



# iSuc-PseOpt: Identifying lysine succinylation sites in proteins by incorporating sequence-coupling effects into pseudo components and optimizing imbalanced training dataset



Jianhua Jia <sup>a, \*\*</sup>, Zi Liu <sup>a</sup>, Xuan Xiao <sup>a, b, \*</sup>, Bingxiang Liu <sup>a</sup>, Kuo-Chen Chou <sup>b, c, \*\*\*</sup>

<sup>a</sup> Computer Department, Jing-De-Zhen Ceramic Institute, Jing-De-Zhen 333403, China

<sup>b</sup> Gordon Life Science Institute, Boston, MA 02478, USA

<sup>c</sup> Center of Excellence in Genomic Medicine Research (CEGMR), King Abdulaziz University, Jeddah 21589, Saudi Arabia

## ARTICLE INFO

### Article history:

Received 10 November 2015

Received in revised form

2 December 2015

Accepted 11 December 2015

Available online 23 December 2015

### Keywords:

Lysine succinylation

Sequence-coupling model

PseAAC

Optimize training dataset

Target cross-validation

## ABSTRACT

Succinylation is a posttranslational modification (PTM) where a succinyl group is added to a Lys (K) residue of a protein molecule. Lysine succinylation plays an important role in orchestrating various biological processes, but it is also associated with some diseases. Therefore, we are challenged by the following problem from both basic research and drug development: given an uncharacterized protein sequence containing many Lys residues, which one of them can be succinylated, and which one cannot? With the avalanche of protein sequences generated in the postgenomic age, the answer to the problem has become even more urgent. Fortunately, the statistical significance experimental data for succinylated sites in proteins have become available very recently, an indispensable prerequisite for developing a computational method to address this problem. By incorporating the sequence-coupling effects into the general pseudo amino acid composition and using KNNC (K-nearest neighbors cleaning) treatment and IHTS (inserting hypothetical training samples) treatment to optimize the training dataset, a predictor called iSuc-PseOpt has been developed. Rigorous cross-validations indicated that it remarkably outperformed the existing method. A user-friendly web-server for iSuc-PseOpt has been established at <http://www.jci-bioinfo.cn/iSuc-PseOpt>, where users can easily get their desired results without needing to go through the complicated mathematical equations involved.

© 2015 Elsevier Inc. All rights reserved.

One of the most efficient biological mechanisms for expanding the genetic code and for regulating cellular physiology is the posttranslational modification (PTM) of proteins [1,2]. Owing to the importance of PTM in basic research and drug development, many efforts have been made with the aim of predicting various PTM sites in proteins (see, e.g., Refs. [3–10] and two review articles [11,12] published recently).

**Abbreviations:** PTM, posttranslational modification; PseAAC, pseudo amino acid composition; SVM, support vector machine; KNN, K-nearest neighbors; RF, random forest; 30-D, 30-dimensional; KNNC, K-nearest neighbors cleaning; IHTS, inserting hypothetical training samples; Acc, overall accuracy; MCC, Mathew's correlation coefficient; Sn, sensitivity; Sp, specificity; ROC, receiver operating characteristic; AUC, area under the curve.

\* Corresponding author. Gordon Life Science Institute, Boston, MA 02478, USA.

\*\* Corresponding author. Computer Department, Jingdezhen Ceramic Institute, Jingdezhen, 333043, China.

\*\*\* Corresponding author. Gordon Life Science Institute, Boston, MA 02478, USA.

E-mail addresses: [jih163yx@163.com](mailto:jih163yx@163.com) (J. Jia), [xxiao@gordonlifescience.org](mailto:xxiao@gordonlifescience.org) (X. Xiao), [kcchou@gordonlifescience.org](mailto:kcchou@gordonlifescience.org) (K.-C. Chou).

<http://dx.doi.org/10.1016/j.ab.2015.12.009>

0003-2697/© 2015 Elsevier Inc. All rights reserved.

The lysine residue in proteins can undergo many types of PTMs, such as methylation, acetylation, biotinylation, ubiquitination, ubiquitin-like modifications, propionylation, and butyrylation, leading to the remarkable complexity of PTM networks.

