METHODOLOGY ARTICLE

A New Algorithm to Train Hidden Markov Models for Biological Sequences with Partial

Labels

Jiefu Li¹
, Jung-Youn Lee^{2,3}
and Li Liao^{1,3*}

*Correspondence:

lliao@cis.udel.edu

¹University of Delaware, Computer and Information Sciences, 101 Smith Hall, Newark, DE 19716 US Full list of author information is available at the end of the article

Abstract

Background: Hidden Markov models (HMM) are a powerful tool for analyzing biological sequences in a wide variety of applications, from profiling functional protein families to identifying functional domains. The standard method used for HMM training is either by maximum likelihood using counting when sequences are labelled or by expectation maximization, such as the Baum-Welch algorithm, when sequences are unlabelled. However, increasingly there are situations where sequences are just partially labelled. In this paper, we designed a new training method based on the Baum-Welch algorithm to train HMMs for situations in which only partial labeling is available for certain biological problems.

Results: Compared with a similar method previously reported that is designed for the purpose of active learning in text mining, our method achieves significant improvements in model training, as demonstrated by higher accuracy when the trained models are tested for decoding with both synthetic data and real data.

Conclusions: A novel training method is developed to improve the training of hidden Markov models by utilizing partial labelled data. The method will impact on detecting de novo motifs and signals in biological sequence data. In particular, the method will be deployed in active learning mode to the ongoing research in detecting plasmodesmata targeting signals and assess the performance with validations from wet-lab experiments.

Keywords: hidden Markov model; partial label

Li et al. Page 2 of 21

¹Background

²Hidden Markov model [1, 2, 3, 4, 5] is a well known probabilistic model in the field ² of machine learning, suitable for detecting patterns in sequential data, such as plain ⁴texts, biological sequences, and time series data in the stock market. For all these ⁵applications, successful learning depends, to a large degree, on the amount and, ⁵ ⁶more importantly, the quality of the data. In text mining problem, though the data ⁷amount is huge, careful labelling tasks consume massive human labor [6]. In bio-⁸logical sequence analysis, discovering de novo signal remains challenging because a ⁹precise full labeling via wet-lab experiments demand even more resources and time, ⁹ and hence it is considered unfeasible in general. Therefore, research is necessary in ¹¹data handling with different labelling quality for applied machine learning commu-¹² nity. In this paper, we focus on designing a Baum-Welch-algorithm based learning ¹² ¹³method for HMMs to handle the biological problems when only partial labeling is ¹³ ¹⁴available in training data. This work is inspired by our recent research on detecting de novo plasmodes mata 15 targeting signals in Arabidopsis Plasmodesmata-located proteins (PDLPs). PDLPs 16 ¹⁷ are type I transmembrane proteins, which are targeted to intercellular pores called ¹⁷ plasmodesmata that form at the cellular junctions in plants [7]. In our study [8], by building a 3-state HMM, we predicted the presence of two different plasmodesmata¹⁹ ²⁰ targeting signals (named alpha and beta) in the juxta membrane region of PDLPs. ²⁰ ²¹While all the predicted signals were successfully verified in wet-lab experiments 22 so far, some predicted signals contain residues that do not conform to the true ²³ signal; wet-lab experiments showed that those residues alone was not sufficient to target the protein to plasmodesmata. Because both the cost and time are high for 24 wet-lab experiments, an improved HMM would be highly desirable. However, due 25 ²⁶ to the limitation in the number of the training examples – Arabidopsis genome encodes only eight PDLP members, further improvements of the model can be 27 hardly achieved. It would require to fully utilize the current wet-lab experimental 28 results to train the model, i.e., by labeling the residues that have been already shown to be either part of the signals or not part of the signals, given that labels 30 $^{31} \mathrm{are}$ not available for all the residues due to limited experimental results. In a related work by Tamposis et al, a semi-supervised approach is developed to handle a mixture of training sequences that contains a subset of fully labelled

Li et al. Page 3 of 21

¹sequences, with the remaining sequences having no labels at all or partial labels [9]. ¹ ²Their method uses the fully labelled sequences to train the parameters for HMMs² ³and then use Viterbi algorithm to predict the missing labels followed by training the³ ⁴model again with the predicted labels. This process is iterated until a convergence ⁴ ⁵condition is met. Instead, we are specifically interested in situations where no fully ⁵ ⁶labelled sequences are available and often the partial labeling is also sparse. In the ⁶ ⁷text mining field, HMM training algorithm of handling partial label was developed ⁷ ⁸especially for active learning purposes and designed to fit into text mining special⁸ ⁹situation: no label scenario, or in other words, no meaningful label can be assigned ¹⁰[6]. However the unit of observation in text mining and information retrieval is ¹⁰ ¹¹a word, instead of a single letter, corresponding to individual amino acid residue ¹¹ ¹²as in biological sequences. So, in order to deal with the partial labeling aforemen-¹² ¹³tioned, we have designed a novel Baum-Welch based HMM training algorithm to ¹³ ¹⁴leverage partial label information with techniques of model selection through par-¹⁴ ¹⁵tial labels. Besides the difference in the observation unit, our algorithm differs from ¹⁵ ¹⁶[6] primarily in how to calculate the expected value for a given partial label at a¹⁶ ¹⁷given position: our method sums over hidden state paths that must be subject to ¹⁷ ¹⁸constraints anywhere given partial labels in the training sequence. In contrast, in ¹⁸ ¹⁹[6] the expected value for a given partial label at a given position is calculated by ¹⁹ ²⁰summing over paths that are only constrained at the position being considered, and ²⁰ ²¹anywhere else in the sequence the hidden paths are free to go through all possible ²¹ ²²states (labels) even at positions where partial labels are given. Moreover, this differ-²² ²³ence affects how the expected value for a transition is calculated, regardless whether ²³ ²⁴the transition happens to involve one partial label, two partial labels, or no partial ²⁴ 25 labels at all. The comparison between our method and the method described in $[6]^{25}$ ²⁶showed that our method outperformed in both synthetic and real data for decoding ²⁶ ²⁷task in biological problems. The rest of this paper is organized as follows. First, the relevant background 29 knowledge of HMM is briefly reviewed, and notations are introduced. Then, our 30 method of training HMM when only partial label sequences are available is described in details. This is followed with experiments and results to examine and demonstrate 32 the modelling power of the novel algorithm. Discussion and conclusion are given at 33 the end.

