
C³TS: CONFIDENCE-GUIDED LEARNING PROCESS FOR CONTINUOUS CLASSIFICATION OF TIME SERIES *

Chenxi Sun

Key Laboratory of Machine Perception
(Ministry of Education), Peking University
Beijing, China
School of Artificial Intelligence,
Peking University
Beijing, China
sun_chenxi@pku.edu.cn

Derun Cai

Key Laboratory of Machine Perception
(Ministry of Education), Peking University
Beijing, China
School of Artificial Intelligence,
Peking University
Beijing, China
cdr@stu.pku.edu.cn

Shenda Hong*

National Institute of Health Data Science,
Peking University
Beijing, China
Institute of Medical Technology,
Health Science Center of Peking University
Beijing, China
hongshenda@pku.edu.cn

Moxian Song

Key Laboratory of Machine Perception
(Ministry of Education), Peking University
Beijing, China
School of Artificial Intelligence,
Peking University
Beijing, China
songmoxiani@pku.edu.cn

Baofeng Zhang

Key Laboratory of Machine Perception
(Ministry of Education), Peking University
Beijing, China
School of Artificial Intelligence,
Peking University
Beijing, China
boffinzhang@stu.pku.edu.cn

Hongyan Li*

Key Laboratory of Machine Perception
(Ministry of Education), Peking University
Beijing, China
School of Artificial Intelligence,
Peking University
Beijing, China
leehy@pku.edu.cn

ABSTRACT

In the real world, the class of a time series is usually labeled at the final time, but many applications require to classify time series at every time. e.g. the outcome of a critical patient is only determined at the end, but he should be diagnosed at all times to facilitate timely treatment. Thus, in this paper, we propose a new concept: Continuous Classification of Time Series (CCTS). There are two open problems about CCTS: (1) Data arrangement. Time series is a kind of dynamic data. It evolves multiple distributions over time. The division of multi-distribution will directly affect the classification accuracy; (2) Model training strategy. When a model learns multi-distributed data, it always forgets the old distribution or overfits in the main distribution. Different data learning orders will result in different model performances. We find that the process of model learning multiple distributions can be similar to the process of human learning multiple knowledge. Thus, we propose a novel Confidence-guided method for CCTS to arrange the data and schedule the training, named C³TS. It imitates the objective-confidence and the alternating self-confidence of humans in their

**Citation*: Chenxi Sun, Moxian Song, Derun Cai, Baofeng Zhang, Shenda Hong, and Hongyan Li. 2022. Confidence-Guided Learning Process for Continuous Classification of Time Series. In Proceedings of the 31st ACM International Conference on Information and Knowledge Management (CIKM '22), October 17–21, 2022, Atlanta, GA, USA. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3511808.3557565>

learning process, which is described by the Dunning-Kruger Effect. Specifically, we define an importance-based objective-confidence to arrange and replay data, and design an uncertainty-based self-confidence to control the training duration. Experimental results on four real-world datasets show that our method is more accurate than all baselines at every time.

Keywords Continuous Classification of Time Series · Confidence · Neural Network Training

1 Introduction

In the real world, many applications need to classify time series at every time. For example, in the Intensive Care Unit (ICU), most detected vital signs change dynamically with the evolution of diseases. The status perception is needed at any time as the real-time diagnosis provides more opportunities for doctors to rescue lives [1]. But patient labels, e.g., mortality and morbidity, are unknown in early stages but available only at the onset time. In response to the current demand, we propose a new concept – Continuous Classification of Time Series (CCTS), to classify time series at every early time before the final labeled time, as shown in Figure 1 left.

CCTS task requires the method to learn data from different advanced stages. But for most practical time series, changed data characteristics lead to the evolved data distribution, and finally produce the data form of multi-distribution. As shown in Figure 1 middle, data distributions of blood pressure of sepsis patients vary among early, middle, and later stages, bringing a triple-distribution. Thus, to achieve CCTS, the method should model such multi-distributed data. And there are two open problems:

How to prepare multi-distribution before training the model? As CCTS requires the mode of continuous classification, all subsequences of a time series had better be learned. The intuition is to get multi-distribution according to time stages. But the optimal time interval is unclear: small time stages produce overlapping distributions while large time stages produce distinct distributions. They all affect model learning by worsening the forgetting or overfitting [2]. Some methods define distributions by the data complexity [3], e.g., the number of objects [4] or rare words [5]. But compared with the data form of image and text, time series is more abstract, the complexity is hard to tell. Besides, data that is difficult for people may not be difficult for machines. e.g., human-identifiable stable vital signs are more likely to confuse the model [6].

How to train the model under multi-distribution? A single model, like deep neural network, is lack of ability to learn all distributions simultaneously [7] as they are restricted by the premise of independent identically distributed (i.i.d) data [8]. As shown in Figure 3, frequent learning of new knowledge will inevitably lead to the forgetting of old ones [2], and too much training on one distribution may make the parameter fall into the local solution, resulting in poor generalization [9]. Besides, a model will receive different parameter matrices according to different training orders [10]. As shown in Figure 4, training a model from early time series to late and from late time series to early will lead to different result accuracy and convergence speed. Thus, feeding examples in a meaningful order but randomly is critical [10].

To solve both two problems, we propose a novel Confidence²-guided method for CCTS, named C³TS. As we assume that the process of model learning multiple distributions can be similar to the process of human learning multiple knowledge: people will first arrange the learning order (data arrangement), then control the progress of learning and review according to their mastery (model learning). The mastery of knowledge is usually assessed by human confidence, they will relearn the knowledge that they are not confident about and improve its priority in the learning order. This human behavior is the Dunning-Kruger effect [11, 12] as shown in Figure 1 right. It is an alternating process with experiences of ignorance, overconfidence, disappointment and development:

When a human learns a new field of knowledge, he will first scratch and grasp the overall framework and therefore has a lot of confidence. Then, he begins to find that he really had little in-depth knowledge and loses his confidence. Over time, he studies deeply, becoming more and more experienced and confident [13].

To imitate human confidence-guided learning, we define the confidence from two aspects of objective-confidence and self-confidence. Specifically, we design the importance-based objective-confidence by using the importance coefficient. It arranges data and makes the model relearn important samples to consolidate memory loosely; We design the uncertainty-based self-objective by defining the total uncertainty. It controls the training duration under the current distribution and schedules the training order among distributions. Experimental results on four real-world datasets show that C³TS is more accurate than all baselines at every time.

²The confidence in this paper represents the human cognition of their ability, different from the confidence in Statistics, the probability of the result in the confidence interval.

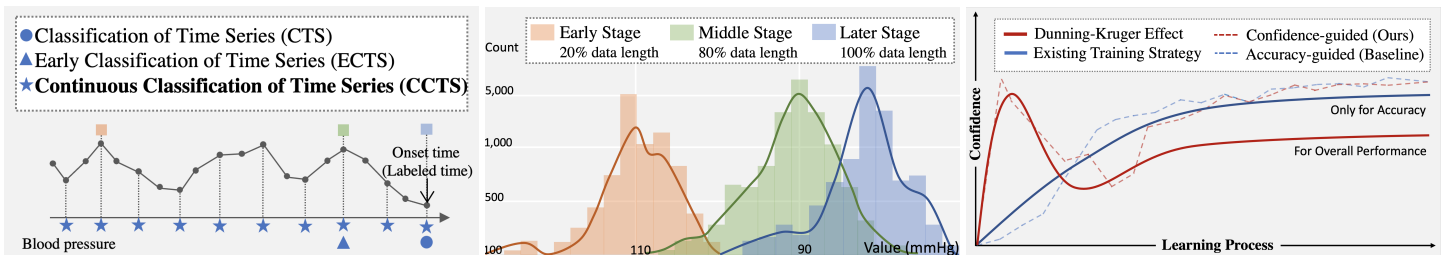


Figure 1: Continuous Classification of Time Series (CCTS) & Multi-distribution in Time Series & Confidence Change during Learning

The left figure shows that different from existing classification tasks of time series (CTS and ECTS), CCTS needs the continuous classification mode; The middle figure shows that the statistics [14] of blood pressure of 2,000 sepsis patients represent multiple distributions among early, middle and later stages; The right figure shows the confidence changes of human with an alternating process: ignorance→overconfidence→disappointment→development.

