ is related to $c_K$ and $c_A$, but it is difficult to establish a mechanism model. We propose a data-driven method to estimate the $c_C$ using another SCN model with the $c_K$ and $c_A$ as part of the inputs.

The remainder of this paper is organized as follows: Sect. 2 describes the production process of Bayer alumina and formulates some relevant problems in the component concentration measurement. Section 3 introduces the SCN framework to support the following hybrid modelling. Section 4 details the mechanism model and the learning-based error compensation model for $c_K$ and $c_A$, and data-driven model based on SCN for $c_C$. Section 5 reports our experimental results in Henan Branch of China Aluminum Co. Ltd, and Sect. 6 concludes this paper.

# 2 Problem description of component concentration measurement in sodium aluminate liquor

## 2.1 Introduction of Bayer alumina production

Bayer alumina production is one of the most popular methods in alkaline alumina process, and the principle of Bayer alumina production is shown in Fig. 1: The bauxite, lime and recycled liquor are mixed by a certain proportion and grinded into original ore pulp. Under the conditions of high temperature and high pressure, the caustic solution is used to dissolve the alumina in the bauxite to produce sodium aluminate liquor. The prepared sodium aluminate liquor is further added with aluminium hydroxide as a seed to be decomposed under cooling and stirring conditions to obtain aluminium hydroxide, and the aluminium hydroxide is roasted to obtain alumina. The remaining mother liquor is evaporated and then used to dissolve a new batch of bauxite. The impurities, such as silica, become into red mud and will be discharged after washing or used in the sintering process.

Throughout the Bayer alumina production process, sodium aluminate liquor exists almost throughout the circulation. Many procedures require real-time detection of sodium aluminate liquor component concentrations. For example, they play an important role in the original ore pulp preparation procedure for the control of liquid–solid ratio, in the high-pressure digestion procedure for the control of dissolution rate, in the seed decomposition
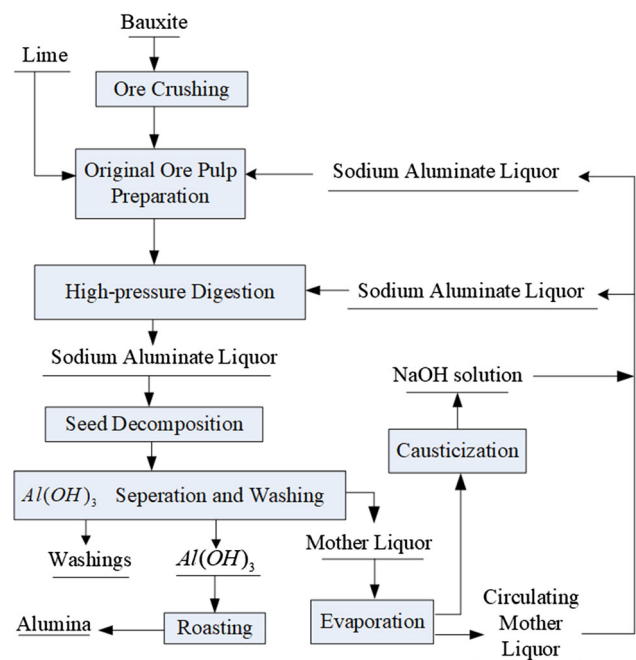


**Fig. 1** Flowchart of Bayer alumina production

procedure for the control of decomposition rate, in the evaporation procedure for the control of alumina/caustic ratio and so on. They are also crucial for inhibiting red mud expansion, and improving the sedimentation performance of red mud. Thus, real-time detection of the component concentrations in sodium aluminate liquor for each procedure has great significance for integrated automation. It will reduce the workers intensity and save the laboratory cost, and realize a stable and high-quality production in Bayer alumina. As a preliminary study, the first important step in the whole process which called original ore pulp preparation is chosen for our experiment.

## 2.2 Component concentration measurement in original ore pulp preparation

The original ore pulp preparation process uses a certain proportion of bauxite, lime and sodium aluminate liquor to prepare the original ore pulp with the chemical composition and physical properties meeting the requirements of the high-pressure digestion process. First, bauxite is crushed to meet the requirements of the particle size, and the crushed aluminium ore and lime are fed from a feed belt to a silo, which is then fed into a lattice mill using a plate feeder, and the amount of aluminium ore and lime fed is adjusted by controlling the speed of the plate feeder. On the other hand, the sodium aluminate liquor (including NaOH solution and the circulating mother liquor returned by the process of evaporation or seed decomposition, etc.) enters the mother liquor tank and is sent to the lattice mill and the classifier by the pump which is controlled by the frequency converter to adjust the flow rate of the circulating mother liquor. After the aluminium ore, lime and circulating mother liquor enter the lattice mill, they are mixed and finely ground to form a mixed pulp with certain fineness, proportion and uniformity. The ground original ore pulp is sent to a classifier at the same time as a certain amount of circulating mother liquor, and the mixed particles of different sizes, different shapes and different specific gravities are classified. After the classification process, the granules that meet the requirements are sent to the buffer tank for the high-pressure digestion process, and the undesired coarse particles are sent back to the lattice mill for re-grinding.

The liquid–solid ratio is an important indicator of this process. If the component concentration measurement is delayed, the flow of the liquor added to the lattice mill will be not accurate, and it will be difficult to prepare a qualified original ore pulp. In addition, the source of the sodium aluminate liquor in this process is scattered and uncertain, and the concentration difference between them is large. It is difficult to measure the component concentration using commonly used variables such as pressure, flow, etc.

Therefore, it is a right choice to use soft measurement methods to solve this problem according to its physical and chemical properties.

## 2.3 Modelling principle for component concentrations in sodium aluminate liquor

According to the factors that affect the electrolyte conductivity and the composition of sodium aluminate liquor, the conductivity, temperature and caustic hydroxide concentration $c_K$, alumina concentration $c_A$ and sodium carbonate concentration $c_C$ have the following relationship [3]:

$$d = f(T, c_K, c_A, c_C), \tag{1}$$

where $f(.)$ is a nonlinear function, $d$ is the conductivity of sodium aluminate liquor, $T$ is the temperature. From this relationship, it can be concluded that if we observe the changes of temperature and conductivity, we will receive the changes of component concentrations in sodium aluminate liquor.

It is also known that the relationship between temperature and conductivity is an approximate linear relationship [3], as shown in Fig. 2.

It can be seen that the conductivity of sodium aluminate liquor increases with the temperature rising, and these lines are corresponding to different component concentrations of $c_K$, $c_A$ and $c_C$. The slope and intercept of the approximate straight line vary with different component concentrations. Combination with (1), we can know that a certain proportion of sodium aluminate liquor (viz. $c_K$, $c_A$, $c_C$ are fixed), its temperature and conductivity have such relationship:
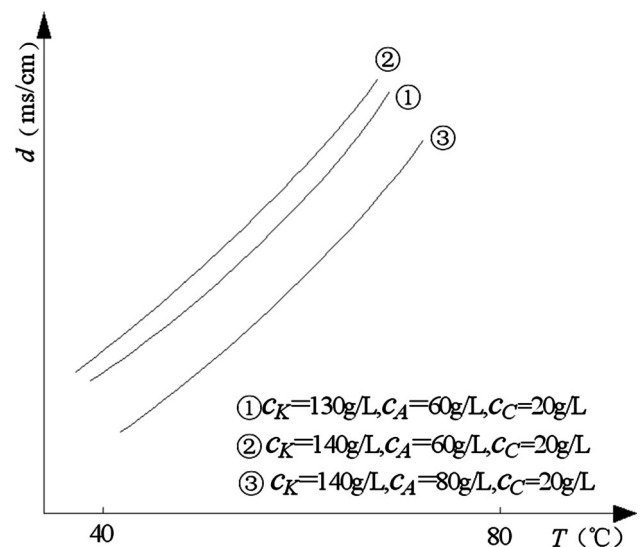


Fig. 2 Relationship between $T$ and $d$

$$d = k(c_K, c_A, c_C)T + b(c_K, c_A, c_C), \qquad (2)$$

where $k$ is the slope, $b$ is the intercept of the straight line and they are nonlinear function of $c_K$, $c_A$, $c_C$. It is known that $c_C$ can be negligible when its value is below 40 g/l, and in industry field $c_C$ is around 30 g/l. Thus, Eq. (2) can be rewritten as

$$d = k(c_K, c_A)T + b(c_K, c_A). \qquad (3)$$

Besides, there is also a linear relation between $d$ and $c_A$ (when $T$, $c_K$, $c_C$ are fixed), as shown in Fig. 3. It can be seen that the conductivity decreases as the alumina concentration increases, and the higher the temperature, the faster the decrease. It is a linear function between the two variables.