Li et al. Page 4 of 21

¹Methods ²Hidden Markov model review ³In general, a HMM consists of a set of states S_i , i=1 to N, and a set of alphabets ${}^{4}K$ that can be emitted from these states with various frequencies; $b_{i}(k)$ stands for ⁵the frequency of letter $k \in K$ being emitted from state S_i , and we use B to denote ⁵ ⁶the emission matrix of dimension $N \times K$, containing $b_i(k)$ as elements. Transitions ⁷among states can be depicted as a graph, often referred as model architecture or ⁷ 8 model structure: each state is represented as a node, and transition from state $S_i^{\ 8}$ ⁹to state S_i is represented by a directed edge, with a weight a_{ij} being the transition ⁹ probability, and we use A to denote the transition matrix of dimension $N \times N$, ¹⁰ "containing a_{ij} as elements. Hereafter, we often refer to a state S_i by its index i. Given a HMM, let θ stand for collectively all its parameters, namely the emission ¹² ¹³ frequencies $b_i(k)$ and transition probabilities a_{ij} . Given a sequence of observation ^{14}O , and its elements $O_t \in K$, where t = 1..T, a main assumption of using HMM 14 15 is that each letter in the sequence is emitted from a state of the model, so correspondingly there is a state sequence, forming a Markov chain, which is hidden from ¹⁶ ¹⁷ direct observation, hence the name: hidden Markov model. One task (decoding) is, ¹⁷ therefore, to find the most probable state sequence (also called hidden path) X^* : 18 $^{19}X^* = \operatorname{argmax}_X Pr(O, X|\theta)$, among all possible state sequences that can emit the observation sequence O. The second task is to train the model on a set of m training sequences. This task is accomplished by adjusting model parameters θ to maximize the likelihood $\sum_{s=1}^{m} Pr(O^{s}|\theta)$ of observing the given training sequences O^{s} , where O^{s} s = 1...m [10]. The decoding task is well studied and straightforward and is solved by Viterbi 24 ²⁵ algorithm efficiently [11]. The technique guarantees to return the optimal answer. ²⁶Note that, in the work by Bagos et al [12], a modified Viterbi algorithm is developed ²⁶ $^{\rm 27}$ to incorporate prior topological information as partial labels to improve predictions, whereas our focus is instead on how to use the partial labels in training the model. ²⁸ ²⁹ However, the second task, or the training of a HMM is not guaranteed to reach 30 optimum when labels are not given for the training sequences. The major training algorithms of HMM are the following three in general: maxi- 32 mum likelihood, Baum-Welch algorithm, and Viterbi training [13]. Maximum like- 33 lihood is used when label information is available fully, and it returns the optimal 33

Li et al. Page 5 of 21

¹ solution. The latter two algorithms are used when no label information is available. ¹
$^2 \mathrm{Interested}$ readers can find a gentle introduction and tutorial for hidden Markov^2
³ models in [10]. For the purposes of comparison, we adopt notations in [6] for future ³
4 discussion of both the background knowledge and our method. The description of 4
⁵ notations is shown in Table 1.
6 In this paper, we focus on a special case for training HMMs when only $partial^{6}$
⁷ label is available. Or in other words, we aimed at finding model θ so that $Pr(O \theta)^7$
$^8{\rm is}$ maximized (locally) and the resulting decoded state sequence must satisfy the 8
⁹ partial labels given in the training sequences at the same time. 10
Training hidden Markov model with partial label sequences
$_{12}\mathrm{As}$ introduced in the previous section, when no labels are available, Baum-Welch $_{12}$
$_{13}{\rm algorithm}$ is typically used to train HMM and Viterbi training is sometimes used $_{13}$
$_{_{14}}\!\!$ for speed and simplicity; when all label information is given, training HMM is $_{_{14}}$
$_{15}\mathrm{straight}$ forward by maximum likelihood approach. Currently, training HMM with $_{15}$
$_{16} \mathrm{partial}$ label is mainly studied in the field of text mining, with a particular focus on $_{16}$
$_{17}{\rm active}$ learning problems, such as the work done in [6], with which we compare our $_{17}$
₁₈ proposed method.
$_{19}$ Our proposed method is a novel approach to this partial label training problem $_{19}$
$_{20}\mathrm{with}$ modification of Baum-Welch algorithm (called constrained Baum-Welch algo- $_{20}$
$_{21}\mathrm{rithm})$ and a model selection technique, which helps our algorithm leverage available $_{21}$
$_{22} \mathrm{information}$ and improve the training and performance in decoding task. In the next $_{22}$
$_{23} \mathrm{two}$ subsections, we discuss in detail our constrained Baum-Welch algorithm and $_{23}$
$_{24}\mathrm{the}$ model selection methods respectively and how to combine the two for model $_{24}$
$_{24}{\rm the~model~selection~methods~respectively~and~how~to~combine~the~two~for~model}_{24}$ $_{25}{\rm training.}$
25 training. 25
25 training. 25 26 Constrained Baum-Welch algorithm 27
$_{25}^{\rm training.}$ $_{25}^{\rm 26}$ Constrained Baum-Welch algorithm $_{\rm 27}^{\rm 27}$ The standard Baum-Welch algorithm is an Expectation-Maximization approach to $_{\rm 28}^{\rm 27}$
25 training. 26 Constrained Baum-Welch algorithm 27 The standard Baum-Welch algorithm is an Expectation-Maximization approach to 28 maximizing likelihood when the system contains latent variables, which are the state 28
25 training. 26 Constrained Baum-Welch algorithm 27 The standard Baum-Welch algorithm is an Expectation-Maximization approach to 28 maximizing likelihood when the system contains latent variables, which are the state 29 sequences for hidden Markov models when training sequences are not labelled. Our 30
25 training. 26 Constrained Baum-Welch algorithm 27 The standard Baum-Welch algorithm is an Expectation-Maximization approach to 28 maximizing likelihood when the system contains latent variables, which are the state 29 sequences for hidden Markov models when training sequences are not labelled. Our 30 constrained Baum-Welch algorithm (cBW) is similar to the standard Baum-Welch 31
25 training. 26 Constrained Baum-Welch algorithm 27 The standard Baum-Welch algorithm is an Expectation-Maximization approach to 28 maximizing likelihood when the system contains latent variables, which are the state 29 sequences for hidden Markov models when training sequences are not labelled. Our 29 constrained Baum-Welch algorithm (cBW) is similar to the standard Baum-Welch 30

Li et al. Page 6 of 21

¹cBW algorithm is identical to standard Baum-Welch's. The difference is the E-step, ¹
²computing forward and backward matrices. The forward matrix α is of $N \times T$, ²
³where N is the number of states and T is the sequence length. An element $\alpha_i(t)$ is ⁴the probability of the observed sequence up to and including O_t , with the symbol ⁴
⁵ O_t being emitted from state i. The backward matrix β is of $N \times T$ dimension has ⁵
⁶element $\beta_i(t)$ as the probability of the observed sequence from position t onto the ⁶
⁷end, with the symbol O_t being emitted from state i. The formulas of computing α ⁷
⁸and β are shown as following respectively.