2 Related Work

The popularity of time series classification has attracted increasing attention in many practical fields [15]. The existing work can be summarized into two categories. And more detailed work and concept distinction are in Appendix.

2.1 One-shot classification: classifying at a fixed time

The foundation is Classification of Time Series (CTS), classifying the full-length time series [16]. But in time-sensitive applications, Early Classification of Time Series (ECTS) is more critical, making classification at an early time [6]. For example, early diagnosis helps for sepsis outcomes [17]. However, both CTS and ECTS give the one-shot classification, only classifying once and just lean a single data distribution. They have good performances on i.i.d data at a fixed time, like early 6 hours sepsis diagnosis [18], but fail for multi-distribution. In fact, continuous classification is composed of multiple one-shot classifications as shown in Figure 1 left.

2.2 Continuous classification: classifying at every time

To achieve CCTS, learning multi-distribution is essential. Most methods use multi-model to model multi-distribution, like SR [19] and ECEC [20]. They divide data according to time stages and design different classifiers for different distributions. But they only focus on the data division with no analysis or experiment to show the rationality. Besides, the operation of classifier selection in multi-model will result in additional losses.

A model will face the catastrophic forgetting problem when learning multi-distributed data. Recently, continual learning [7] methods aim to address the issue of static models incapable of adapting their behavior for new knowledge. It learns a new task at every moment and each task corresponds to one data distribution. Replay-based methods [14, 21] retrain the model with old data to consolidate memory and focus on the storage limitation. But in their settings, the old task and new task are clear so that the multi-distribution is fixed. This is different from the scenario of CCTS, where the distributions are not determined and need to be divided and defined firstly. Besides, one continual learning contains multiple tasks, their data distributions are not particularly similar. But in CCTS, it aims at one task, thus the replay easily causes the over-fitting problem.

Training a model by multiple distributions should also consider the training order. A subdiscipline named curriculum learning [22], where an instinctive strategy presents an easy-to-hard order, referring to human courses from simple to difficult. Many methods define the difficulty through the complexity of data, e.g., the number of object in an image [4], the number of rare words in a sentence [5]. But compared with the data form of image and text, time series is more abstract, and the complexity is hard to define. Some features, such as Shaplets [23] and frequency [24], merely consider one aspect of time series which fails to fully cope with the data difficulty for a model. Besides, the difficulty of human definition may not match the machine. For example, stable vital signs are more likely to confuse the model and lead to classification errors [6].

The existing human behavior-imitated methods assume that the confidence monotonically increases in the training process [25], which does not conform to the Dunning-Kruger effect [11, 12]. As ignorance more frequently begets confidence³, the increasing confidence may not lead to growing knowledge.

³Quote from Charles Robert Darwin (1809-1882)

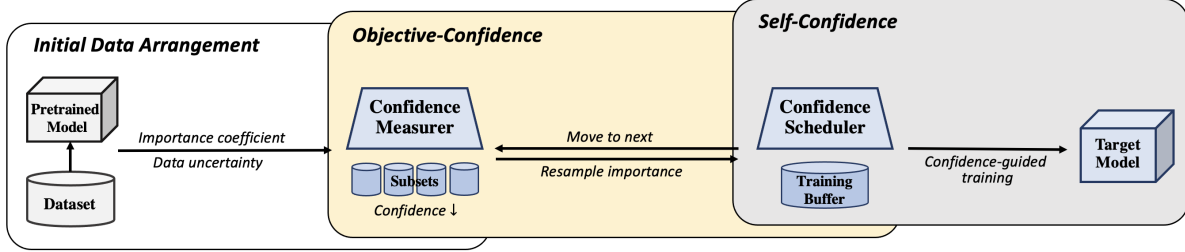


Figure 2: Confidence-Guided Training Process for Continuous Classification of Time Series (C³TS)

The uncertainty [26] is similar to the confidence [27]. It can measure if the model parameters could describe the data distribution [28]. The uncertainty can be quantified by using the variance of probabilistic distribution of model parameters [12]. Some researchers design more general methods, using dropout mechanism [29] and prior networks [30] to assist in the evaluation. Some curriculum learning works use data uncertainty to arrange the initial training order, like data competence [31, 32] and joint difficulty [12]. In CCTS, in addition to the uncertainty, we need to combine the feature of time series data form, avoid forgetting problem to design the model confidence, and consider both distribution order and learning duration.

3 Problem Formulation

Without the loss of generality, we use the univariate time series to present the problem. Multivariate time series can be described by changing x_t to x_t^i . i is the i -th dimension.

Definition 1 (Continuous Classification, CC). *A time series $X = \{x_1, \dots, x_T\}$ is labeled with a class $C \in \mathcal{C}$ at the final time T . CC classifies X at every t with loss $\sum_{t=1}^T \mathcal{L}(f(X_{1:t}), C)$.*

Note that unlike continuous classification, CTS and ECTS are just one-shot classification, where they optimize the objective with a single loss $\mathcal{L}(f(x), c)$. Our CCTS should consider the multi-distribution and classify more times. Thus, we define CCTS as:

Definition 2 (Continuous Classification of Time Series CCTS). *A dataset \mathcal{S} contains N time series. Each time series $X = \{x_1, \dots, x_T\}$ is labeled with a class $C \in \mathcal{C}$ at the final time T . As time series varies among time, it has a subsequence series with M different distributions $\mathcal{D} = \{\mathcal{D}^1, \dots, \mathcal{D}^M\}$, each \mathcal{D}^m has subsequence $X_{1:t^m}$. CCTS learns every \mathcal{D}^m and introduces a task sequence $\mathcal{M} = \{\mathcal{M}^1, \dots, \mathcal{M}^M\}$ to minimize the additive risk $\sum_{m=1}^M \mathbb{E}_{\mathcal{M}^m} [\mathcal{L}(f^m(\mathcal{D}^m; \theta), C)]$ with model f and parameter θ . f^m is the model f after being trained for \mathcal{M}^m . When the model is trained for \mathcal{M}^m , its performance on all observed data cannot degrade: $\frac{1}{m} \sum_{i=1}^m \mathcal{L}(f^i, \mathcal{M}^i) \leq \frac{1}{m-1} \sum_{i=1}^{m-1} \mathcal{L}(f^i, \mathcal{M}^i)$.*

4 Confidence-guided Method

We introduce the confidence from objective-confidence and self-confidence. They are achieved by the importance coefficient and the uncertainty during the model training process.

4.1 Objective-Confidence — Replay by Importance Coefficient

People gain the objective-confidence through regular examinations or tests, where they can know their weak knowledge through test results and learn it again. Thus, we use the classification accuracy as the test results of the model and apply the importance coefficient to find the weak knowledge of the model and replay it again to improve the objective-confidence of the model.

We focus on the adaptive method to explore a wider space, where the replayed data is dynamic and determined according to the current state. We introduce an importance-based replay method. In each round, it only re-trained the model by some important samples. The importance of each sample is learned from the importance coefficient in the objective of an additive loss function.

The model learns the importance of a time series X_i to it by the coefficient β_i of X_i 's loss \mathcal{L}_i . The overall loss is the sum of the loss of each sample in current distribution \mathcal{D} :

$$\mathcal{L} = -\frac{1}{|\mathcal{D}|} \sum_{X_i \in \mathcal{D}} \beta_i^2 (-\mathcal{L}_i) + \lambda(\beta_i - 1)^2 \quad (1)$$

The importance coefficient β_i of each sample X_i is learned/updated by the gradient descent $\beta_i \leftarrow \beta_i - \frac{\partial \mathcal{L}}{\partial \beta_i}$. If a sample X_i is hard to classify, its loss \mathcal{L}_i will be larger and $-\mathcal{L}_i$ will be smaller. In order to minimize the overall loss, its β_i will be larger. Based on this, the important samples for the model are those difficult to learn with $\beta > \epsilon$. Meanwhile, as β is the coefficient of loss, if $\beta = 0$, the loss are hard to be optimized. Thus, inspired by [14], we introduce a regularization term $(\beta - 1)^2$ and initialize $\beta = 1$ to penalize it when rapidly decaying toward 0.