Recently, a new type of PTM, called lysine succinylation, was identified by mass spectrometry and protein sequence alignment. It has been shown that lysine succinylation responds to different physiological conditions and is evolutionary conserved [13]. In 2013, Park and coworkers [14] identified 2565 succinylation sites from 779 proteins and revealed that lysine succinylation has potential impacts on enzymes involved in mitochondrial metabolism, including amino acid degradation, tricarboxylic acid (TCA) cycle, and fatty acid metabolism [14]. Lysine succinylation also occurs in histones, suggesting that it may play an important role in regulating chromatin structures and functions as well [15,16]. Accordingly, identification of lysine succinylation sites in proteins is no doubt a crucial topic in cellular physiology and pathology, which can

provide very useful information for both biomedical research and drug development.

It is time-consuming and expensive to determine the succinylation residues by purely using the experimental techniques alone. In particular, facing the explosive growth of protein sequences in the postgenomic age, it is highly critical to develop computational tools for timely and effectively identifying the succinylation sites in proteins.

Actually, some computational methods have been proposed (see, e.g., Ref. [17]) for the aforementioned purpose. However, because of the importance of the topic, as well as the urgency of demanding more powerful high-throughput tools in this area, further efforts are definitely needed to enhance the prediction quality. The current study was initiated in an attempt to address this problem by developing a more powerful predictor via incorporating a vectorized sequence-coupling model [18] into the general form of pseudo amino acid composition (PseAAC) [19].

As shown in a series of recent publications [20–27] in compliance with Chou's five-step rule [19], to establish a really useful sequence-based statistical predictor for a biological system, we should logically follow the five guidelines below and make them crystal clear: (i) how to construct or select a valid benchmark dataset to train and test the predictor, (ii) how to formulate the biological sequence samples with an effective mathematical expression that can truly reflect their intrinsic correlation with the target to be predicted, (iii) how to introduce or develop a powerful algorithm (or engine) to operate the prediction, (iv) how to properly perform cross-validation tests to objectively evaluate its anticipated accuracy, and (v) how to establish a user-friendly web-server that is accessible to the public. Below, we address the aforementioned five procedures one by one.

## Materials and methods

### Benchmark dataset

The benchmark dataset used in this study was derived from the CPLM, a protein lysine modification database [28]. It contains 2521 lysine succinylation sites and 24,128 non-succinylation sites determined from 896 proteins [28]. All of the corresponding protein sequences were derived from the UniProt [29] database. For facilitating description later, Chou's peptide formulation was adopted. It was used for studying signal peptide cleavage sites [30], HIV protease cleavage sites [18], and protein–protein interaction [25]. According to Chou's scheme, a potential succinylation site-containing peptide sample can be generally expressed by

$$P_{\xi}(\mathbb{K}) = R_{-\xi}R_{-(\xi-1)} \cdots R_{-2}R_{-1}\mathbb{K}R_{+1}R_{+2} \cdots R_{+(\xi-1)}R_{+\xi}, \quad (1)$$

where the center  $\mathbb{K}$  represents “lysine,” the subscript  $\xi$  is an integer,  $R_{-\xi}$  represents the  $\xi$ -th upstream amino acid residue from the center,  $R_{+\xi}$  represents the  $\xi$ -th downstream amino acid residue, and so forth. The  $(2\xi + 1)$ -tuple peptide sample  $P_{\xi}(\mathbb{K})$  can be further classified into the following categories:

$$P_{\xi}(\mathbb{K}) \in \begin{cases} P_{\xi}^{+}(\mathbb{K}), & \text{if its center is a succinylation site} \\ P_{\xi}^{-}(\mathbb{K}), & \text{otherwise} \end{cases}, \quad (2)$$

where  $P_{\xi}^{+}(\mathbb{K})$  denotes a true succinylation segment with lysine at its center,  $P_{\xi}^{-}(\mathbb{K})$  denotes a false succinylation segment with lysine at its center, and the symbol  $\in$  means “a member of” in the set theory.