⁹ Given the model $\theta = (\pi, A, B)$, where π is a N dimension vector, with π_i being ¹⁰ the probability that any hidden state path would start with state i. Then, the ¹¹ initial values of forward matrix α for one given training sequence $O = (O_1, ..., O_T)^{11}$ is computed as follows.

$$\alpha_i(1) = \pi b_i(O_1) \tag{1)14}$$

15

 $_{16}$ After calculating the initial values of α , by dynamic programming, the remaining $_{16}$ values at any position for any state are calculated recursively by summing over the $_{17}$ $_{18}$ possible state paths $X=(X_1,...,X_T)$, allowed by the model, that lead to the point $_{18}$ whose α value is being calculated. However, since we now have partial labels for the $_{19}$ $_{20}$ training sequence O, care must be taken to satisfy the constraints at each position $_{20}$ $_{21}O_t$ imposed by the partial label, $L(O_t) \in S \cup \{0\}$, where a value zero means no label $_{21}$ $_{22}$ available. Specifically,

23
$$\alpha_i(t+1) = \begin{cases} b_i(O_{t+1}) \sum_{j=1}^N \alpha_j(t) a_{ji}, & \text{if } L(O_{t+1}) = 0 \text{ or } i \\ 0 & \text{if } L(O_{t+1}) \neq 0 \text{ and } L(O_{t+1}) \neq i \end{cases}$$
23
$$(2)^{24}$$
25

In the above equation, the first case is when position O_{t+1} is either unconstrained ²⁶ $^{27}(0)$ or constrained to be state i by the partial label. In such a case, the α value ²⁷ is computed in the same way as the standard Baum-Welch algorithm, though the ²⁸ actual value can still be affected by partial labels at earlier positions via recursion. ²⁹ The second case is when the position t+1 is constrained by the partial label to be ³⁰ a state other than i. In this case, $\alpha_i(t+1) = 0$. This latter case is what makes the ³¹ algorithm different from the standard Baum-Welch algorithm in order to "honor" ³² the partial labels.

Li et al. Page 7 of 21

The backward matrix β is initialized as the following.

2

1

$$\beta_i(T) = 1 \tag{3}$$

4

 $^5{\rm Then},$ similarly, a recursive procedure is applied for the remaining of backward $^6{\rm matrix}.$

7

$$\beta_{i}(t) = \begin{cases} \sum_{j=1}^{N} \beta_{j}(t+1)a_{ij}b_{j}(O(t+1)), & \text{if } L(O_{t}) = 0 \text{ or } i \\ 0 & \text{if } L(O_{t}) \neq 0 \text{ and } L(O_{t}) \neq i \end{cases}$$

$$(4)_{9}$$

¹¹Note that, while the α is calculated the same way as the modified Forward algorithm. ¹² in [12] but the β is calculated differently from their modified Backward algorithm. ¹² ¹³After the calculations of α and β , then we can calculate γ variable, where $\gamma_i(t)$ is ¹³ ¹⁴the probability of observing the training sequence O from all possible state paths ¹⁴ that are allowed by hidden Markov model θ as constrained by the partial labels and ¹⁵ ¹⁶go through state i at position t. $\gamma_i(t)$ is computed as follows.

$$P(X(t)) = t \cdot O(t)$$

1

₂₂where the last equal sign holds because $P(O|\theta) = \sum_{j=1}^{N} \alpha_j(t)\beta_j(t)$. The next step₂₂
₂₃is to compute $\xi_{ij}(t)$, which is the probability of of observing the training sequence₂₃
₂₄O from all possible state paths that are allowed by hidden Markov model θ as₂₄
₂₅constrained by the partial labels and go through state i at positive t and transition₂₅
₂₆to state j at position t+1:

28
$$\xi_{ij}(t) = \frac{P(X(t) = i, X(t+1) = j, O|\theta)}{P(O|\theta)}$$
 28

$$= \frac{\alpha_i(t)a_{ij}\beta_j(t+1)b_j(O(t+1))}{P(O|\theta)}$$
(6)²⁹

³¹Finally, with γ , ξ , the M-step is to update the initial probability π^* , every elements ³² of the transition matrix A^* : a_{ij}^* , and every elements of the emission matrix B^* : ³³ $b_i^*(o_k)$.

Li et al. Page 8 of 21

2 2

$$\pi(i)^* = \gamma_i(1) \tag{7}$$

$$a_{ij}^* = \frac{\sum_{t=1}^{T-1} \xi_{ij}(t)}{\sum_{t=1}^{T-1} \gamma_i(t)}$$

$$(8)^4$$

$$b_i^*(o_k) = \frac{\sum_{t=1}^{T-1} \gamma_i(t) I_{O(t)=o_k}}{\sum_{t=1}^{T-1} \gamma_i(t)}$$

$$(9)_6$$

where $I_{O(t)=o_k}$ stands for indicator function, which equals to 1 if $O(t)=o_k$, and 0 otherwise. Then, for the case of multiple sequences, each sequences indexed by s, total number of sequences of m, The only changing is the updating of π^*, A^* , and $_{11}B^*$ as follows.

$$\pi(i)^* = \frac{\sum_{s=1}^m \gamma_i^s(1)}{m}$$

$$13$$

$$14$$

$$a_{ij}^* = \frac{\sum_{s=1}^m \sum_{t=1}^{T^s-1} \xi_{ij}^s(t)}{\sum_{s=1}^m \sum_{t=1}^{T^{s-1}} \gamma_i^s(t)}$$

$$(10)^{12}$$

$$(11)_{14}$$

$$a_{ij}^* = \frac{\sum_{s=1}^m \sum_{t=1}^{T^s-1} \xi_{ij}^s(t)}{\sum_{s=1}^m \sum_{t=1}^{T^s-1} \gamma_i^s(t)}$$
(11)₁₄

$$b_i^*(o_k) = \frac{\sum_{s=1}^m \sum_{t=1}^{T^s-1} \gamma_i^s(t) I_{O^s(t) = o_k}}{\sum_{s=1}^m \sum_{t=1}^{T^s-1} \gamma_i^s(t)} \tag{12}$$