4.2 Self-Confidence — Approximate by Uncertainty Evaluation

People gain the current self-confidence through their abilities and the problem difficulty. Correspondingly, we approximate the self-confidence of the model by defining the total uncertainty according to the model uncertainty and the data uncertainty.

The high self-confidence of the model shows that the model believes it already has good classification ability, which means the uncertainty is low. A small score of total uncertainty indicates the model is confident that the current training data has been well learned, and the termination of the decline in scores represents the signal to shift to the next training stage.

$$Confidence = \frac{1}{U_{total}} \quad (2)$$

The existing work shows that the total uncertainty = epistemic uncertainty + aleatoric uncertainty [30]:

$$\underbrace{\mathcal{H}[\mathbb{E}_{P(\theta|\mathcal{D})}[P(y|x^*, \theta)]]}_{Total} = \underbrace{\mathcal{I}[y, \theta|x^*, \mathcal{D}]}_{Epistemic} + \underbrace{\mathbb{E}_{P(\theta|\mathcal{D})}[\mathcal{H}[P(y|x^*, \theta)]]}_{Aleatoric}$$

Inspired by this idea, we define the total uncertainty as:

$$U_{total} = U_{model} + U_{data} \quad (3)$$

First, we define the model uncertainty as the variance σ of a distribution of K classification probabilities for data X . The classification results are predicted by the model with K different parameter disturbances. For reasons of computational efficiency, we adopt widely used Monte Carlo Dropout [29] to approximate Bayesian inference and achieve the parameter disturbances. For the current distribution $\mathcal{D}^m = \{X_{1:t^m}, C\}$, the model uncertainty of model $f(\theta)$ is:

$$U_{model}(f(\theta), \mathcal{D}^m) = \frac{1}{|\mathcal{D}^m|} \sum_{(x_{1:t^m}, c) \in \mathcal{D}^m} \sigma_{k=1}^K(p(c|x_{1:t^m}, \theta_k)) \quad (4)$$

Second, we define a novel importance-aware data uncertainty function. It is based on the entropy of the predictive distribution [33]. It behaves similar to max probability, but represents the uncertainty encapsulated in the entire distribution. We also inject the learned importance coefficient β of the data x to data uncertainty. Thus, the most uncertain data for the model is that has the largest combination value of entropy and importance. The uncertainty of a data x is:

$$U_{data}(x) = \frac{\beta}{|\mathcal{C}|} \sum_{c \in \mathcal{C}} p(c|x) \log(1 - p(c|x)) \quad (5)$$

Finally, according to the model uncertainty and the data uncertainty, we can get the total uncertainty for the current model f^m :

$$\begin{aligned} U_{total}(f^m) &= U_{model}(f^m, \mathcal{D}^m) + U_{data}(\mathcal{D}^m) \\ &= \frac{1}{|\mathcal{C}|} \sum_{c \in \mathcal{C}, x \in \mathcal{D}^m} U_{model}(f^m, (x, c)) \\ &\quad + \frac{1}{|\mathcal{D}^m|} \sum_{x \in \mathcal{D}^m} U_{data}(x) \end{aligned} \quad (6)$$

Algorithm 1 C³TS

Input: Training set $\mathcal{S} = \{(X^i, C^i)\}_{i=1}^N$;
 An untrained target model f^1 ;
 An untrained same model as the target model f' .

Output: A well-trained target model f^M .

- 1: // DATA ARRANGEMENT (PRE-TRAIN)
- 2: Extend $\mathcal{S} \leftarrow \{(X_{1:t}, C)\}_{t=1}^T$
- 3: Train f' by \mathcal{S} using loss in Equation 1
- 4: Get U_{data} for each $X \in \mathcal{S}$ of f' using Equation 5
- 5: Split \mathcal{S} into M datasets $\mathcal{D} = \{\mathcal{D}^m\}_{m=1}^M$ by U_{data} and define M baby steps $\mathcal{M} = \{\mathcal{M}^m\}_{m=1}^M$ by \mathcal{D} with increasing U_{data}
- 6: //OBJECTIVE-CONFIDENCE (TASK SCHEDULING)
- 7: Initialize current training set buffer $\mathcal{B} = \mathcal{D}^1$
- 8: **for** $m = 1$ to M **do**
- 9: // SELF-CONFIDENCE (DURATION SCHEDULING)
- 10: **while** not early stop of Equation 2 **do**
- 11: Train f^m by \mathcal{B} using loss U_{total} in Equation 6
- 12: **end while**
- 13: Get β for each $X \in \mathcal{B}$ by f^m using Equation 1
- 14: Update $\mathcal{B} \leftarrow \{\mathcal{D}^{m+1}, \{X_i | \beta_i > \epsilon\}\}$
- 15: **end for**

4.3 Confidence-Guided Training Process

C³TS consists of three interrelated and cooperative modules: Initial data arrangement module, objective-confidence scheduling module, and self-confidence scheduling module.

- The *initial data arrangement module* gives an initial data learning order for the model based on confidence. It imitates the fact that before starting a study, students will arrange the learning order according to the knowledge difficulty.

It is a pre-train process as shown in the white part of Figure 2. We first expand the original data set \mathcal{S} by extracting T subsequences $\{X_{1:t}\}_{t=1}^T$ from each time series $X = \{x_1, \dots, x_T\}$. The subsequence $X_{1:t}$ is an early $T - t$ time data. We train the model from early $T - 1$ time dataset to early 0 time dataset with loss in Equation 1. Then, we can get the importance coefficient β of each sample in \mathcal{S} . The initial data learning order is based on the model's confidence in data, which is calculated by the importance-aware data uncertainty in Equation 5. After obtaining the U_{data} of each sample, we split \mathcal{S} into M datasets $\mathcal{D} = \{\mathcal{D}^m\}_{m=1}^M$ according to U_{data} and sort \mathcal{D} from small U_{data} to large U_{data} . So far, we get the data learning order for the model. The model is trained to converge on each distribution \mathcal{D}^m , which can be regarded as the model solving M tasks $\mathcal{M} = \{\mathcal{M}^m\}_{m=1}^M$.

- The *objective-confidence scheduling module* controls the overall learning process of each task \mathcal{M} and dataset \mathcal{D} . It determines the new data to learn and the old data to review based on the objective-confidence and aims to solve problems of catastrophic forgetting and overfitting. It imitates the fact that students will decide what to review based on their test scores.

This module is mainly controlled by the confidence measurer as shown in the yellow part of Figure 2. The confidence measurer organizes and arranges datasets $\mathcal{D} = \{\mathcal{D}^m\}_{m=1}^M$, which is initially obtained from data arrangement module. After the model learning the task \mathcal{M}^m with dataset \mathcal{D}^m , the confidence measurer will determine the dataset \mathcal{D}^{m+1} that the model will learn next. \mathcal{D}^{m+1} contains samples in new distribution and samples that should be re-learned. The confidence measurer finds samples to review according to objective-confidence in Equation 1. Sample with larger β need to be learned again. The partial relay will alleviate the catastrophic forgetting without causing the over fitting caused by over training.

- The *self-confidence scheduling module* controls the duration of each training stage with a task \mathcal{M}^m . It determines the training direction of the model on the dataset \mathcal{D}^m and the evidence of model convergence by the self-confidence method. It imitates the fact that students decide whether they have mastered the knowledge or not through their confidence in the current knowledge.

This module is mainly controlled by the confidence scheduler as shown in the gray part of Figure 2. The confidence scheduler determines whether the model converges on the current dataset \mathcal{D}^m . If the model converges, it will obtain

\mathcal{D}^{m+1} provided by the confidence measurer and let the model start learning the new task \mathcal{M}^{m+1} ; If not, it will continue training the model on the current dataset \mathcal{D}^m . The confidence scheduler judges whether the model converges according to Equation 2. Training stops when confidence is high. Specifically, like the early stop based on the loss [34], the model is stopped training if the self-confidence is no longer increased in several epochs.