Using least squares regression, we have the following equations

$$k = \left( \frac{\partial k}{\partial c_A} \right)_{c_k} c_A + k_0, \qquad (4)$$

$$b = \left( \frac{\partial b}{\partial c_A} \right)_{c_k} c_A + b_0. \qquad (5)$$

It has been proved that a quadratic least squares regression gives the best fit when regression these coefficients against $c_K$ [3].

$$k = (K_1 c_K^2 + K_2 c_K + K_3)c_A + (K_4 c_K^2 + K_5 c_K + K_6), \qquad (6)$$

$$b = (B_1 c_K^2 + B_2 c_K + B_3)c_A + (B_4 c_K^2 + B_5 c_K + B_6), \qquad (7)$$

where $K_1$–$K_6$ and $B_1$–$B_6$ are unknown coefficients.

It can be seen from the derivation of the above equations, the input–output relationship of our model can be established as follows:
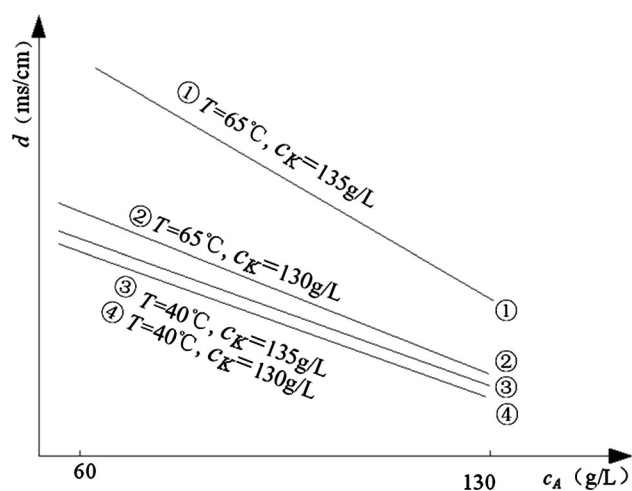


**Fig. 3** Relationship between $c_A$ and $d$

$$\begin{cases} \hat{y}_K = f_1[k(\boldsymbol{T}, \boldsymbol{d}), b(\boldsymbol{T}, \boldsymbol{d})], \\ \hat{y}_A = f_2[k(\boldsymbol{T}, \boldsymbol{d}), b(\boldsymbol{T}, \boldsymbol{d}), \hat{y}_K], \\ \hat{y}_c = f_3[k(\boldsymbol{T}, \boldsymbol{d}), b(\boldsymbol{T}, \boldsymbol{d}), \hat{y}_K, \hat{y}_A], \end{cases} \qquad (8)$$

where $\hat{y}_K$, $\hat{y}_A$, $\hat{y}_C$ is the model output of component concentrations, $\boldsymbol{T} = [T_1, T_2, T_3]$ and $\boldsymbol{d} = [d_1, d_2, d_3]$ are different temperatures ( heating, cooling and mixing) and the corresponding conductivities of sodium aluminate liquor.

Based on this principle, a hybrid online modelling strategy is proposed, as shown in Fig. 4. First, a device is needed for measuring the temperatures and conductivities of sodium aluminate liquor. After data sampling and preprocessing, a mechanism model for $c_K$ and $c_A$ is proposed, and the modelling error is compensated by an SCN model. $y_K$ and $y_A$ is the artificial laboratory value of $c_K$ and $c_A$, $y_{Km}$ and $y_{Am}$ is the output of mechanism model of $c_K$ and $c_A$, $e_K$ and $e_A$ is the error between them. $\hat{e}_K$ and $\hat{e}_A$ is the output of SCN-based compensation model. The proposed data-driven modelling method for $c_C$ is also based on SCN.

# 3 Introduction of SCN for industrial data modelling

This section provides a fast and effective machine learning methodology to support the proposed hybrid modelling method for problem-solving.

Hybrid modelling methods address a fusion technology of mechanism analysis and data-driven modelling techniques which is effective and widely used in many applications of industrial processes [4, 10, 12, 13, 17]. Due to the approximate linear relationship between the input variables and the output variables described in (2) and (3), modelling error occurs and some compensation strategies should be in place for improved modelling performance. In industrial hybrid modelling, there is mostly a combination of simple mechanism model and neural networks model. Unfortunately, the BP compensation model in [22] suffers from the sensitive setting of the hidden nodes number and learning rate, local minima and very slow convergence. To overcome this problem, random vector functional link networks (RVFL) are widely used and they perform reasonably good in terms of the learning performance and the predictability, compared to neural networks with optimization-based learning algorithms. Unfortunately, RVFL networks lack practical learning schemes that ensure the learning capability for a given dataset. The key issue presented in RVFL networks is the scope setting of the random weights and biases that may result in a failure of function approximation or data modelling [9]. Thanks to an advancement of randomized learning techniques, stochastic configuration networks (SCNs) were proposed in [18],
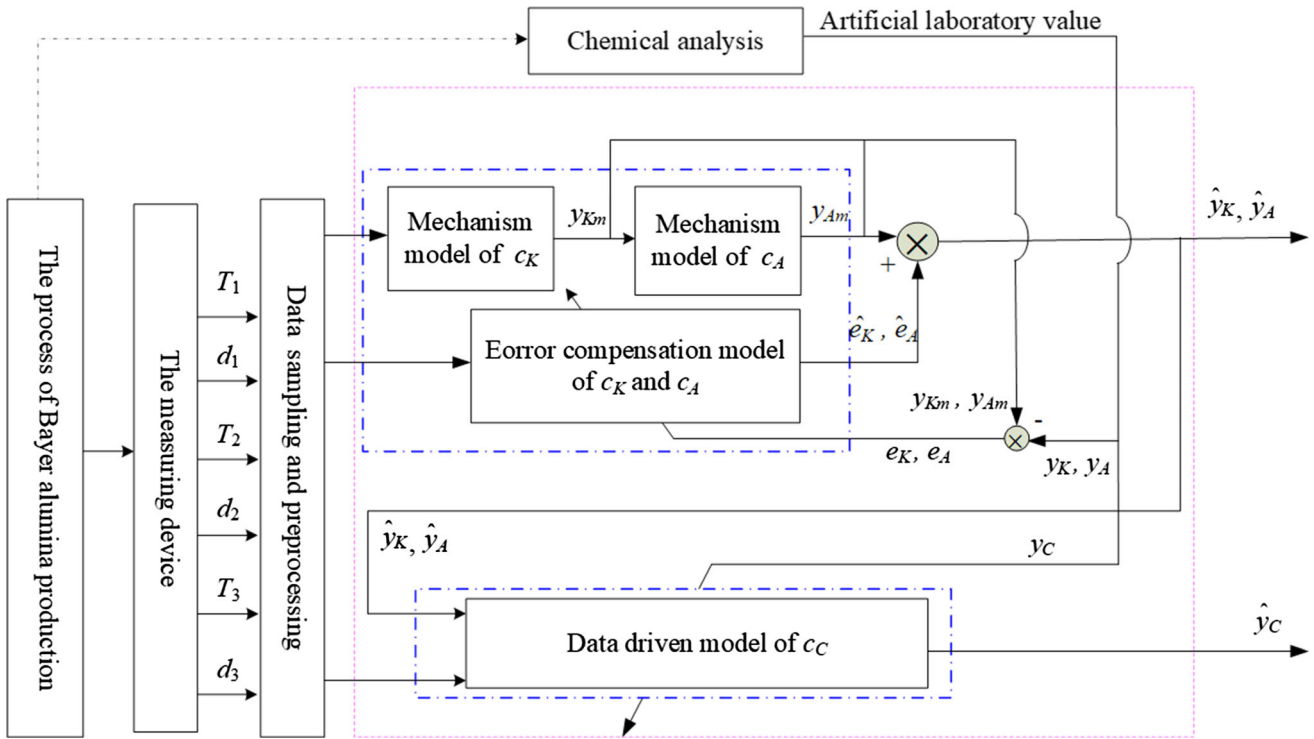
me

**Fig. 4** Structure of the modelling strategy

where a supervisory mechanism is firstly presented to guarantee the universal approximation property (i.e., learning capability). The following subsection briefly describes the SCN concept for fast industrial data modelling.