As elaborated in a comprehensive review [31], there is no need at all to separate a benchmark dataset into a training dataset and a testing dataset if the predictor to be developed will be tested by the

jackknife test or the subsampling (K-fold) cross-validation test because the outcome obtained in this way is actually from a combination of many different independent dataset tests. Thus, the benchmark dataset  $S_{\xi}$  for the current study can be formulated as

$$S_{\xi} = S_{\xi}^{+} \cup S_{\xi}^{-}, \quad (3)$$

where the positive subset  $S_{\xi}^{+}$  contains only the samples of true succinylation segments  $P_{\xi}^{+}(\mathbb{K})$  and the negative subset  $S_{\xi}^{-}$  contains only the samples of false succinylation segments  $P_{\xi}^{-}(\mathbb{K})$  (see Eq. (2)), whereas  $\cup$  represents the symbol for “union” in the set theory.

Because the length of peptide sample  $P_{\xi}(\mathbb{K})$  is  $2\xi + 1$  (see Eq. (1)), the benchmark dataset with a different  $\xi$  value will contain peptide segments with a different number of amino acid residues, as illustrated below:

The length of peptide samples in  $S_{\xi}$

$$= \begin{cases} 19 \text{ amino acid residues,} & \text{if } \xi = 9 \\ 23 \text{ amino acid residues,} & \text{if } \xi = 11 \\ 27 \text{ amino acid residues,} & \text{if } \xi = 13 \\ 31 \text{ amino acid residues,} & \text{if } \xi = 15 \\ 35 \text{ amino acid residues,} & \text{if } \xi = 17 \\ \vdots & \vdots \end{cases}. \quad (4)$$

The detailed procedures to construct  $S_{\xi}$  are as follows. First, as done in Ref. [32], slide the  $(2\xi + 1)$ -tuple peptide window along each of the 896 protein sequences taken from Ref. [28], and only those peptide segments that have K (Lys or lysine) at the center (see Eq. (1)) were collected. Second, if the upstream or downstream in a protein sequence was less than  $\xi$  or greater than  $L - \xi$  ( $L$  is the length of the protein sequence concerned), the lacking amino acid was filled with its mirror image (Fig. 1). Third, the peptide segment samples obtained in this way were put into the positive subset  $S_{\xi}^{+}$  if their centers have been experimentally annotated as the succinylation sites; otherwise, they were put into the negative subset  $S_{\xi}^{-}$ . Fourth, using the CD-HIT software [33], the aforementioned samples were further subject to a screening procedure to winnow those that had  $\geq 40\%$  pairwise sequence identity to any other in a same subset. By following the above procedures, we obtained a series of benchmark datasets with different  $\xi$  values.

But preliminary tests had indicated that it would be most promising when  $\xi = 15$ . Accordingly, for further study below, instead of Eq. (3) we shall consider

$$S_{\xi=15} = S_{\xi=15}^{+} \cup S_{\xi=15}^{-}, \quad (5)$$

where the benchmark dataset  $S_{\xi=15}$  contains 4720  $(2\xi + 1) = 31$ -tuple peptide samples, of which 1167 belong to the positive subset  $S_{\xi=15}^{+}$  and 3553 belong to the negative subset  $S_{\xi=15}^{-}$ . For readers' convenience, the detailed sequences of the aforementioned positive and negative samples are given in Supporting

### (A) Mirror image for N terminus

$$R_{-1}R_{-2} \cdots R_{-(\xi-1)}R_{-\xi} \Leftrightarrow R_{-\xi}R_{-(\xi-1)} \cdots R_{-2}R_{-1}$$

### (B) Mirror image for C terminus

$$R_{L-\xi}R_{L-\xi+1} \cdots R_{L-1}R_L \Leftrightarrow R_LR_{L-1} \cdots R_{L-\xi+1}R_{L-\xi}$$

**Fig.1.** Schematic illustration to show the mirror images of the  $\xi$  residues for the N terminus (A) and the C terminus (B). The red symbol  $\Leftrightarrow$  represents a mirror, and the real peptide segment is colored in black, whereas its mirror image is colored in blue. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

Download English Version:

<https://daneshyari.com/en/article/1173056>

Download Persian Version:

<https://daneshyari.com/article/1173056>

[Daneshyari.com](https://daneshyari.com)