The overall algorithm is in Algorithm 1. It first pre-train the model and get the initial \mathcal{M} and \mathcal{D} . Then, it trains the model from \mathcal{M}^1 to \mathcal{M}^M . For each \mathcal{M}^m , it controls the training duration based on the self-confidence. Between \mathcal{M}^m and \mathcal{M}^{m+1} , it dynamically adjusts \mathcal{D}^{m+1} according to the objective-confidence.

4.4 Complexity Analysis

Assuming that the computational complexity of updating a neural network model by one sample is $\mathcal{O}(d)$, then training a model by D time series with length T and E epochs usually costs $\mathcal{O}(TEDd)$. C^3TS contains a pre-training process and a training process. The pre-training process with loss in Equation 1 has the complexity of $\mathcal{O}(TEDd)$, being the same as the normal method. In the training process, C^3TS trains a model with N distributions \mathcal{D} with N training tasks \mathcal{M} , assuming there are E' epochs and S retained sample in each \mathcal{M}^n , the complexity will be $\mathcal{O}(NE'(D+S)d)$. The overall complexity is $\mathcal{O}((TE+NE')(2D+S)d)$. As $S \ll D$ and N is a small constant with $N < T$, the complexity of C^3TS approximates $\mathcal{O}(TEDd)$, almost being the same as the general training strategy.

5 Experiments

More detailed experiments and analyses are in Appendix.

5.1 Datasets

- UCR-EQ dataset [35] has 471 earthquake records from UCR time series database archive. It is the univariate time series of seismic feature value. Natural disaster early warning, like earthquake warning, helps to reduce casualties and property losses [36].
- USHCN dataset [37] has the daily meteorological data of 48 states in U.S. from 1887 to 2014. It is the multivariate time series of 5 weather features. Rainfall warning is not only the demand of daily life, but also can help prevent natural disasters [38].
- COVID-19 dataset [39] has 6,877 blood samples of 485 COVID-19 patients from Tongji Hospital, Wuhan, China. It is the multivariate time series of 74 laboratory test features. Mortality prediction helps for the personalized treatment and resource allocation [40].
- SEPSIS dataset [41] has 30,336 patients' records, including 2,359 diagnosed sepsis. It is the multivariate time series of 40 related patient features. Early diagnose of sepsis is critical to improve the outcome of ICU patients [42].

Not that for each time series in the above four datasets, every time point is tagged with a class label, which is the same as its outcome, such as earthquake, rain, mortality, sepsis.

5.2 Baselines

Based on Section 2, we use four types of baselines. They are all training strategies, so we use the same LSTM model for them.

The first is early classification of time series (ECTS-based).

- SR [19]. It has multiple models trained by the full-length time series. The classification result is the fusion result of its models.
- ECEC [20]. It has multiple models. Different models are trained by the data in different time stages. When classifying, it selects the classifier based on the time stage of the input data.

The second is replay continual learning (Replay-based).

- CLEAR [43]. It uses the reservoir sampling to limit the number of stored samples to a fixed budget assuming an i.i.d. data stream.
- CLOPS [14]. It trains a base model by replaying old tasks with importance-guided buffer storage to avoid catastrophic forgetting.

Table 1: Baselines Classification Accuracy (AUC-ROC \uparrow) for 4 Real-world Datasets at 10 Time Steps.
 *k% means the current classification is based on k% of the length of the time series; Bold font indicates the highest accuracy.

Dataset	Method	10%*	20%	30%	40%	50%	60%	70%	80%	90%	100%
UCR-EQ	SR	0.711±0.015	0.736±0.014	0.802±0.018	0.863±0.015	0.879±0.012	0.888±0.017	0.920±0.015	0.928±0.005	0.938±0.008	0.941±0.004
	ECEC	0.710±0.013	0.738±0.018	0.809±0.015	0.865±0.014	0.881±0.013	0.890±0.015	0.922±0.014	0.929±0.007	0.937±0.010	0.940±0.009
	CLEAR	0.725±0.011	0.768±0.018	0.821±0.017	0.874±0.016	0.882±0.009	0.895±0.014	0.919±0.008	0.923±0.002	0.930±0.010	0.933±0.003
	CLOPS	0.728±0.013	0.767±0.017	0.819±0.013	0.876±0.016	0.877±0.014	0.885±0.012	0.920±0.015	0.929±0.008	0.932±0.007	0.934±0.004
	DIF	0.724±0.018	0.770±0.015	0.825±0.015	0.880±0.013	0.886±0.013	0.904±0.012	0.913±0.014	0.923±0.004	0.929±0.008	0.932±0.005
	UNCERT	0.720±0.012	0.773±0.016	0.821±0.016	0.878±0.016	0.886±0.012	0.902±0.015	0.908±0.011	0.917±0.006	0.921±0.011	0.925±0.005
	DROPOUT	0.721±0.014	0.765±0.017	0.823±0.013	0.870±0.019	0.883±0.010	0.901±0.011	0.905±0.013	0.920±0.010	0.923±0.009	0.930±0.011
	TCP	0.729±0.011	0.771±0.012	0.827±0.012	0.881±0.010	0.881±0.011	0.898±0.015	0.907±0.010	0.921±0.009	0.925±0.010	0.933±0.008
	STL	0.731±0.017	0.770±0.013	0.830±0.017	0.875±0.012	0.879±0.008	0.906±0.005	0.910±0.009	0.918±0.006	0.928±0.007	0.934±0.005
	c³TS		0.736±0.013	0.774±0.014	0.840±0.014	0.882±0.014	0.899±0.010	0.907±0.009	0.925±0.010	0.933±0.008	0.939±0.006
USHCN	SR	0.719±0.020	0.730±0.022	0.750±0.020	0.761±0.023	0.802±0.020	0.836±0.016	0.889±0.012	0.902±0.013	0.923±0.010	0.933±0.009
	ECEC	0.721±0.019	0.736±0.024	0.752±0.021	0.760±0.025	0.800±0.019	0.837±0.016	0.890±0.011	0.906±0.017	0.921±0.011	0.931±0.009
	CLEAR	0.720±0.022	0.736±0.025	0.770±0.023	0.798±0.024	0.801±0.020	0.834±0.016	0.883±0.016	0.896±0.017	0.914±0.011	0.926±0.007
	CLOPS	0.722±0.025	0.728±0.026	0.758±0.024	0.781±0.023	0.805±0.018	0.838±0.013	0.885±0.013	0.899±0.010	0.923±0.009	0.928±0.005
	DIF	0.725±0.021	0.738±0.025	0.756±0.023	0.784±0.024	0.798±0.020	0.837±0.010	0.875±0.011	0.879±0.012	0.912±0.010	0.921±0.004
	UNCERT	0.724±0.023	0.740±0.024	0.751±0.019	0.781±0.025	0.800±0.023	0.835±0.016	0.873±0.012	0.877±0.011	0.907±0.007	0.919±0.007
	DROPOUT	0.719±0.017	0.738±0.016	0.750±0.015	0.782±0.014	0.797±0.018	0.829±0.013	0.877±0.013	0.885±0.012	0.908±0.013	0.920±0.015
	TCP	0.722±0.016	0.738±0.012	0.756±0.018	0.785±0.013	0.806±0.017	0.839±0.011	0.895±0.011	0.900±0.011	0.920±0.011	0.929±0.014
	STL	0.718±0.015	0.730±0.016	0.757±0.017	0.788±0.015	0.802±0.016	0.836±0.017	0.871±0.011	0.876±0.013	0.916±0.012	0.925±0.013
	c³TS		0.728±0.016	0.742±0.017	0.759±0.018	0.791±0.019	0.811±0.017	0.841±0.012	0.897±0.014	0.910±0.014	0.930±0.012
COVID-19	SR	0.730±0.024	0.831±0.022	0.867±0.016	0.900±0.018	0.913±0.013	0.923±0.010	0.937±0.007	0.946±0.006	0.960±0.004	0.962±0.005
	ECEC	0.732±0.028	0.834±0.020	0.870±0.016	0.904±0.014	0.913±0.012	0.924±0.015	0.940±0.012	0.948±0.007	0.958±0.008	0.963±0.007
	CLEAR	0.769±0.015	0.837±0.019	0.875±0.019	0.888±0.018	0.916±0.017	0.923±0.014	0.936±0.012	0.940±0.013	0.955±0.009	0.956±0.008
	CLOPS	0.779±0.017	0.835±0.018	0.869±0.015	0.885±0.021	0.914±0.017	0.924±0.017	0.936±0.011	0.939±0.010	0.951±0.007	0.953±0.005
	DIF	0.785±0.019	0.830±0.021	0.862±0.017	0.879±0.016	0.915±0.014	0.926±0.014	0.936±0.010	0.941±0.007	0.947±0.006	0.952±0.006
	UNCERT	0.775±0.013	0.841±0.016	0.871±0.015	0.900±0.013	0.915±0.013	0.925±0.015	0.935±0.009	0.940±0.007	0.950±0.005	0.954±0.006
	DROPOUT	0.740±0.017	0.831±0.020	0.860±0.013	0.885±0.012	0.912±0.011	0.924±0.018	0.932±0.011	0.939±0.007	0.943±0.004	0.945±0.005
	TCP	0.786±0.012	0.832±0.019	0.872±0.018	0.895±0.016	0.915±0.014	0.927±0.011	0.941±0.010	0.942±0.007	0.948±0.004	0.947±0.008
	STL	0.770±0.012	0.833±0.017	0.870±0.011	0.895±0.014	0.916±0.013	0.925±0.013	0.941±0.012	0.944±0.011	0.948±0.007	0.950±0.008
	c³TS		0.790±0.021	0.843±0.019	0.877±0.020	0.901±0.015	0.919±0.015	0.927±0.012	0.945±0.011	0.960±0.011	0.967±0.010
SEPSIS	SR	0.652±0.017	0.701±0.016	0.728±0.017	0.749±0.025	0.821±0.026	0.825±0.027	0.827±0.020	0.837±0.018	0.845±0.014	0.866±0.023
	ECEC	0.655±0.013	0.702±0.014	0.731±0.019	0.752±0.025	0.823±0.016	0.825±0.019	0.829±0.014	0.839±0.015	0.849±0.016	0.863±0.014
	CLEAR	0.667±0.022	0.711±0.022	0.733±0.023	0.763±0.020	0.827±0.026	0.830±0.026	0.838±0.024	0.847±0.017	0.848±0.015	0.854±0.016
	CLOPS	0.665±0.026	0.709±0.023	0.730±0.024	0.765±0.025	0.826±0.023	0.831±0.025	0.836±0.028	0.846±0.015	0.849±0.014	0.853±0.012
	DIF	0.666±0.025	0.710±0.021	0.732±0.024	0.755±0.027	0.825±0.025	0.832±0.027	0.839±0.028	0.844±0.012	0.847±0.010	0.848±0.016
	UNCERT	0.665±0.026	0.705±0.019	0.733±0.025	0.759±0.028	0.824±0.029	0.831±0.027	0.838±0.026	0.846±0.015	0.850±0.017	0.857±0.018
	DROPOUT	0.660±0.021	0.766±0.015	0.720±0.014	0.748±0.021	0.820±0.020	0.825±0.022	0.832±0.021	0.835±0.018	0.840±0.015	0.850±0.011
	TCP	0.662±0.021	0.705±0.016	0.722±0.020	0.758±0.025	0.826±0.027	0.827±0.027	0.839±0.025	0.840±0.017	0.843±0.011	0.862±0.016
	STL	0.660±0.023	0.709±0.016	0.727±0.021	0.760±0.024	0.824±0.026	0.833±0.028	0.836±0.015	0.840±0.016	0.845±0.016	0.855±0.017
	c³TS		0.671±0.025	0.715±0.024	0.734±0.021	0.768±0.026	0.831±0.023	0.840±0.024	0.840±0.018	0.851±0.012	0.853±0.012