SCNs are a class of feedforward neural networks built by randomized algorithm [18]. The core contribution of the SCN framework lies in a supervisory mechanism for the random parameter assignment, and the learner model is incrementally built. This constructive approach for building SCNs guarantees the universal approximation property. For details of the SCN theory and learning algorithms, readers can refer to [18]. We outline the stochastic configuration algorithm as follows for our purpose in this study.

Given a training data set with $N$ sample pairs $\{(\mathbf{X}_p, \mathbf{Y}_p), p = 1, 2, \ldots, N\}$, where $\mathbf{X}_p = [x_1, x_2, \ldots, x_m]^{\mathrm{T}} \in R^m$ and $\mathbf{Y}_p = [y_1, y_2, \ldots, y_n]^{\mathrm{T}} \in R^n$. Let $X \in R^{N \times m}$ and $Y \in R^{N \times n}$ represent the input and output data matrix, respectively; $e_{L-1}(X) \in R^{N \times n}$ be the residual error matrix for the SCN model with $L - 1$ terms, where each column $e_{L-1,q}(X) = [e_{L-1,q}(\mathbf{X}_1), \ldots, e_{L1,q}(\mathbf{X}_N)]^{\mathrm{T}} \in R^N, q = 1, 2, \ldots, n$. Denote the output vector of the $L$-th hidden node $\phi_L$ for the input $X$ by

$$h_L(X) = [\phi_L(w_L^{\mathrm{T}}X_1 + b_L), \ldots, \phi_L(w_L^{\mathrm{T}}X_N + b_L)]^{\mathrm{T}}. \quad (9)$$

Thus, the hidden layer output matrix of the SCN model can be expressed as $H_L = [h_1, h_2, \ldots, h_L]$. Let

$$\xi_{L,q} = \frac{\left(e_{L-1,q}^{\mathrm{T}}(X) * h_L(X)\right)^2}{h_L^{\mathrm{T}}(X) * h_L(X)} - (1 - r_L)e_{L-1,q}(X), \quad (10)$$

$$q = 1, 2, \ldots, n.$$

With these notations, the SC Algorithm proposed in [18] can be summarized as follows:

*Step 1.* Set up learning parameters, including a set of scope $[-\lambda_i, \lambda_i]$, $i = 1, 2, \ldots, s$, where $0 < \lambda_1 < \lambda_2 < \cdots < \lambda_s$, and an increasing sequence $r_1 < r_2 < \cdots < r_t < 1$; also, set up two termination conditions, that is, either the maximum number of the hidden nodes $L_{\max}$ or the error tolerance $\tau$.

*Step 2.* Take random parameters $w_L$ and $b_L$ from adjustable interval $[-\lambda, \lambda]$ for Nc (a user specified integer) times, and check out the following inequalities with $r_i$, $i = 1, 2, \ldots, t$ (starting from $r_1$)

$$\xi_{L,q} \geq 0, q = 1, 2, \ldots, n. \quad (11)$$

If (11) holds, define the set of random parameters $w_L$ and $b_L$ such that $\xi_L = \sum_{q=1}^{n} \xi_{L,q}$ takes the maximum.

*Step 3.* Evaluate the output weight matrix $\beta$ by solving the following least means square problem:

$$\beta^* = \arg\min_{\beta} \|H_L\beta - Y\|_F^2 = H_L^+ Y, \tag{12}$$

where $H_L^+$ is the Moore–Penrose generalized inverse of the matrix $H_L$, and $\|.\|_F$ represents the Frobenius norm.

To speed up the procedure of building SCN models, we could add the top $n_B$ (an end-user specified integer) ranked (according to the values of $\xi_L$) candidate nodes as a batch in each incremental loop of the SC algorithm. Some detailed discussions can be found in [5].

# 4 Hybrid modelling method for component concentrations

## 4.1 Hybrid modelling for $c_K$ and $c_A$

Simplify (6) and (7), $y_{Km}$ can be evaluated by the following equation:

$$\begin{aligned}
y_{Km}^4 &+ m_1 y_{Km}^3 + (m_2 k + m_3 b + m_4) y_{Km}^2 \\
&+ (m_5 k + m_6 b + m_7) y_{Km} \\
&+ (m_8 k + m_9 b + m_1 0) = 0,
\end{aligned} \tag{13}$$

and $y_{Am}$ can be evaluated by

$$y_{Am} = \frac{(k - b) - (n_4 y_{Km}^2 + n_5 y_{Km} + n_6)}{n_1 y_{Km}^2 + n_2 y_{Km} + n_3}, \tag{14}$$

where $m_1 - m_{10}$, $n_1 - n_6$ are the coefficients to be determined. The concrete steps for mechanism model of $c_K$ and $c_A$ are as follows:

*Step 1*: variable conversion

As described in Fig. 2, we use the approximate linear relationship between temperature and conductivity and convert the variables $T$ and $d$ into $k$ and $b$ for the next step, that is,

$$\begin{cases} d_1 = kT_1 + b, \\ d_2 = kT_2 + b, \\ d_3 = kT_3 + b. \end{cases}$$

Then, by using the least squares method, we can easily get

$$\begin{bmatrix} k \\ b \end{bmatrix} = \left( \begin{bmatrix} 1 & 1 & 1 \\ T_1 & T_2 & T_3 \end{bmatrix} \begin{bmatrix} 1 & T_1 \\ 1 & T_2 \\ 1 & T_3 \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & 1 \\ T_1 & T_2 & T_3 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix}.$$

*Step 2*: orthogonal experimental design for model parameters determination

The raw materials used in the orthogonal experiment are composed of evaporated liquor, pure sodium hydroxide and pure sodium carbonate. The composition and concentration of the evaporated liquor are shown in Table 1. The variation range of the orthogonal experimental factors is shown

**Table 1** Composition and concentration of evaporated liquor

| Composition | NaOH | Al(OH)$_3$ | Na$_2$CO$_3$ |
|---|---|---|---|
| Concentration (g/l) | 272.7 | 137.7 | 26.55 |

**Table 2** Variation range of orthogonal experimental factors

| Factors | $c_K$ (g/l) | $c_A$ (g/l) | $c_C$ (g/l) | T (°C) |
|---|---|---|---|---|
| Upper bound | 230 | 120 | 36 | 96 |
| Lower bound | 190 | 80 | 24 | 70 |

in Table 2. The factor and level coding is shown in Table 3.

According to the orthogonal table of four factors by quadratic regression, 15 kinds of solution need to be configured. The temperature and conductivity of these solutions for each group are measured by continuous heating, and the experimental equipment and measuring process are shown in Fig. 5.

The data of this orthogonal experiment are collected, and the parameters of the mechanism model are determined by these data. The algorithm is as follows:

(a) Using the orthogonal experimental data $c_K$, $k$ and $b$, the undetermined coefficient $m_1 - m_{10}$ in (13) is estimated by least square regression.