The procedure above is repeated till either the $\sum_{s}^{m} log(P(O^{s}|\theta))$ converge or reaching maximum iteration numbers set by the user. As mentioned in the Introduction section, a key difference between our method and [6] lies in the E-step for calculating the expected value for a given emission or transition. Our method $\dot{}$ handles the partial label constraints recursively for the α and β , whereas [6] calculates α and β without using the partial labels and only uses the partial labels²² in resetting γ at each partial labelled position independently, as if partial labels 24 elsewhere would have no effect for the position being considered. Since E-step in $^{25} \mbox{Baum-Welch}$ algorithm invokes forward and backward algorithms, which are essentiated as $^{25} \mbox{Baum-Welch}$ tially a dynamic programming to more efficiently calculate the likelihood: $Pr(O|\theta)$ $L = \sum_{X \in \Gamma} Pr(O, X | \theta)$ with Γ being the set of all hidden paths, and hence should give the same result when the likelihood is computed by exhaustively summing over probability for all possible hidden state paths. Therefore, we believe, the partial labels would restrict the possible hidden state paths, $Pr(O|\theta) = \sum_{X \in \Gamma'} Pr(O, X|\theta)^{30}$ with Γ' being the set of all hidden paths constrained by partial labels and such constraints should be handled recursively in the dynamic programming. Figure 1 shows an example for the forward/backward dynamic programming table construcLi et al. Page 9 of 21

¹tion. Another advantage of our method comparing with the method in [6] is that ¹
²our training method can keep the topology of the initial transition and emission ²
³guesses for the model as standard Baum-Welch does. In other words, if prior knowl-³
⁴edge is available for the model topology, our training method for partial label data ⁴
⁵can keep the knowledge to the end of training.

⁵
⁶

⁷Model selection based on partial label information

The second part of our method is model selection based on partial label informa-8 9tion. The rationale is straightforward: while the constrained Baum-Welch algorithm9 10 increases the log-likelihood of the given training sequences (with partial labels), it-10 11 eration after iteration monotonically as ensured by EM approach, there is no direct 11 12 guarantee that the increased log-likelihood would necessarily lead to higher decod-12 13 ing accuracy. Therefore, at each iteration of constrained Baum-Welch algorithm, 13 14 decoding accuracy for the partially labelled training sequence can be calculated 14 15 and factored into model selection.

Specifically, after reaching convergence condition or maximum number of iter-16 17 ations, the total number of iteration is Q and the i^{th} iteration's model and the 17 18 corresponding log-likelihood can be denoted as θ_i and $\sum_s^m Log(P(O^s|\theta_i))$ respec-18 19 tively and let the decoding accuracy denote as $Accuracy(\theta_i, O, X)$. The final model 19 20 returned by the algorithm can be expressed as:

21
$$\underset{\theta^*}{\operatorname{argmin}} \Pr(O|\theta^* \equiv \{\underset{\theta_i \in \theta_1 \dots Q}{\operatorname{argmax}} \operatorname{Accuracy}(\theta_i, O, X)\}) \tag{13}_{22}$$

Notice that θ^* is a set of models in general. Finally, combining the constrained 24 Baum-Welch and the model selection described above, the overall algorithm of our 25 proposed method is given in Algorithm 1. In next section, Table 2 to Table 5 will 26 show the usefulness of both this model selection method and the ability of keeping 27 correct topology of cBW method.

²⁹Results

In this section, we set up experiments using both real biological data and synthetic data to test our method for decoding task and compared the results with those from using the method in [6]. It has been reported that [14, 15] posterior decoding in general performs better than Viterbi algorithm. So, in order to evaluate how our

Li et al. Page 10 of 21

```
Algorithm 1: Constrained Baum-Welch with model selection based on de-
                                                                                                2
  coding accuracy.
   Data: sequences: O, partially label X
   Result: bestA,bestB
                                                                                                4
   initialization: \theta = (\pi, A, B);
   bestAccuracy \leftarrow 0;
    while not reaching maximum iteration nor convergent do
        calculate \alpha, \beta by Eq. 1 to 4;
        calculate \gamma, \xi by Eq. 5 to 6;
        update \pi^*, A^*, B^* by Eq. 7 to 9;
8
        X^* \leftarrow \mathsf{Decoding}(\pi^*, A^*, B^*, O);
                                                                                                8
        myAccuracy \leftarrow Accuracy(X^*,X);
        if myAccuracy > bestAccuracy then
            bestAccuracy \leftarrow myAccuracy;
10
                                                                                                10
            \pi \leftarrow \pi^*;
11
                                                                                                11
            \mathsf{bestA} \leftarrow A^* \; ;
            \mathsf{bestB} \leftarrow B^*:
12
                                                                                                12
        end
13
                                                                                                13
   end
14
                                                                                                14
training method can impact on decoding, we carried out the decoding on the testing _{16}
sequences with the trained model using both the standard Viterbi algorithm \begin{bmatrix} 10 \end{bmatrix}_{47}
and Posterior-Viterbi algorithm described in [15], and the accuracy was computed _{\mbox{\tiny 18}}
by comparing the predicted label with the ground truth label at each position to
determine the number of correct predictions:
                                                                                                20
21
                        Accuracy = \frac{\#of\ correct\ predicted\ labels}{\#of\ total\ labels}
22
                                                                                                22
^{23}The results of these experiments show that our method outperforms Scheffer et al's
method in model training, as evidenced in the improved decoding accuracy, regard-
^{25} less which decoding algorithm is used. Specifically, on average, decoding accuracy
^{26} improves by 33% with Viterbi algorithm, 36% with Posterior-Viterbi algorithm
^{27} in real data, and improves by 7.35\% to 14.06\% with Viterbi algorithm, 7.08\% to
^{28}13.89\% with Posterior-Viterbi algorithm in synthetic data with significant p-values.
Note that, in two cases when the sequences are either almost fully labelled (95\%) or
wery sparsely labelled (5%), the differences between various algorithms are insignif-
 icant. This phenomenon is no surprising though, as it is expected that the benefit
^{32} from making good use of partial labels diminishes when labels are extremely sparse,
 which makes the various algorithms converge to Baum-Welch algorithm, or when
```