The third is curriculum learning (CL-based).

- DIF [44]. It arranges curriculum (data learning order/task order) by loss/difficulty getting from the teacher model.
- UNCERT [14]. It arranges curriculum (data learning order/task order) by sentence uncertainty and model uncertainty.

The fourth is about the confidence research (Confidence-based).

- DROPOUT [29]. It uses the Monte Carlo Dropout method to approximate Bayesian inference and gets the model uncertainty.
- TCP [45]. It designs a true class probability to estimate the confidence of the model prediction.
- STL [46]. It gives the uncertainty about the signal temporal logic, which is more suitable for time series.

Table 2: The Performances of Solving Catastrophic Forgetting Problem.
The left part is BWT \uparrow results of baselines, the right part is FWT \uparrow results of baselines.

Method \ Dataset	SR	ECEC	CLEAR	CLOPS	DIF	STL	c³TS	SR	ECEC	CLEAR	CLOPS	DIF	STL	c³TS
UCR-EQ	+0.019	+0.021	+0.053	+0.052	+0.022	+0.020	+0.058	+0.121	+0.129	+0.312	+0.301	+0.124	+0.120	+0.345
USHCN	+0.028	+0.034	+0.063	+0.074	+0.017	+0.023	+0.084	+0.212	+0.128	+0.335	+0.301	+0.205	+0.216	+0.342
COVID-19	-0.001	+0.010	+0.009	+0.014	+0.002	+0.006	+0.020	+0.126	+0.221	+0.427	+0.439	+0.220	+0.232	+0.455
SEPSIS	-0.019	-0.017	+0.030	+0.032	-0.010	-0.008	+0.035	+0.095	+0.165	+0.401	+0.397	+0.106	+0.124	+0.410

Table 3: The Performances of Solving Over Fitting Problem.

Sepsis classification accuracy with non-uniform training sets and validation sets. \downarrow means the accuracy is greatly reduced (top 3).

Subset	SR	ECEC	CLEAR	CLOPS	DIF	UNCERT	DROPOUT	TCP	STL	c³TS
Male	0.844 \pm 0.017	0.851 \pm 0.019	0.853 \pm 0.017	0.855 \pm 0.016	0.852 \pm 0.015	0.858 \pm 0.024	0.853 \pm 0.022	0.850 \pm 0.021	0.856 \pm 0.022	0.860 \pm 0.022
Female	0.817 \pm 0.026 \downarrow	0.827 \pm 0.024	0.835 \pm 0.024	0.830 \pm 0.025 \downarrow	0.832 \pm 0.021	0.823 \pm 0.021 \downarrow	0.849 \pm 0.019	0.848 \pm 0.021	0.850 \pm 0.020	0.853 \pm 0.021
30-	0.845 \pm 0.021	0.853 \pm 0.016	0.863 \pm 0.012	0.867 \pm 0.016	0.861 \pm 0.015	0.854 \pm 0.024	0.848 \pm 0.021	0.850 \pm 0.022	0.852 \pm 0.024	0.864 \pm 0.022
30+	0.814 \pm 0.026	0.812 \pm 0.024 \downarrow	0.820 \pm 0.025 \downarrow	0.824 \pm 0.026	0.809 \pm 0.027 \downarrow	0.838 \pm 0.024	0.838 \pm 0.021	0.840 \pm 0.020	0.832 \pm 0.023	0.843 \pm 0.025

Table 4: Model Performances (Result Accuracy, Model Convergence, Result Uncertainty) under Different Learning Order.

All strategies train a same LSTM model to classify Sepsis, and their difference is reflected in the data learning order. Accuracy \uparrow is evaluated by using the full-length time series; Epochs \downarrow is the training rounds when the model converges; α \downarrow is the expected non-coverage probability for results.