(b) Substitute orthogonal experimental data $k$ and $b$, solvie the equation with four order about $c_K$, and get the calculated value $y_{Km}$ of the mechanism model.

(c) Using the orthogonal experimental data $k$, $b$ and the calculated value of $y_{Km}$, the undetermined coefficients $n_1 - n_6$ in (14) are regressed to obtain $c_A$.

(d) By using the established approximate mechanism model with known coefficients, both the calculated values $y_{Km}$ and $y_{Am}$ can be obtained by new field data.

The inputs of compensation model based on SCNs are heating temperature $T_1$ and conductivity $d_1$, cooling temperature $T_2$ and conductivity $d_2$, mixing temperature $T_3$ and conductivity $d_3$, mechanism model outputs $y_{Km}$ and $y_{Am}$, and the neural network outputs are the errors $e_K$ and $e_A$ between the mechanism calculation values $y_{Km}$, $y_{Am}$ and real laboratory analysis values $y_K$ and $y_A$. The structure of this network is shown in Fig. 6.

The predicted results $\hat{e}_K$ and $\hat{e}_A$ of neural network are used to compensate the error of mechanism model. The output of this hybrid model can be described as follows:

$$\hat{y}_K = y_{Km} + \hat{e}_K, \tag{15}$$

and

**Table 3** Factor and level coding of orthogonal experiment ($r = 1.414$)

| Variables | $n_1$ (caustic hydroxide) | $n_2$ (alumina) | $n_3$ (sodium carbonate) | $n_4$ (temperature) |
|---|---|---|---|---|
| Factor | $c_K$ (g/l) | $c_A$ (g/l) | $c_C$ (g/l) | $T$ (°C) |
| Benchmark level | 210 ($Z_1$) | 100 ($Z_2$) | 30 ($Z_3$) | 83 ($Z_4$) |
| Change distance | 14.14 ($\Delta_1$) | 14.14 ($\Delta_2$) | 4.24 ($\Delta_3$) | 9.19 ($\Delta_4$) |
| Upper level (+1) | 224.14 | 114.14 | 34.24 | 92.19 |
| Lower level (−1) | 195.86 | 85.86 | 25.76 | 73.81 |
| Upper asterisk arm | 230 | 120 | 36 | 96 |
| Lower asterisk arm | 190 | 80 | 24 | 70 |



**Fig. 5** Process of orthogonal experiment



**Fig. 7** Structure of data-driven model for $c_C$



**Fig. 6** Structure of error compensation model based on SCNs

$$\hat{y}_A = y_{Am} + \hat{e}_A, \tag{16}$$

where $\hat{y}_K$ and $\hat{y}_A$ are the final model output of $c_K$ and $c_A$.

## 4.2 SCN-based data-driven model for $c_C$

The proposed SCN-based data-driven model structure for $c_C$ is shown in Fig. 7.

The input variables are $T_1$, $T_2$, $T_3$, $d_1$, $d_2$, $d_3$ and $\hat{y}_K$, $\hat{y}_A$. The neural network output is artificial laboratory value $y_C$. The SC algorithm for $c_C$ model is almost the same as the compensation model introduced above. The difference is
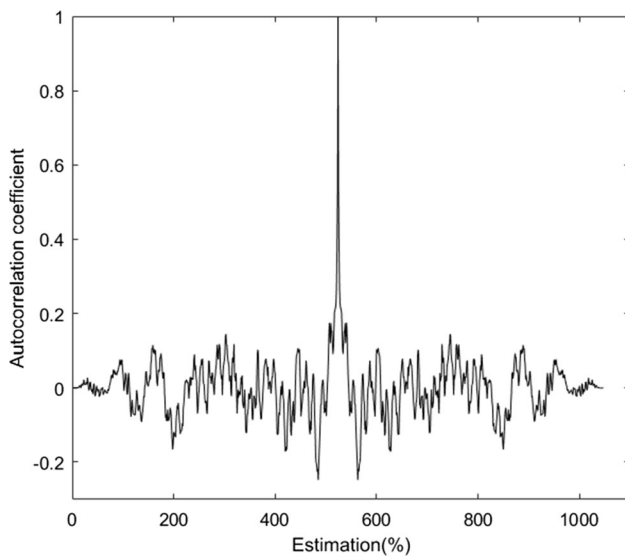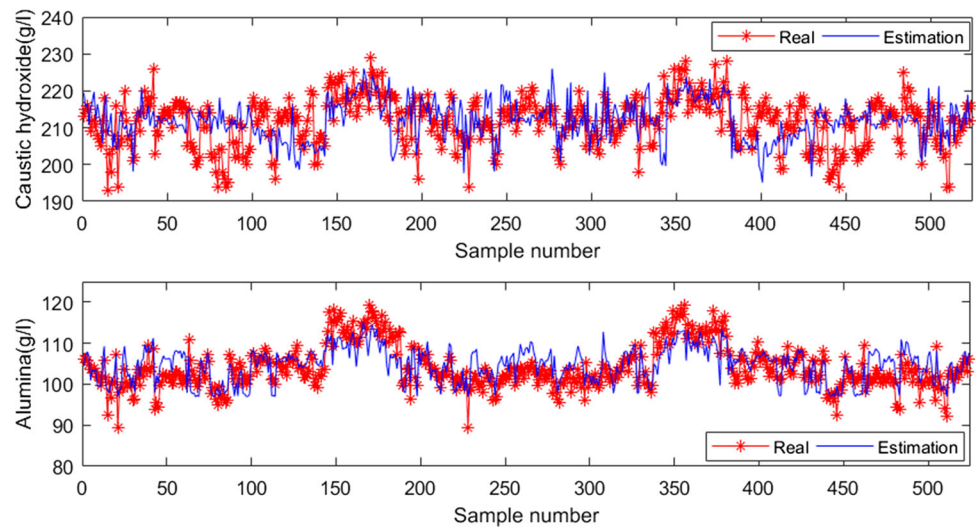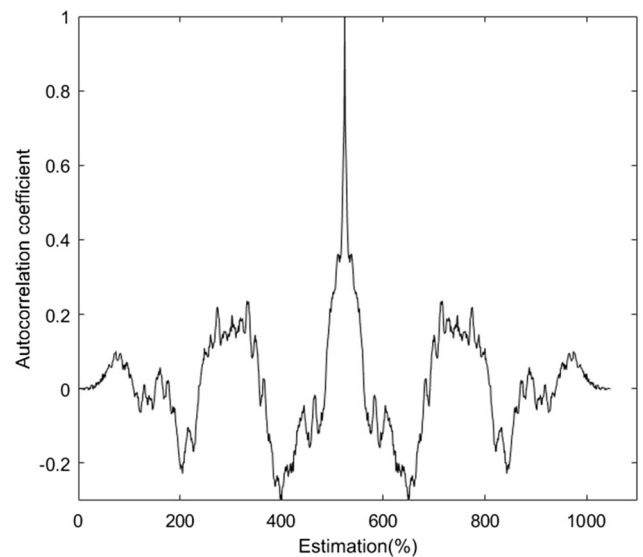


**Fig. 8** Hardware of measurement device

the training data set with $N$ sample pairs $\{(\boldsymbol{x}_{Cn}, \boldsymbol{y}_{Cn}), n = 1, 2, \ldots, N\}$ where $\boldsymbol{x}_{Cn} \in R^8, \boldsymbol{y}_{Cn} \in R^1$, and the input and output data matrix is $X_C \in R^{N \times 8}, Y_C \in R^{N \times 1}$, respectively.