Li et al. Page 11 of 21

¹sequences are almost fully labelled, which makes the various algorithms converge to ¹ ²the maximum likelihood. Therefore, our evaluations are divided into two settings² ³ for synthetic data. Setting 1 has partial label information from 5% to 95%. Setting ³ ⁴2 has partial label information from 10% to 90%. ⁶Synthetic data ⁷The method described in [6] is mainly focused on handling text mining problems ⁸using synthetic data. To make the comparison fair, we have also performed exper-⁸ ⁹ iments using synthetic data, which allowed us to observe our method's different ¹⁰performance in different situations. In the experiments with synthetic data, the ¹⁰ ¹¹data is generated from ground truth HMMs, which are also generated randomly ¹¹ ¹² with predefined connections. For each experiment, the size for initial guess of tran-¹² ¹³sition and emission matrices are identical to the corresponding ground truth model. ¹³ ¹⁴We fixed the number of symbols in hidden Markov model to be 20 to mimic the 20¹⁴ ¹⁵amino acids in protein sequences. To test how model complexity may impact the ¹⁵ ¹⁶training, we chose three different numbers of states: 3, 5, and 7. Moreover, differ-¹⁶ ¹⁷ent levels of training sample size were also considered as an experimental variable. ¹⁷ ¹⁸Each experiment (with fixing state number and training examples) was evaluated ¹⁸ ¹⁹ for different levels of partial label and repeated for 50 times, and the corresponding ¹⁹ ²⁰paired p-values were also calculated to assess the statistical significance of the per-²⁰ ²¹ formance difference between our method and the other method. Since our method ²¹ ²²can maintain the topology of initial guess of transition matrix, experiments were ²² ²³divided into two different groups. One was initialized the transition matrix with the ²³ ²⁴same connectivity as the ground truth model, and the other was initialized with ²⁴ ²⁵fully connected transition matrix. Three sets of experimental results with fully connected transition matrix as initial²⁶ $^{\rm 27}$ guess are shown in Figure 2 to Figure 4. Additional results are shown in Table 2 to ²⁸Table 5 for comparison. Conducted using different numbers of states, training examples, and different decoding algorithms, the results show that our method outperforms the method by 31 Scheffer et al by 7.08% to 14.06% across different percentage of unlabelled data, with significant p-value (< 0.05) for majority of the experiments. While both methods achieve a performance closer to that of the ground truth model as the level of partial³³

Li et al. Page 12 of 21

¹labels increases, the improvement of our method over the method of Scheffer et al's ² is more pronounced when partial labels are sparse, namely the level of unlabelled ³data is high, as shown in the X-axis of the Figures. For example, in Figure 2, ³ ⁴with Viterbi decoding, at the level of 70% unlabelled data, i.e., 30% partial labels, ⁴ ⁵our method reaches an accuracy of 62%, which is 98% of the ground truth model⁵ ⁶accuracy, whereas Scheffer et al's reaches accuracy of 54%, which is 85% of the ⁶ 7 ground truth model accuracy. Similar trends hold true for Figure 3 and Figure 4^7 ⁸when the model has 5 and 7 states respectively regardless of the decoding algorithm⁸ ⁹used. 10 10 ¹¹Real data For the real biological data, we adopted data from [16]. The data contains 83 multipass transmembrane proteins with complete label information. The topology of ¹³ multi-pass transmembrane protein is shown in Figure 5. The label for each sequence contain three different values: i, o, M. They stand for the region of protein sequence inside, outside cell membrane, and the transmembrane domain respec- $^{17}\mathrm{tively}.$ While much more sophisticated hidden Markov models have been used for modeling transmembrane protein topology [16, 17, 18, 19], a simple HMM is used in this study to primarily evaluate the new training algorithm for partial labels. The architecture of the HMM is shown in Figure 6, in which a redundant M' node is introduced as a simple mechanism to avoid a state path, such as iiiimmmmiii or **oooommmoooo**, that does not correspond to the real topology of transmembrane protein, in which a membrane domain has to be flanked by i on one side and oon the other side. Therefore, the transition matrix is 4 by 4, corresponding to the four states. Note that the amino acid emission frequencies for the transmembrane state are calculated by lumping together counts or expectation from both M and $^{27}\mathrm{M}^{\circ}$ states. We set up two different experiments with different initial conditions: (1) ²⁸Transition matrix has correct zeros as ground truth model. (2) Transition matrix 29 is fully connected. We set up experiments for condition (2) because the method in [6] cannot enforce initial zeros to remain zeros during the training, therefore, condition (2) gives more fair comparison of the two methods when no prior knowledge is available. The HMM is trained by these two different methods in a 10-fold 32 cross validation scheme. Different levels of unlabelled data in training examples are

Li et al. Page 13 of 21

¹ actuated by selecting locations randomly to be unlabelled for each sequence. Since ¹
2 no ground truth model is available, maximum likelihood method with fully labelled 2
3 training data is used to mimic the role of the ground truth model in experiments 3
⁴ with synthetic. ⁴
$^{5}~$ For condition (1), the result shown in Figure 7 demonstrates that our method (5
$^6\mathrm{constrained}$ Baum-Welch with model selection) outperforms other method (Schef- 6
$^7 {\rm fer}$ et al) by 33.59% with Viterbi Algorithm and 36.16% with Posterior-Viterbi 7
8 algorithm. For condition (2), the result shown in Figure 8 attests that our method 8
$^9\mathrm{outperforms}$ other method by 33.20% with Viterbi Algorithm and 36.32% with 9
$^{10}\mathrm{Posterior\text{-}Viter bi}$ algorithm. For both conditions, the performance of our method 10
$^{11}\mathrm{with}$ or without model selection technique and maximum likelihood are very close. $^{12}\mathrm{cm}$
12
¹³ Discussion
$^{14}\mathrm{From}$ the results of experiments with synthetic data in Table 2 to Table 5, they show: $^{14}\mathrm{From}$
$^{15}(1)$. constrained Baum-Welch algorithm with or without model selection technique
16 achieve significant better performance than Scheffer et al in [6]; (2). constrained
$^{17}\mathrm{Baum\text{-}Welch}$ benefit from having correct topology (comparisons between the $4\mathrm{th}^{17}$
$^{18}\mathrm{columns}$ of Table 2 and Table 3); (3). constrained Baum-Welch algorithm performs 18
$^{19}\mathrm{better}$ when model selection technique used, especially when the task is hard (15
20 comparisons between 2nd and 4th column in Tables); (4). disregarding the training 20
$^{21}\mathrm{methods},$ Posterior-Viterbi always outperforms standard Viterbi (Shown in Figure 2^{21}
²² to Figure 4, Figure 7, and Figure 8).
From the results of experiments with real data, performance of constrained Baum- 23
24 Welch with or without model selection are very close to maximum likelihood ap- 24
25 proach across different percentages of partial label. However, the performance of 25
$^{26}\mathrm{Scheffer}$ et al's drops dramatically after the percentage of unlabelled data is greater 26
27 than 10. The reason behind this is the method done by Scheffer et al cannot enforce
28 the correct topology even the initial guess is correct. For this problem in particular, 28
have a HMM with correct topology is key for higher accuracy.
Moreover, there are a few points worth mentioning for the benefits of those who
31 may consider using this method for their applications. First, the ability of keeping 32
32 correct topology makes cBW method compatible with more complex HMM, such as 32
profile HMMs. However, as a trade-off, the training time can significantly increase.