	Random	Time \uparrow	Time \downarrow	Difficulty \uparrow	Difficulty \downarrow	Uncertainty \uparrow	Uncertainty \downarrow	Confidence \uparrow	Confidence\downarrow
Accuracy	0.831 \pm 0.018	0.854 \pm 0.013	0.861 \pm 0.012	0.862 \pm 0.015	0.852 \pm 0.014	0.866 \pm 0.017	0.853 \pm 0.022	0.853 \pm 0.017	0.872\pm0.016
Epochs		637	581	532	496	525	475	558	457
α		0.88	0.67	0.71	0.65	0.67	0.45	0.52	0.54

5.3 Evaluation Metrics

The classification accuracy is evaluated by the Area Under Curve of Receiver Operating Characteristic (AUC-ROC). The performance of solving forgetting&overfitting is evaluated by Backward Transfer (BWT) and Forward Transfer (FWT), the influence that learning a current has on the old/future. $R \in \mathbb{R}^{|\mathcal{M}| \times |\mathcal{M}|}$ is an accuracy matrix, $R_{i,j}$ is the accuracy on \mathcal{M}^j after learning \mathcal{M}^i . \bar{b} is the accuracy with random initialization.

$$BWT = \frac{1}{|\mathcal{M}| - 1} \sum_{i=1}^{|\mathcal{M}|-1} R_{|\mathcal{M}|,i} - R_{i,i} \quad (7)$$

$$FWT = \frac{1}{|\mathcal{M}| - 1} \sum_{i=2}^{|\mathcal{M}|} R_{i-1,i} - \bar{b}_{i,i} \quad (8)$$

The result uncertainty is evaluated by prediction intervals (PI). We use the conditional PI specifically [47]. CP is the conditional probability function that can construct an interval of $(1 - \alpha)$ confidence level. α is the expected non-coverage probability.

$$PI = [CP^{-1}(\frac{\alpha}{2}), CP^{-1}(1 - \frac{\alpha}{2})] \quad (9)$$

5.4 Results and Analysis

5.4.1 Analysis of multi-distribution

Before discussing method performances, we show the basic scenario of CCTS: multi-distribution. The data in different time stages have distinct statistical characteristics and finally form multiple distributions. Comparing Figure 1, Figure 5 and Figure 6, different data arrangement will lead to different multi-distribution. The fundamental goal of the following experiment is to model them.

Table 5: The percentage (%) of different samples between C^3TS and other methods in each data distribution of 4 training steps

Training / baby step	1	2	3	4	5
Time↓	0.14	0.24	0.27	0.20	0.10
Difficulty↑ (DIF)	0.17	0.28	0.22	0.18	0.09
Uncertainty↑ (UNCERT)	0.13	0.25	0.23	0.16	0.08

Table 6: C^3TS Performance with Different Distribution Number

N	2	4	6	8	10
Accuracy	0.839±0.01	0.861±0.01	0.816±0.01	0.797±0.01	0.784±0.01
BWT	+0.153	+0.156	+0.139	+0.135	+0.125
FWT	+0.398	+0.424	+0.379	+0.365	+0.364
Batches	685	657	632	578	594

Table 7: C^3TS Performance with Different Retrained Samples

M	1	2	3	4	5
Accuracy	0.861±0.01	0.841±0.01	0.826±0.01	0.813±0.01	0.815±0.01
BWT	+0.156	+0.144	+0.127	+0.130	+0.128
FWT	+0.424	0.443	+0.389	+0.376	+0.368

Figure 3: Catastrophic Forgetting and Over Fitting Case.

Accuracy decreases on Sepsis distribution 1 after learning Sepsis distribution 2. Much retraining on distribution 1 lets model over fit in it. While retraining partial samples of distribution 1 can alleviate these two problems.

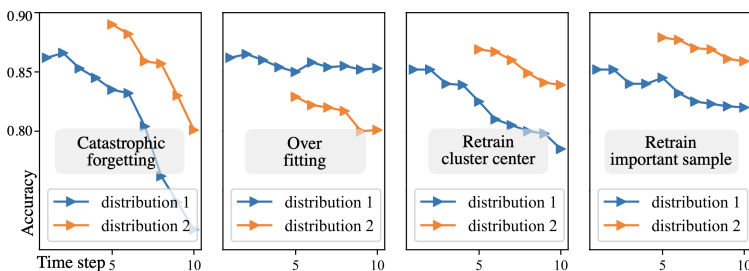
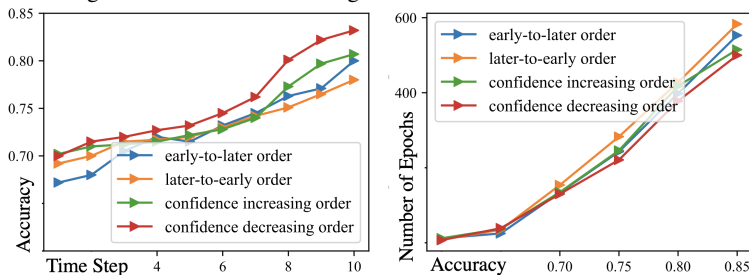


Figure 4: Model Performances under Different Training Orders.

For a LSTM model, its Sepsis classification accuracy and average training epochs are different by training orders of early-to-later, later-to-early, confidence-increasing and confidence-decreasing.



5.4.2 Continuous classification accuracy

Our method C^3TS is significantly better than all baselines. In Bonferroni-Dunn test, as shown in Table 1, $k = 10$, $n = 4$, $m = 5$ are the number of methods, datasets, cross-validation fold, then $N = n \times m = 20$, correspondingly, $q_{0.05} = 2.773$

Figure 5: Multi-distribution in Sepsis Dataset

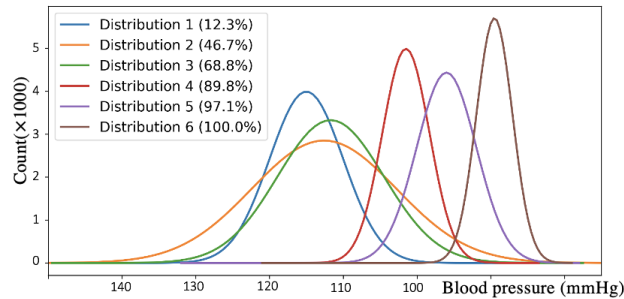
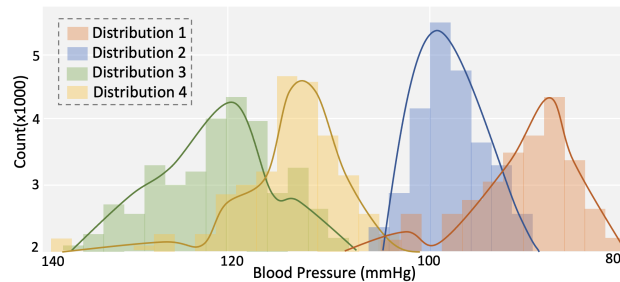


Figure 6: Four Distributions in Sepsis Dataset



and $CD = q_{0.05} \sqrt{\frac{k(k+1)}{6N}} = 2.655$, finally $rank(CCTS) = 1 < CD$. Thus, the accuracy is significantly improved. Specifically, CCTS can classify more accurately at every time than 9 baselines. Take sepsis diagnosis as an example, compared with the best baseline, our method improves the accuracy by 1.2% on average, 2.1% in the early 50% time stage when the key features are unobvious. Each hour of delayed treatment increases sepsis mortality by 4-8% [42]. With the same accuracy, we can predict 0.965 hour in advance.

5.4.3 Difficulty 1: Catastrophic forgetting and over fitting

Our method C^3TS can alleviate this problem. As shown in Table 1, C^3TS has the best performance on the early time series, showing the ability to alleviate catastrophic forgetting. As shown in Table 2, C^3TS has the highest BWT, meaning it has the lowest negative influence that learning the new tasks has on the old tasks. C^3TS has the highest FWT, meaning it has the highest positive influence that learning the former data distributions has on the task, especially for Sepsis and COVID-19 datasets. Figure 3 shows the case study that the partial replay of importance samples trades off forgetting and overfitting. Meanwhile, our method can avoid model overfitting and guarantee certain model generalizations. In Table 3, for most baselines, the accuracy on the validation set is much lower than that on the training set. But confidence-based methods, DROPOUT, TCP, STL and C^3TS , have relatively better generalization performance, which shows the potential of confidence-guided training strategies.