**Table 4** Partial of the original data

| Variables | 1 record | 2 records | 3 records | ... | 523 records | 524 records |
|---|---|---|---|---|---|---|
| $T_1$ (°C) | 86.18 | 90.21 | 88.71 | ... | 92.86 | 92.94 |
| $D_1$ (ms/cm) | 577.73 | 593.59 | 578.44 | ... | 661.09 | 660.70 |
| $T_2$ (°C) | 72.29 | 72.92 | 66.82 | ... | 75.28 | 74.25 |
| $D_2$ (ms/cm) | 455.95 | 456.09 | 400.86 | ... | 493.91 | 486.33 |
| $T_3$ (°C) | 77.71 | 80.87 | 79.82 | ... | 78.95 | 78.95 |
| $D_3$ (ms/cm) | 482.11 | 507.34 | 489.77 | ... | 526.02 | 521.41 |
| $c_K$ (g/l) | 213.00 | 216.00 | 215.00 | ... | 218.00 | 218.00 |
| $c_A$ (g/l) | 106.24 | 106.24 | 106.23 | ... | 103.94 | 104.27 |
| $c_C$ (g/l) | 30.60 | 30.60 | 29.60 | ... | 29.00 | 29.80 |



**Fig. 9** Results of mechanism model for $c_K$ and $c_A$



**Fig. 10** Autocorrelation function of mechanism modelling error for $c_K$



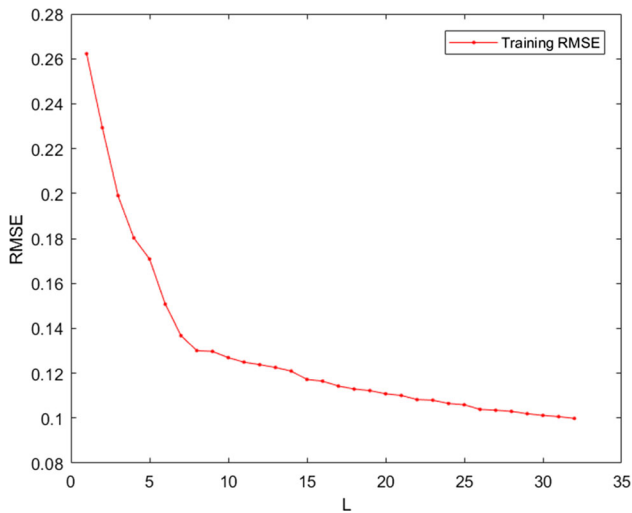**Fig. 11** Autocorrelation function of mechanism modelling error for $c_A$

**Fig. 12** RMSE of training process for SCN

# 5 Experiments

## 5.1 System setup

A measurement device is designed to implement the proposed modelling algorithm. As shown in Fig. 8, the hardware system can be divided into three main parts: sampling equipment, instrument control system, IPC (Industrial Personal Computer) and peripheral equipment. The sampling device is used to obtain sodium aluminate solution and collect temperature and conductivity data online. The device consists of pipes, solution tanks, heater, coolers, solenoid valves, manual valves, etc. They are made of stainless steel and are equipped with three temperature and conductivity probes in the solution tanks. PLC controller is used for switch control of heater, cooler, solenoid valves and other equipment. IPC and peripherals are used for parameter monitoring, model calculation and report printing.

The working principle of this device is as follows: The pipe inlet introduces the sodium aluminate liquor sample into the system through the main valve and adjusts the optimum flow rate by control the valve opening. After the liquor is heated by the heater, a portion of the liquor enters the first tank through the heating pipe, and a portion of the liquor is cooled by the cooler and passed through the cooling pipe to the second tank. The two liquors are mixed in the mixing tank and flowed into the third tank. Heating temperature and conductivity, cooling temperature and conductivity, mixing temperature and conductivity were measured during this process. The liquor passes through three pipes and is eventually returned from the device outlet to the on-site solution pipe.

In order to ensure the quality of modelling data, the required heating temperature range is about 85–95 °C, the cooling temperature range is about 65–75 °C, and the mixing temperature range is about 75–85 °C.

## 5.2 Data sampling and parameter selection in mechanism models

The data sampling period is set as 5 s in the measurement device; after the median filtering, it is about 65 s to store a set of data. Manual sampling interval during our

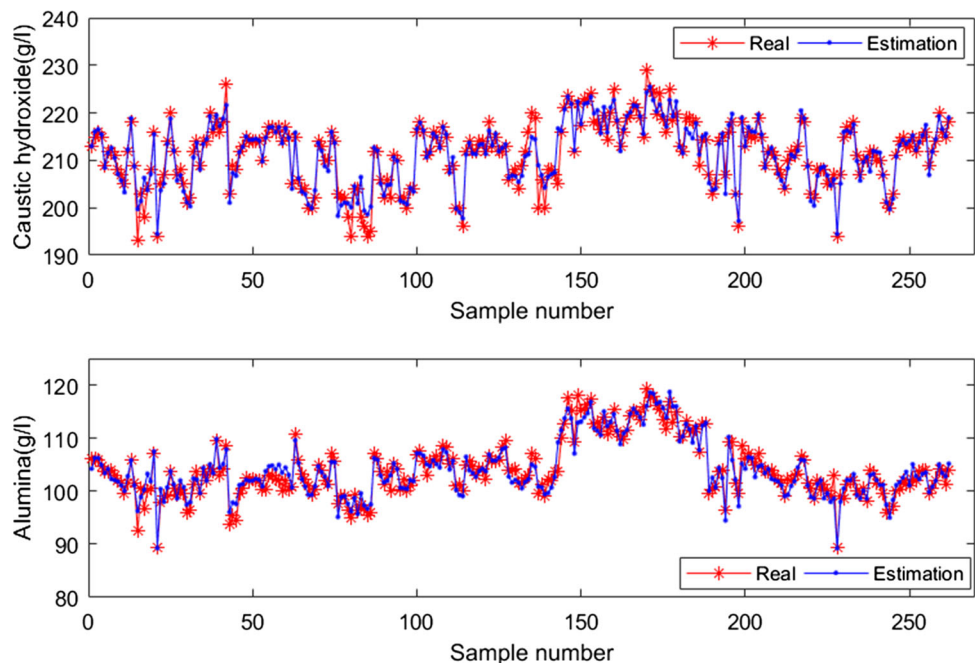**Fig. 13** Training results of $c_K$ and $c_A$ after SCN-based compensation

**Fig. 14** Test results of $c_K$ and $c_A$ before and after SCN-based compensation
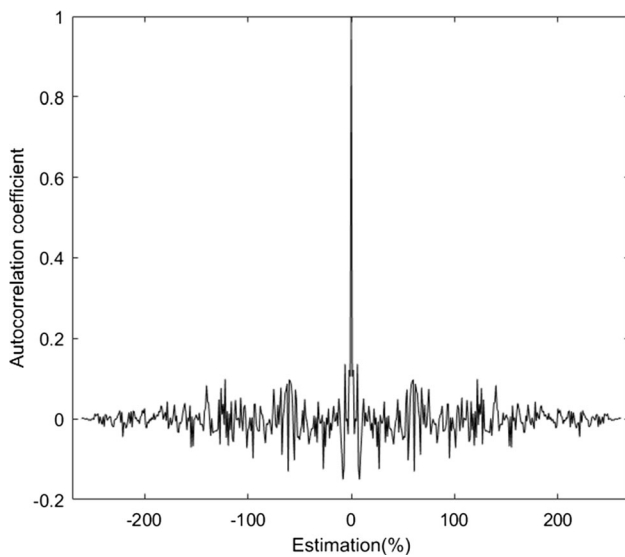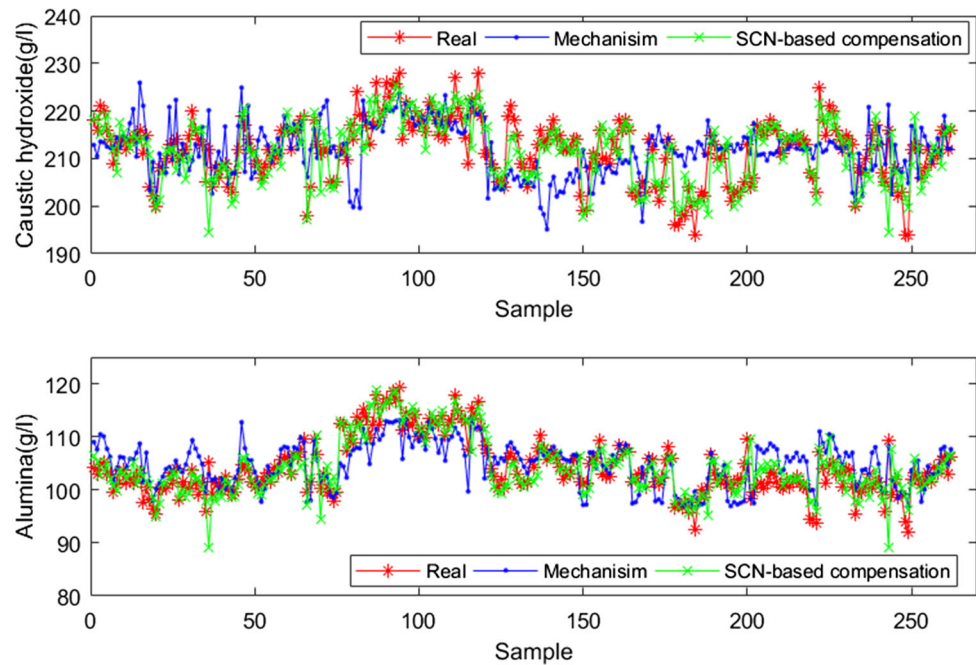