Li et al. Page 14 of 21

Second, model selection technique, although optional, is highly recommended to	be [*]
² used with Posterior-Viterbi instead of standard Viterbi for best decoding perfo	or- ²
³ mance. Lastly, our method is designed especially for the task of detecting de no	vo ³
⁴ targeting signals, which assumes no fully labelled sequence is available in general	$al.^4$
⁵ For the cases with relaxing constrain: some fully labelled sequences are available	le, ⁵
⁶ our method is not the only choice, interested readers may also consider method	ds^6
⁷ in [9].	7
8	8
₉ Conclusion	9
10In this work, by modifying the standard Baum-Welch algorithm, we developed	aıc
11 novel training method, which, along with a model selection scheme, enables levera	g-11
12ing the partial labels in the data to improve the training of hidden Markov mode.	ls.12
13Compared with a similar method, our method achieved significant improvemen	ntsia
14in training hidden Markov models as evidenced by better performance in decoding	
15both synthetic data and the real biological sequence data.	15
For future work, we will further investigate the impact of this training method	
170n detecting de novo motifs and signals in biological data. In particular, we pla	
18to deploy the method in active learning mode to the ongoing research in detecting	
19plasmodesmata targeting signals and assess the performance with validations fro	
20wet-lab experiments.	20
21 22 ^{Abbreviations}	21
22 ^{-ABD Collino} HMM: Hidden Markov model	22
23PDLP: Plasmodesmata-located proteins	23
cBW: constrained Baum-Welch algorithm 24 EM: Expectation-Maximization	24
25 Declarations	25
26Competing interests	26
The authors declare that they have no competing interests. 27	27
28Consent for publication	28
Not applicable.	29
³⁰ Ethics approval and consent to participate	30
Not applicable. 31	31
	32
³² Acknowledgements The authors are grateful for the approximate reviewers' valuable comments and suggestions, in particular for bringing	

the Posterior-Viterbi decoding to their attention.

Li et al. Page 15 of 21

thor's contributions	1
L and LL designed the project, JL and LL devised algorithms, and JL implemented algorithms and carried out the	ne ₂
periments with advice from JYL and LL. All authors have read and approved the manuscript.	•
	3
nding	4
e work is funded by National Science Foundation NSF-MCB1820103. The funding agency had no role in the	
sign, collection, analysis, data interpretation and writing of this study.	5
	6
ailability of data and materials	
tasets and source code are freely available on the web at	7
tps://www.cis.udel.edu/~lliao/partial-label-HMMs	8
thor details	9
niversity of Delaware, Computer and Information Sciences, 101 Smith Hall, Newark, DE 19716 US. ² University of	
	10
	11
Discovery Blvd, 19713 Newark, US. ⁵ University of Delaware, Quantitative Biology, Düsternbrooker Weg 20,	11
716 Newark, US.	12
ferences	13
Baum, L.E., Petrie, T.: Statistical inference for probabilistic functions of finite state markov chains. The annal	ls .
of mathematical statistics 37 (6), 1554–1563 (1966)	14
Baum, L.E., Eagon, J.A.: An inequality with applications to statistical estimation for probabilistic functions of	15
markov processes and to a model for ecology. Bulletin of the American Mathematical Society 73(3), 360–363	
(1967)	16
Baum, L.E., Sell, G.: Growth transformations for functions on manifolds. Pacific Journal of Mathematics 27(2	!), ₁₇
	. 18
	f
	19
	00
	20
	21
	te 23
	24
	25
· · · · · ·	
	26
- ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' ' '	27
	-
	28
	inoc
	28
	30
	31 es
	es 32
prediction of the topology of all-heta membrane proteins RMC high-formatics 6(4) 1–7 (2005)	33
	eriments with advice from JYL and LL. All authors have read and approved the manuscript. Inding e work is funded by National Science Foundation NSF-MCB1820103. The funding agency had no role in the ign, collection, analysis, data interpretation and writing of this study. aliability of data and materials tasets and source code are freely available on the web at tape://www.cis.udel.edu/~111ao/partial~1abel~1899s thor details niversity of Delaware, Computer and Information Sciences, 101 Smith Hall, Newark, DE 19716 US. ² University laware, Plant and Soil Sciences, 15 Innovation Way, 19716 Newark, USA. ³ University of Delaware, Delaware technology Institute, 15 Innovation Way, 19716 Newark, USA. ⁴ University of Delaware, Delaware Delaware, 20 Discovery Blvd, 19713 Newark, US. ⁵ University of Delaware, Data Science Institute, 20 Discovery Blvd, 19713 Newark, US. ⁵ University of Delaware, Quantitative Biology, Düsternbrooker Weg 20, 716 Newark, US. Ferences Baum, L.E., Petrie, T.: Statistical inference for probabilistic functions of finite state markov chains. The anna of mathematical statistics 37(6), 1554-1563 (1966) Baum, L.E., Eagon, J.A.: An inequality with applications to statistical estimation for probabilistic functions of markov processes and to a model for ecology. Bulletin of the American Mathematical Society 73(3), 360-363 (1967) Baum, L.E., Petrie, T., Soules, G., Weiss, N.: A maximization technique occurring in the statistical analysis o probabilistic functions of markov chains. The annas of mathematical statistics 41(1), 164-171 (1970) Baum, L.E., An inequality and associated maximization technique in statistical estimation of probabilistic functions of a markov process. Inequalities 3, 1-8 (1972) Baum, L.: An inequality and associated maximization technique in statistical estimation of probabilistic functions of a markov process. Inequalities 3, 1-8 (1972) Baum, L.: An inequality and associated maximization technique in statistical estimation of probabilistic functions of a markov pro

Li et al. Page 16 of 21

Kahsay, R.Y., Gao, G., Liao, L.: An improved hidden markov model for transmembrane protein detection and topology prediction and its applications to complete genomes. Bioinformatics 21(9), 1853–1858 (2005)
 Sonnhammer, E.L., Von Heijne, G., Krogh, A., et al.: A hidden markov model for predicting transmembrane helices in protein sequences. (1998)
 Käll, L., Krogh, A., Sonnhammer, E.L.: Advantages of combined transmembrane topology and signal peptide prediction—the phobius web server. Nucleic acids research 35(suppl.2), 429–432 (2007)
 Hayat, S., Peters, C., Shu, N., Tsirigos, K.D., Elofsson, A.: Inclusion of dyad-repeat pattern improves topology 5 prediction of transmembrane β-barrel proteins. Bioinformatics 32(10), 1571–1573 (2016)

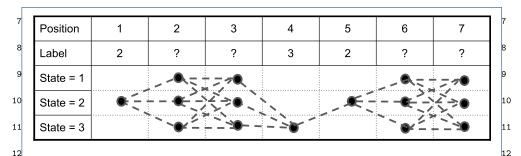


Figure 1 Example of constrained Baum-Welch's forward/backward dynamic programming table construction Position 1 to 2 is the case of labelled position to unlabelled position. Position 3 to 4 is the case of unlabelled position to labelled position. Position 4 to 5 is the case of labelled position to labelled position to labelled position. Dashed lines indicate state transitions. Generated by Google Drawings