5.4.4 Difficulty 2: Optimal data learning order

The learning order based on confidence decline makes the LSTM model perform best in classification accuracy and model convergence, and perform well in result uncertainty, shown in Table 4 and Figure 4. Most existing methods ignore the data learning order and train the model randomly. Nowadays, some work [10] has paid attention to the data learning order, but the existing approaches basically use the difficulty and uncertainty to measure data. Experiments show that it has greater potential to define the order of model learning data by imitating human confidence in knowledge. Internal differences among these methods are analyzed in the next section.

5.4.5 The difference among baselines

Figure 7: Replayed Important Samples

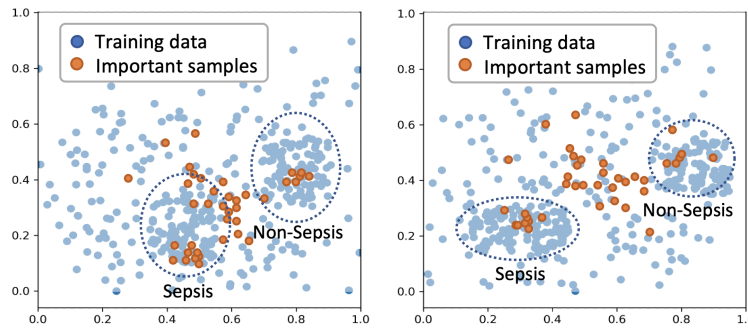
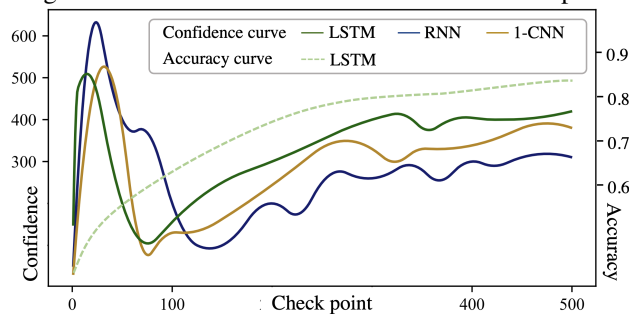


Figure 8: Curves of Confidence at Different Checkpoints



C^3TS and baselines will divide data into different distributions training/baby steps. Their divisions are considerably different. As shown in Table 5, the difference in the middle step, steps 2-4, is greater than that in steps 1 and 5. It shows that the most simple and complex time series quantified by different measures are relatively similar, and the main diversity lies in those samples of which the difficulties hardly to be distinguished. Therefore, we argue that the improvements of C^3TS may be mainly contributed by the differences in these middle steps.

5.4.6 Analysis of hyper-parameters

Different numbers of data distributions or baby steps will result in different model performances. We first construct U_{data} value of training data into a normal distribution $\mathcal{N}_u(\mu, \sigma^2)$, then use its confidence intervals $-n\sigma < \mu < -(n-1)\sigma$ to split the dataset and get initial N baby steps. As shown in Table 6, for sepsis classification, the best number of distribution or baby step is 4. We show these 4 distributions by visualizing blood pressure in Figure 6, which is obviously different from 6 distributions divided by time stages in Figure 5. Meanwhile, Different replayed samples will also result in different model performances. We construct the importance coefficient β into a positive skewed distribution $\mathcal{N}_\beta(m_o, \sigma^2)$ with $m_o < m_e < m$ and make $\epsilon = m$. The important sample is the data with $\beta > \epsilon$. To test the influence of the number of important samples, we replay $\frac{1}{M}$ of the data with larger β . For Sepsis dataset, the result is better with the decrease of M .

5.4.7 Analysis of Replayed Important Data

The important samples include not only the data hard to learn but also the representative data in each class, shown in Figure 7. It might be because the representative data is similar to the most common data, resulting in a greater additive loss, therefore leading to smaller coefficients in Equation 1. Thus, we can redefine the important samples: Important samples are samples that can represent most data of a class and samples that are difficult to distinguish by the model.

5.4.8 Analysis of the confidence change

C³TS can draw similar changing trends of confidence when training different models. As shown in Figure 8, the confidence first increases sharply, then drops and rises, eventually balances. It matches the overall trend of the Dunning-Kruger Effect curve and shows that we have realized the human confidence-imitated learning process. At the beginning of training, model weights are very random and the classification results with arbitrarily invalid weights are almost the same. As the confidence is based on the uncertainty, it will lead to high self-confidence although the model is ignorant. Then the model weights are gradually formed and mature, dropping weights randomly will easily lead to different classification results, so the confidence is reduced. Finally, The weight is robust enough to avoid the result change caused by weight failure. As shown in Figure 8, low accuracy corresponds to blindly high confidence, fairly good accuracy corresponds to confused low confidence, and high accuracy corresponds to robust high confidence.

6 Conclusion

In this paper, we propose a new concept of Continuous Classification of Time Series (CCTS) to meet real needs. To solve its two problems about data arrangement and model training, we design a novel confidence-guided approach C³TS. It can imitate the human behavior of object-confidence and self-confidence by the importance coefficient method and uncertainty evaluation. It arranges data iteration and reviews data according to importance-based object-confidence, and schedules training duration and training order according to uncertainty-based self-confidence. We test the method on four real-world datasets based on perspectives of classification accuracy, solving two difficulties, the difference among baselines, hyper-parameters, and analysis of data distributions and confidence change. The results show that our method is better than all baselines. It proves that the confidence-guided training strategy is an effective and self-adaptive indicator to guide the training and is more worthy of future research.

Acknowledgments

This work was supported by the National Key Research and Development Program of China under Grant 2021YFE0205300, and the National Natural Science Foundation of China (No.62172018, No.62102008).

References

- [1] W. Chen, J. Wang, Q. L. Fe Ng, S. C. Xu, and L. Ba. The treatment of severe and multiple injuries in intensive care unit: report of 80 cases. *European Review for Medical & Pharmacological Sciences*, 18(24):3797, 2014.
- [2] German Ignacio Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, 2019.
- [3] Yinpei Dai and Hangyu Li. Preview, attend and review: Schema-aware curriculum learning for multi-domain dialogue state tracking. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 879–885, 2021.
- [4] Yunchao Wei, Xiaodan Liang, Yunpeng Chen, Xiaohui Shen, Ming-Ming Cheng, Jiashi Feng, Yao Zhao, and Shuicheng Yan. STC: A simple to complex framework for weakly-supervised semantic segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(11):2314–2320, 2017.
- [5] Benfeng Xu, Licheng Zhang, Zhendong Mao, Quan Wang, and Hongtao Xie. Curriculum learning for natural language understanding. In *Proceedings of Association for Computational Linguistics (ACL)*, pages 6095–6104, 2020.
- [6] Ashish Gupta, Hari Prabhat Gupta, Bhaskar Biswas, and Tanima Dutta. Approaches and applications of early classification of time series: A review. *IEEE Trans. Artif. Intell.*, 1(1):47–61, 2020.
- [7] Matthias Delange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Ales Leonardis, Greg Slabaugh, and Tinne Tuytelaars. A continual learning survey: Defying forgetting in classification tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021.
- [8] Dongsub Shim, Zheda Mai, Jihwan Jeong, and Scott Sanner. Online class-incremental continual learning with adversarial shapley value. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 9630–9638, 2021.
- [9] Gobinda Saha and Isha Garg. Gradient projection memory for continual learning. In *International Conference on Learning Representations (ICLR)*, 2021.
- [10] Xin Wang, Yudong Chen, and Wenwu Zhu. A comprehensive survey on curriculum learning. *CoRR*, abs/2010.13166, 2020.