**Fig. 15** Autocorrelation function of estimation error for $c_K$ after SCN-based compensation
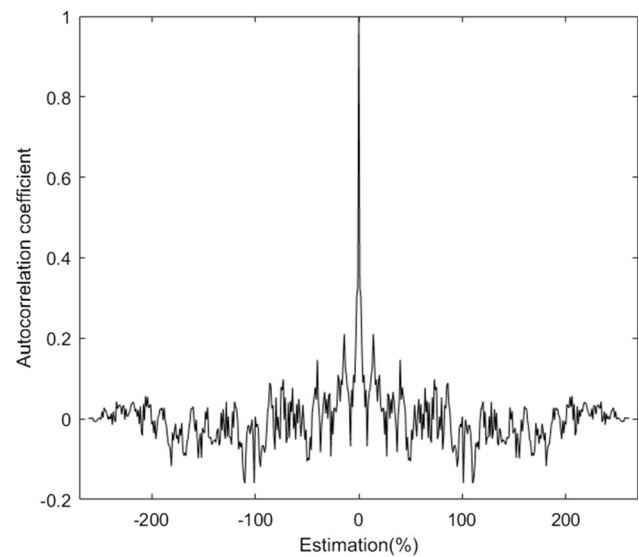


**Fig. 16** Autocorrelation function of estimation error for $c_A$ after SCN-based compensation



experiment is 30 min daytime (sampling interval in industrial locale is 2 h). We collected 524 records from the industrial locale in about 2 months. Half samples are used for model training, and the other half samples for test. Partial of the original data are shown in Table 4.

Based on the orthogonal experimental data, the mechanism model parameters in (13) are regressed, and using the calculated value of $y_{Km}$, the parameters in (14) are regressed. Then, the calculated values $y_{Am}$ can be obtained. The calculated values $y_{Km}$ and $y_{Am}$ will be used for the compensation model.

## 5.3 Results and discussion for $c_K$ and $c_A$

The results of mechanism model for $c_K$ and $c_A$ are shown in Fig. 9, and the autocorrelation function of mechanism modelling error for $c_K$ and $c_A$ is shown in Figs. 10 and 11.

From Fig. 9, we can see the trend of mechanism model estimation value is almost the same as the real value, indicating that the mechanism model is effective. From Figs. 10 and 11, we can see that the error of mechanism model should be improved and a compensation model is needed.

**Table 5** Test results of mechanism and compensation model

| Algorithm | Mechanism model | Mechanism and SCN-based compensation model |
| --- | --- | --- |
| RMSE of $c_K$ | 6.99 | 1.81 |
| RMSE of $c_A$ | 3.62 | 1.51 |

Using the proposed SCN-based compensation model algorithm, the training RMSE is shown in Fig. 12 and the model is built with 32 hidden nodes. The training and test results of $c_K$ and $c_A$ after compensation are shown in Figs. 13 and 14, and the autocorrelation function of estimation error for $c_K$ and $c_A$ after compensation is shown in Figs. 15 and 16. From these figures, we can see that good performance has been achieved after compensation, and the



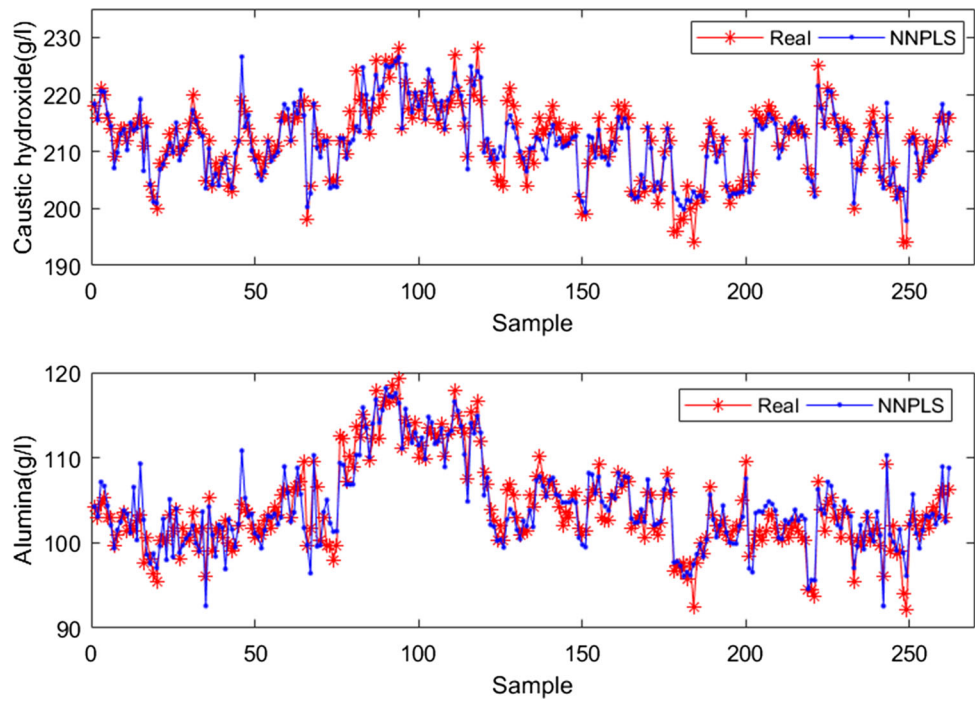**Fig. 17** Test results of $c_K$ and $c_A$ after NNPLS compensation



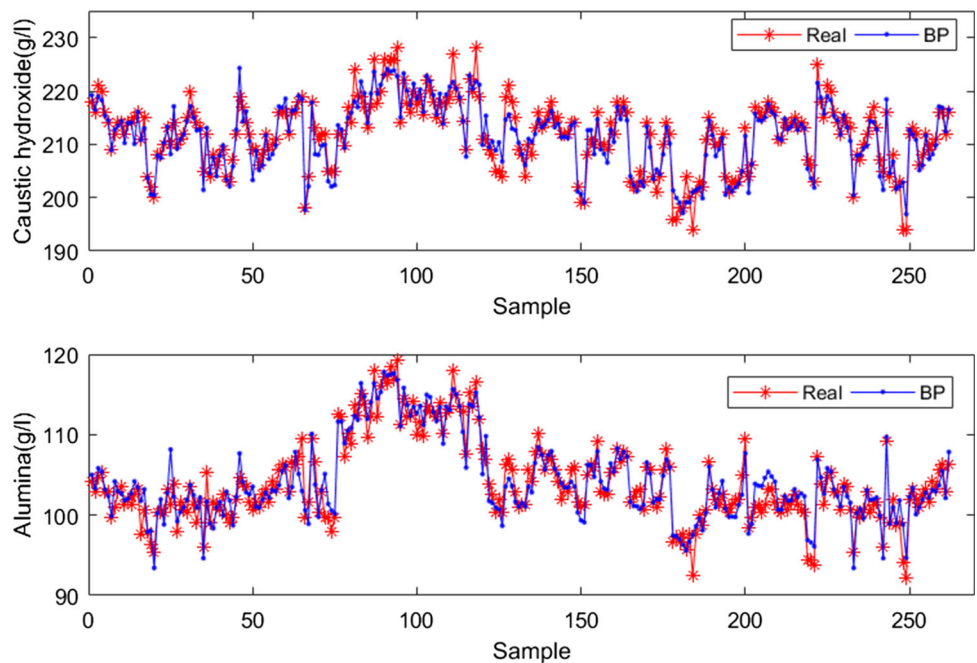**Fig. 18** Test results of $c_K$ and $c_A$ after BP compensation