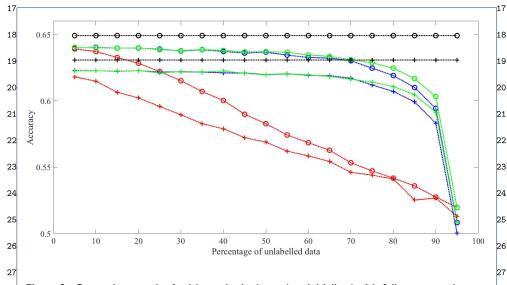
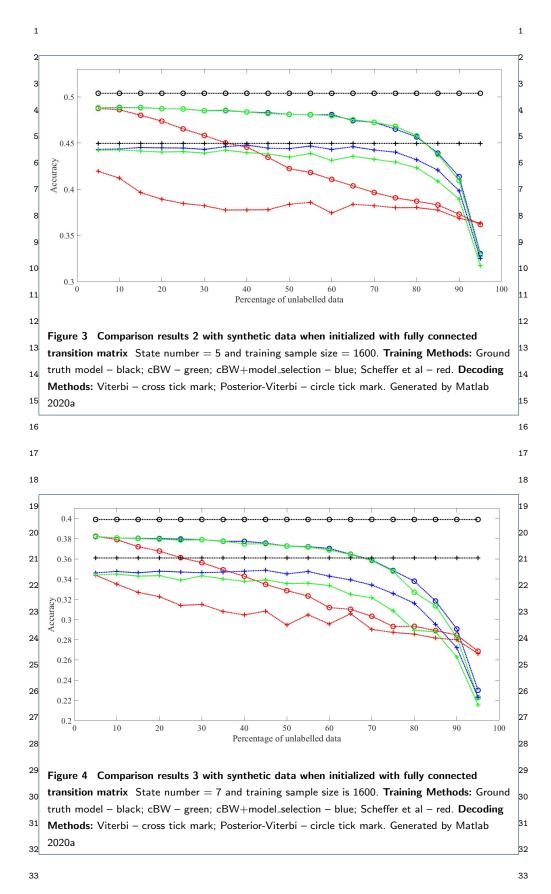
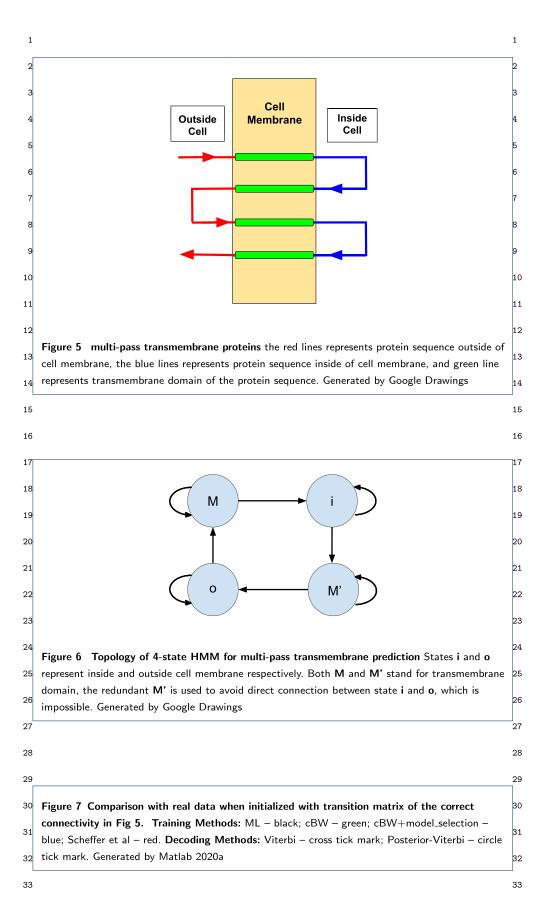


Figure 2 Comparison results 1 with synthetic data when initialized with fully connected transition matrix state number = 3 and training sample size = 1600. Training Methods: Ground truth model – black; cBW – green; cBW+model_selection – blue; Scheffer et al – red. Decoding Methods: Viterbi – cross tick mark; Posterior-Viterbi – circle tick mark. Generated by Matlab 2020a

Li et al. Page 17 of 21



Li et al. Page 18 of 21



Page 19 of 21 Li et al.

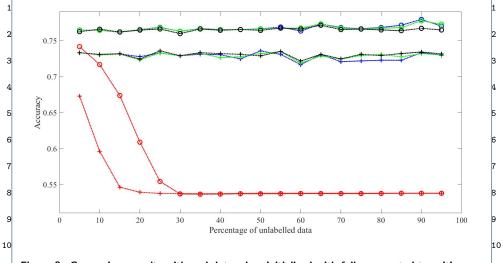


Figure 8 Comparison results with real data when initialized with fully connected transition matrix. Training Methods: ML - black; cBW - green; cBW+model_selection - blue; Scheffer et al – red. Decoding Methods: Viterbi – cross tick mark; Posterior-Viterbi – circle tick mark. Generated by Matlab 2020a

 $_{14}$ Table 1 Notations

14		
	Symbols	Explanations
15	θ	Hidden Markov model: $\theta = (\pi, A, B)$
16	N	States' number in hidden Markov model.
	K	Symbolic Number in hidden Markov model.
17	A	Transition matrix with dimension $N \times N$.
18	a_{ij}	Probability of state i transition to state j.
	В	Emission matrix with dimension $N \times K$.
.9	$b_j(k)$	Probability of state j emitted from symbol k.
20	π	Initial probability of states with dimension $N \times 1$.
	O^s	The s^{th} sequence with length T^s
21	X^s	State sequence of ${\cal O}^s$
22	m	Total number of sequences

 $\textbf{Table 2} \ \ \textbf{Improvements of cBW} + \textbf{model selection, cBW alon vs Scheffer et al, with Fully Connected}$ $^{\rm 23}$ Initial Transition Matrix for Synthetic Data with Viterbi Algorithm