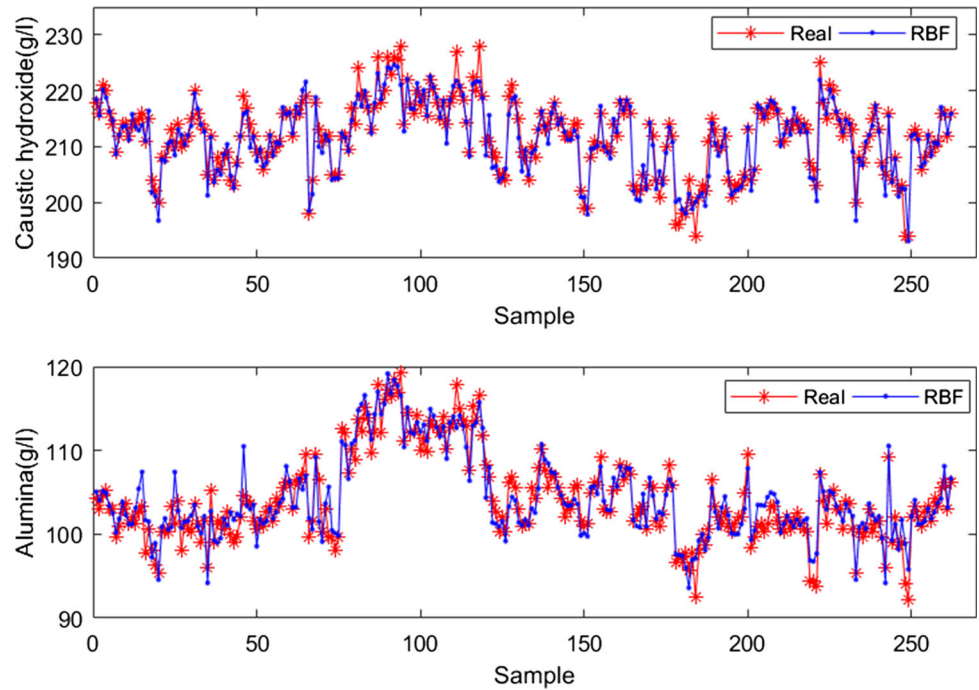**Fig. 19** Test results of $c_K$ and $c_A$ after RBF compensation



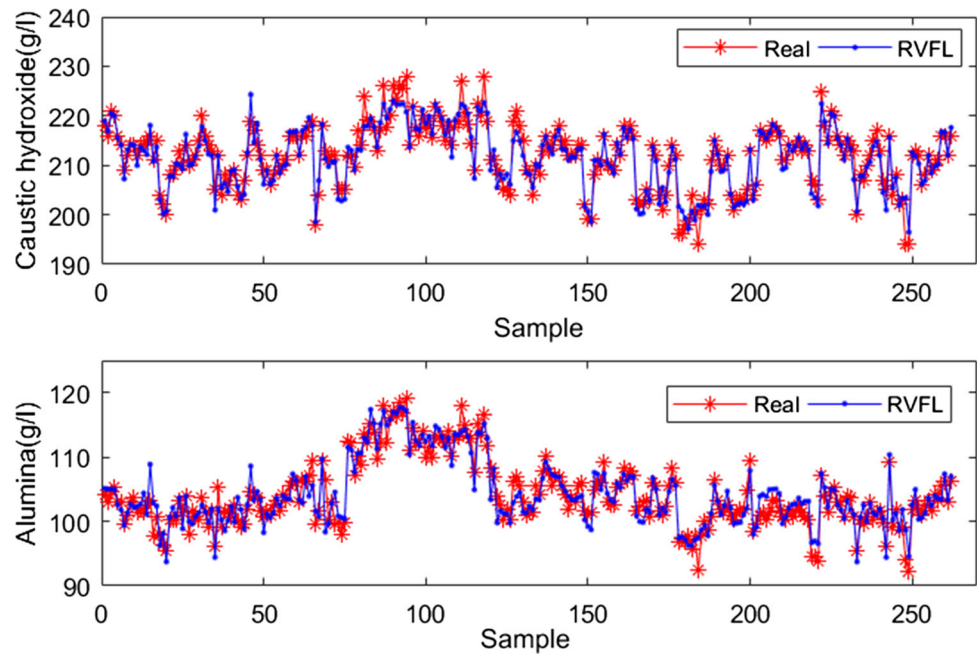**Fig. 20** Test results of $c_K$ and $c_K$ after RVFL compensation



**Table 6** Per cent variance captured by NNPLS model

| $c_K$ | X-Block | | Y-Block | | $c_A$ | X-Block | | Y-Block | |
|---|---|---|---|---|---|---|---|---|---|
| LV # | This LV | Total | This LV | Total | LV # | This LV | Total | This LV | Total |
| 1 | 15.48 | 15.48 | 76.08 | 76.08 | 1 | 26.23 | 26.23 | 49.33 | 49.33 |
| 2 | 27.09 | 42.57 | 8.24 | 84.33 | 2 | 25.82 | 52.45 | 11.54 | 60.87 |
| 3 | 23.32 | 65.89 | 2.03 | 86.36 | 3 | 11.24 | 63.69 | 7.16 | 68.04 |
| 4 | 16.89 | 82.78 | 1.14 | 87.5 | 4 | 9.45 | 73.15 | 4.79 | 72.83 |
| 5 | 9.99 | 92.77 | 0.17 | 87.67 | 5 | 14.75 | 87.9 | 0.35 | 73.18 |

**Table 7** Comparison of different compensation algorithms for $c_K$ and $c_A$

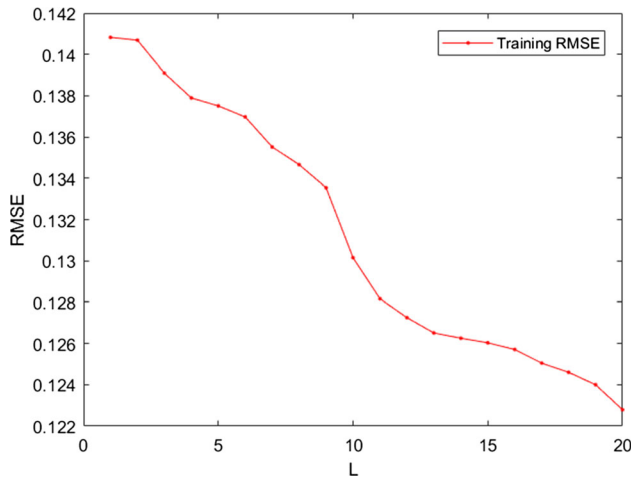| Algorithms | Training | | Test | | Runing time | Number of hidden nodes |
|---|---|---|---|---|---|---|
| | RMSE$c_K$ | RMSE$c_A$ | RMSE$c_K$ | RMSE$c_A$ | | |
| NNPLS | 2.5337 | 1.8782 | 2.2003 | 1.965 | 0.95 | – |
| BP | 2.2442 | 1.553 | 2.083 | 1.6335 | 11.04 | 18 |
| RBF | 1.7795 | 1.3976 | 1.7897 | 1.6949 | 6.01 | 260 |
| RVFL | 2.092 | 1.6421 | 1.9251 | 1.7856 | 0.93 | 32 |
| SCN | 1.8056 | 1.5116 | 1.8063 | 1.5111 | 1.39 | 32 |



**Fig. 21** RMSE of SCN training process for $c_C$

estimation accuracy has been improved significantly. There is RMSE (root mean square error) defined as follows:

$$\text{RMSE} = \sqrt{\frac{1}{N}\sum_{k=1}^{N}(\hat{y}(k) - y(k))^{\text{T}}(\hat{y}(k) - y(k))}. \quad (17)$$

The result comparison is shown in Table 5.