24	state # /	Average improvements	Average p-value of	Average improvements	24
25	training	of cBW+model selection	cBW+model selection	of cBW alone	25
25	sample	in setting 1 / 2	in setting 1 / 2	in setting 1 / 2	25
26	3 /1600	7.35% / 8.29%	2.1E-02 / 6.2E-05	7.61% / 8.49%	26
	3 /2000	7.99% / 8.82%	3.9E-02 / 2.6E-09	8.31% / 9.00%	
27	3 /2400	8.47% / 9.13%	2.8E-03 / 2.4E-10	8.71% / 9.18%	27
28	3 /2800	8.51% / 9.24%	5.8E-03 / 6.9E-10	8.65% / 9.24%	28
	5 /1600	12.97% / 14.75%	7.8E-05 / 3.3E-05	11.13% / 12.82%	
29	5 /2000	14.63% / 16.32%	2.5E-03 / 1.2E-06	12.65% / 14.11%	29
30	5 /2400	14.69% / 16.54%	7.8E-04 / 6.2E-06	12.73% / 14.50%	30
	5 /2800	15.56% / 17.22%	1.6E-02 / 1.3E-07	13.72% / 15.22%	
31	7 /1600	8.61% / 10.42%	3.5E-02 / 1.1E-02	5.56% / 7.20%	31
32	7 /2000	10.56% / 12.40%	6.4E-03 / 6.2E-03	7.87% / 9.52%	32
	7 /2400	11.35% / 13.21%	8.2E-03 / 7.8E-03	8.71% / 10.43%	
33	7 /2800	12.16% / 14.06%	1.0E-03 / 1.1E-04	9.68% / 11.41%	33

Li et al. Page 20 of 21

 $^{
m 3}$ Table 3 Improvements of cBW + model selection, cBW alone vs Scheffer et al in Cases of Correct 4Initial Transition Matrix for Synthetic Data with Viterbi Algorithm

5	state # /	Average improvements	Average p-value of	Average improvements
	training	of cBW $+$ model selection	cBW + model selection	of cBW alone
6	sample	in setting $1\ /\ 2$	in setting $1\ /\ 2$	in setting 1 / 2
7	3 /1600	8.07% / 8.57%	2.3E-04 / 4.7E-07	8.11% / 8.56%
•	3 /2000	8.58% / 9.05%	1.9E-06 / 6.5E-09	8.63% / 9.03%
8	3 /2400	8.93% / 9.24%	1.6E-07 / 1.0E-09	8.97% / 9.20%
9	3 /2800	8.87% / 9.31%	1.3E-08 / 1.5E-09	8.94% / 9.26%
9	5 /1600	11.99% / 13.24%	1.7E-02 / 4.5E-06	11.76% / 13.08%
10	5 /2000	13.07% / 14.20%	4.1E-02 / 3.4E-06	12.87% / 14.11%
11	5 /2400	13.22% / 14.59%	2.0E-02 / 1.2E-05	12.94% / 14.35%
11	5 /2800	13.89% / 15.20%	4.1E-02 / 1.6E-07	13.85% / 15.16%
12	7 /1600	7.85% / 9.37%	6.7E-02 / 3.5E-02	6.04% / 7.34%
10	7 /2000	9.75% / 11.28%	5.6E-03 / 4.1E-03	7.93% / 9.32%
13	7 /2400	10.50% / 12.10%	1.8E-02 / 1.9E-02	8.99% / 10.53%
14	7 /2800	11.39% / 12.95%	1.4E-03 / 1.9E-04	9.75% / 11.29%

20 Table 4 Improvements of cBW + model selection, cBW alon vs Scheffer et al, with Fully Connected 20 Initial Transition Matrix for Synthetic Data with Posterior-Viterbi Algorithm

	state # /	Average improvements	Average p-value of	Average improvements
22	training	of cBW+model selection	cBW+model selection	of cBW alone
23	sample	in setting 1 / 2	in setting 1 / 2	in setting 1 / 2
23	3 /1600	7.08% / 8.02%	3.6E-02 / 1.9E-04	7.49% / 8.35%
24	3 /2000	7.93% / 8.57%	9.5E-03 / 3.9E-05	8.32% / 8.88%
25	3 /2400	8.21% / 8.85%	1.7E-03 / 2.5E-05	8.55% / 9.03%
25	3 /2800	8.43% / 9.15%	1.8E-03 / 1.1E-06	8.82% / 9.36%
26	5 /1600	9.13% / 10.62%	2.4E-02 / 3.2E-03	9.08% / 10.58%
07	5 /2000	10.32% / 11.68%	1.4E-02 / 6.2E-05	10.38% / 11.76%
27	5 /2400	11.26% / 12.74%	1.1E-02 / 2.0E-06	11.29% / 12.84%
28	5 /2800	12.48% / 13.76%	9.8E-03 / 2.1E-08	12.56% / 13.86%
	7 /1600	8.40% / 10.08%	6.0E-02 / 2.5E-02	7.77% / 9.50%
29	7 /2000	10.22% / 11.96%	3.2E-02 / 1.3E-04	9.82% / 11.65%
30	7 /2400	11.10% / 12.68%	8.1E-03 / 1.5E-05	10.60% / 12.31%
	7 /2800	12.18% / 13.89%	1.6E-04 / 8.4E-08	11.96% / 13.77%

Li et al. Page 21 of 21

1					1
2					2
3					3
4					4
5					5
6					6
7					7
8					8
9					9
10					10
11					11
Tab	le 5 Improver	ments of cBW $+$ model selec	ction, cBW alone vs Schef	fer et al in Cases of Correct	
12		Matrix for Synthetic Data w			12
13	state # /	Average improvements	Average p-value of	Average improvements	13
14	training	of cBW+model selection	cBW+model selection	of cBW alone	14
	sample	in setting 1 / 2	in setting 1 / 2	in setting 1 / 2	
15	3 /1600	7.46% / 8.01%	2.8E-02 / 1.7E-02	7.64% / 8.11%	15
16	3 /2000	8.09% / 8.52%	3.0E-02 / 1.6E-02	8.25% / 8.60%	16
10	3 /2400	8.32% / 8.67%	3.7E-02 / 5.9E-03	8.49% / 8.73%	10
17	3 /2800	8.58% / 8.99%	4.9E-02 / 1.2E-03	8.79% / 9.09%	17
	5 /1600	9.52% / 10.62%	1.1E-01 / 5.1E-02	9.45% / 10.59%	
18	5 /2000	10.53% / 11.55%	6.0E-02 / 8.2E-03	10.47% / 11.63%	18
19	5 /2400	11.51% / 12.58%	2.0E-02 / 3.5E-04	11.44% / 12.58%	19
	5 /2800	12.49% / 13.55%	1.7E-02 / 8.0E-06	12.55% / 13.61%	
20	7 /1600	8.75% / 10.19%	2.9E-02 / 2.6E-02	8.27% / 9.77%	20
21	7 /2000	10.50% / 11.99%	2.1E-02 / 4.2E-03	9.82% / 11.31%	21
	7 /2400	11.15% / 12.47%	6.9E-02 / 7.8E-04	10.69% / 12.11%	
22	7 /2800	12.36% / 13.75%	5.2E-02 / 1.1E-05	12.05% / 13.57%	22
23	,	,	,	,	23
24					24
25					25
26					26
27					27
28					28
29					29