Compare the proposed SCN-based method with other methods for compensation, such as NNPLS, BP, RBF and RVFL; the best estimation values of each algorithm are shown in Figs. 17, 18, 19, and 20, respectively.

In the algorithm of NNPLS, 5 LVs has been chosen and the per cent variance is shown in Table 6.

The specific results of each algorithm are shown in Table 7.

From the comparison of the curves for different algorithms, NNPLS method has poor accuracy for both training

and test. BP and RBF results are better than NNPLS, but their running time is longer than other methods; it will be even longer for big industrial data. The number of hidden nodes for RBF is also larger. RVFL is faster than BP and RBF, but its accuracy is lower than them. From comprehensive analysis, we can see that SCN has better accuracy and fast learning speed; it also has better generalization than the other algorithms.

## 5.4 Results and discussion for $c_C$

Using the proposed SCN-based data-driven model, the RMSE of training process for $c_C$ is shown in Figs. 20, and 21 hidden nodes are chosen. From the figure, we can see the model estimation value almost can follow the real value trend, indicating the SCN-based data-driven model is effective. Different algorithms to predict $c_C$ are compared, such as NNPLS, BP, RBF and RVFL, with the same hidden nodes, and the results are shown in Table 8.

From the model test of $c_C$, 4 LVs is chosen for NNPLS and 20 hidden nodes are set by each neural network; we can see the conclusion is almost the same as $c_K$ and $c_A$ compensation model. The accuracy and running time of SCN are both better than other neural networks. NNPLS has shorter running time but lower accuracy. BP has lower generalization and longer running time. RBF has sound accuracy but longer running time. RVFL has shorter running time but lower accuracy.

In summary, the hybrid model which combined with mechanism analysis and data-driven techniques makes full use of the mechanism knowledge and is more effective than the single data-driven model we proposed for the alumina production process in [21]. Although the dynamic mechanism model of the component concentration is not

**Table 8** Comparison of different algorithms for $c_C$

| Algorithms | Training RMSE$c_C$ | Test RMSE$c_C$ | Runing time (s) | Number of hidden nodes |
|---|---|---|---|---|
| NNPLS | 1.1072 | 1.1075 | 0.47 | – |
| BP | 0.9863 | 1.1236 | 8.19 | 20 |
| RBF | 1.0156 | 1.0986 | 1.71 | 20 |
| RVFL | 1.025 | 1.1229 | 1.13 | 20 |
| SCN | 1.0068 | 1.0987 | 0.69 | 20 |

obtained, the static mechanism approximation model can still obtain better prediction results. However, single data-driven model is straightforward and simple, and it is more versatile than hybrid model for using in other similar industrial processes. The accuracy of both methods can meet the requirements of the alumina production process, and the user can decide which method to choose according to the actual object.

# 6 Concluding remarks

For Bayer alumina production process, better modelling performance of component concentrations in sodium aluminate liquor is significant for process control and optimization. A hybrid modelling strategy with mechanism model and SCN-based data-driven model is proposed in this paper. From the real-world application results, we can see that mechanism model reflects the internal relation and improved performance can be achieved by using the SCN-based compensation model and data-driven model. SCN model has the merits such as less human intervention on the network size setting, the scope adaptation of random parameters, fast learning and sound generalization. It fulfils the industrial requirements and provides a new online approach to measure the component concentrations.

Along with the topic addressed in this work, some further studies are being expected. For instance, data should be sampled continuously all day and the sampling interval should be the same. Also, sample data from locale may be contaminated with noises and outliers, robust stochastic configuration networks and ensemble data modelling techniques can be applied [19, 20]. Dynamic features and time-related variations in the process should also be considered in modelling. Besides, due to the sensors collected more input data, but the output laboratory data are less, semi-supervised learning algorithm should be developed.

# References

1. Broomhead DS, Lowe D (1988) Multi-variable functional interpolation and adaptive networks. Complex Syst 2:321–355
2. Browne G, Finn C (1977) Determination of aluminum content of Bayer liquors by electrical conductivity measurement. Metall Mater Trans B 8(1):349–349
3. Browne G, Finn C (1981) The effects of aluminum content, temperature and impurities on the electrical conductivity of synthetic Bayer liquors. Metall Trans B 12(3):487–492
4. Cong Q, Yu W, Chai T (2010) Cascade process modeling with mechanism-based hierarchical neural networks. Int J Neural Syst 20(01):1–11
5. Dai W, Li DP, Zhou P, Chai T (2019) Stochastic configuration networks with block increments for data modeling in process industries. Inf Sci 484:367–386
6. Igelnik B, Pao YH (1995) Stochastic choice of basis functions in adaptive function approximation and the functional-link net. IEEE Trans Neural Netw 6(6):1320–1329
7. Kadlec P, Gabrys B, Strandt S (2009) Data-driven soft sensors in the process industry. Comput Chem Eng 33(4):795–814
8. Kadlec P, Grbić R, Gabrys B (2011) Review of adaptation mechanisms for data-driven soft sensors. Comput Chem Eng 35(1):1–24
9. Li M, Wang D (2017) Insights into randomized algorithms for neural networks: practical issues and common pitfalls. Inf Sci 382:170–178
10. Ng C, Hussain M (2004) Hybrid neural network prior knowledge model in temperature control of a semi-batch polymerization process. Chem Eng Process 43(4):559–570
11. Pao YH, Takefji Y (1992) Functional-link net computing. IEEE Comput J 25(5):76–79
12. Psichogios DC, Ungar LH (1992) A hybrid neural network-first principles approach to process modeling. AIChE J 38(10):1499–1511
13. Qi H, Zhou XG, Liu LH, Yuan WK (1999) A hybrid neural network-first principles model for fixed-bed reactor. Chem Eng Sci 54(13–14):2521–2526
14. Scardapane S, Wang D (2017) Randomness in neural networks: an overview. Wiley Interdiscip Rev Data Min Knowl Discov 7(2):e1200
15. Tan A, Xiao C (1997) Direct determination of caustic hydroxide by a micro-titration method with dual-wavelength photometric end-point detection. Anal Chim Acta 341(2–3):297–301
16. Tan A, Zhang L, Xiao C (1999) Simultaneous and automatic determination of hydroxide and carbonate in aluminate solutions by a micro-titration method. Anal Chim Acta 388(1–2):219–223
17. Thompson ML, Kramer MA (1994) Modeling chemical processes using prior knowledge and neural networks. AIChE J 40(8):1328–1340
18. Wang D, Li M (2017) Stochastic configuration networks: fundamentals and algorithms. IEEE Trans Cybern 47(10):3466–3479
19. Wang D, Li M (2017) Robust stochastic configuration networks with kernel density estimation for uncertain data regression. Inf Sci 412:210–222
20. Wang D, Cui C (2017) Stochastic configuration networks ensemble with heterogeneous features for large-scale data analytics. Inf Sci 417:55–71
21. Wang W, Chai T, Yu W, Wang H, Su C (2011) Modeling component concentrations of sodium aluminate solution via Hammerstein recurrent neural networks. IEEE Trans Control Syst Technol 20(4):971–982
22. Wang W, Yu W, Zhao L, Chai T (2011) PCA and neural networks-based soft sensing strategy with application in sodium aluminate solution. J Exp Theor Artif Intell 23(1):127–136