- [11] J. Kruger and D. Dunning. Unskilled and unaware of it: How difficulties in recognizing one’s own incompetence lead to inflated self-assessments. *Journal of Personality and Social Psychology*, 1(6):30–46, 2000.
- [12] Yikai Zhou, Baosong Yang, Derek F. Wong, and Yu Wan. Uncertainty-aware curriculum learning for neural machine translation. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 6934–6944, 2020.
- [13] L. Teneyck. *Dunning-Kruger Effect*. Decision Making in Emergency Medicine, 2021.
- [14] D Kiyasseh, T. Zhu, and D Clifton. A clinical deep learning framework for continually learning from cardiac signals across diseases, time, modalities, and institutions. *Nature Communications*, 12(1):4221, 2021.
- [15] Tiago Santos and Roman Kern. A literature survey of early time series classification and deep learning. In *Proceedings Methods in Industry 4.0*, 2016.
- [16] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. *Data Min. Knowl. Discov.*, 33(4):917–963, 2019.
- [17] Bin Liu, Ying Li, Zhaonan Sun, Soumya Ghosh, and Kenney Ng. Early prediction of diabetes complications from electronic health records: A multi-task survival analysis approach. In *Proceedings of AAAI Conference on Artificial Intelligence (AAAI)*, pages 101–108, 2018.
- [18] M. A. Reyna, C. S. Josef, R. Jeter, S. P. Shashikumar, and A. Sharma. Early prediction of sepsis from clinical data: The physionet/computing in cardiology challenge 2019. *Critical Care Medicine*, 48(2):1, 2019.
- [19] Usue Mori, Alexander Mendiburu, Sanjoy Dasgupta, and José Antonio Lozano. Early classification of time series by simultaneously optimizing the accuracy and earliness. *IEEE Trans. Neural Networks Learn. Syst.*, 29(10):4569–4578, 2018.
- [20] Junwei Lv, Xuegang Hu, Lei Li, and Peipei Li. An effective confidence-based early classification of time series. *IEEE Access*, 7:96113–96124, 2019.
- [21] David Isele and Akansel Cosgun. Selective experience replay for lifelong learning. In *AAAI*, pages 3302–3309, 2018.
- [22] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. Curriculum learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 382, pages 41–48, 2009.
- [23] Zhiyu Liang and Hongzhi Wang. Efficient class-specific shapelets learning for interpretable time series classification. *Inf. Sci.*, 570:428–450, 2021.
- [24] Siwei Lai, Liheng Xu, Kang Liu, and Jun Zhao. Recurrent convolutional neural networks for text classification. In *Proceedings of AAAI Conference on Artificial Intelligence (AAAI)*, pages 2267–2273, 2015.
- [25] Joongho Jo and Jongsun Park. Confidence score based mini-batch skipping for CNN training on mini-batch training environment. In *International SoC Design Conference (ISOCC)*, pages 129–130, 2020.
- [26] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 5574–5584, 2017.
- [27] Li Dong, Chris Quirk, and Mirella Lapata. Confidence modeling for neural semantic parsing. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 743–753, 2018.
- [28] Yijun Xiao and William Yang Wang. Quantifying uncertainties in natural language processing tasks. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 7322–7329, 2019.
- [29] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *Proceedings of International Conference on Machine Learning (ICML)*, volume 48, pages 1050–1059, 2016.
- [30] Andrey Malinin and Mark J. F. Gales. Predictive uncertainty estimation via prior networks. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 7047–7058, 2018.
- [31] Xuan Zhang et al. An empirical exploration of curriculum learning for neural machine translation. *CoRR*, abs/1811.00739, 2018.
- [32] Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabás Póczos, and Tom M. Mitchell. Competence-based curriculum learning for neural machine translation. In *Proceedings of Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT*, pages 1162–1172, 2019.
- [33] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 6402–6413, 2017.

- [34] Rich Caruana and Steve Lawrence. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 402–408, 2000.
- [35] Yanping Chen, Eamonn Keogh, Bing Hu, Nurjahan Begum, Anthony Bagnall, Abdullah Mueen, and Gustavo Batista. The ucr time series classification archive, July 2015. www.cs.ucr.edu/~eamonn/time_series_data/.
- [36] C. J. Ammon, A. A. Velasco, T. Lay, and T. C. Wallace. *Earthquake prediction, forecasting, & early warning*. Academic press, 2021.
- [37] M. et al. Menne. Long-term daily and monthly climate records from stations across the contiguous united states. *U.S. Historical Climatology Network*, 2016.
- [38] W. Y. Lee, S. K. Park, and H. H. Sung. The optimal rainfall thresholds and probabilistic rainfall conditions for a landslide early warning system for chuncheon, republic of korea. *Landslides*, 2021.
- [39] Goncalves J et al. Yan L, Zhang H T. An interpretable mortality prediction model for covid-19 patients. *Nature, Machine intelligence*, 2, 2020.
- [40] Chenxi Sun, Shenda Hong, Moxian Song, Hongyan Li, and Zhenjie Wang. Predicting covid-19 disease progression and patient outcomes based on temporal deep learning. *BMC Medical Informatics and Decision Making*, 21:45, 2020.
- [41] Matthew A. Reyna, Christopher Josef, Salman Seyedi, Russell Jeter, Supreeth P. Shashikumar, M. Brandon Westover, Ashish Sharma, Shamim Nemati, and Gari D. Clifford. Early prediction of sepsis from clinical data: the physionet/computing in cardiology challenge 2019. In *46th Computing in Cardiology, CinC 2019, Singapore, September 8-11, 2019*, pages 1–4. IEEE, 2019.
- [42] Christopher W Seymour, Foster Gesten, Hallie C Prescott, and Marcus E Friedrich. Time to treatment and mortality during mandated emergency care for sepsis. *New England Journal of Medicine*, 376(23):2235–2244, 2017.
- [43] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy P. Lillicrap, and Gregory Wayne. Experience replay for continual learning. In *Advances in Neural Information Processing (NeurIPS)*, pages 348–358, 2019.
- [44] Guy Hacohen and Daphna Weinshall. On the power of curriculum learning in training deep networks. In *Proceedings of International Conference on Machine Learning (ICML)*, volume 97, pages 2535–2544, 2019.
- [45] Charles Corbière, Nicolas Thome, Avner Bar-Hen, Matthieu Cord, and Patrick Pérez. Addressing failure prediction by learning model confidence. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 2898–2909, 2019.
- [46] Nasim Baharisangari, Jean-Raphaël Gaglione, Daniel Neider, Ufuk Topcu, and Zhe Xu. Uncertainty-aware signal temporal logic. *CoRR*, abs/2105.11545, 2021.
- [47] Hussain Mohammed Dipu Kabir, Abbas Khosravi, Abdollah Kavousi-Fard, Saeid Nahavandi, and Dipti Srinivasan. Optimal uncertainty-guided neural network training. *Appl. Soft Comput.*, 99:106878, 2021.
- [48] Alistair EW et al. Johnson. MIMIC-III, a freely accessible critical care database. *Scientific data*, 3:160035, 2016.
- [49] Edward Choi, Andy Schuetz, Walter F Stewart, and Jimeng Sun. Using recurrent neural network models for early detection of heart failure onset. *Journal of the American Medical Informatics Association*, 24(2):361–370, 2017.
- [50] James Kirkpatrick and Razvan Pascanu. Overcoming catastrophic forgetting in neural networks. *CoRR*, abs/1612.00796, 2016.
- [51] David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 6467–6476, 2017.
- [52] Lin Chen, Christopher Harshaw, Hamed Hassani, and Amin Karbasi. Projection-free online optimization with stochastic gradient: From convexity to submodularity. In *Proceedings of the International Conference on Machine Learning (ICML)*, volume 80, pages 813–822, 2018.
- [53] Jiahao Xie, Zebang Shen, Chao Zhang, and Boyu Wang. Efficient projection-free online methods with stochastic recursive gradient. In *The AAAI Conference on Artificial Intelligence (AAAI)*, pages 6446–6453. AAAI Press, 2020.

Appendix for C³TS: Confidence-guided Learning Process for Continuous Classification of Time Series

A Concept

A.1 Continuous Classification of Time Series (CCTS)

Continuous Classification of Time Series (CCTS) aims to classify as accurately as possible with the evolution of time series at every time.

Time series is one of the most common data forms, the popularity of time series classification has attracted increasing attention in many practical fields, such as healthcare and industry. In the real world, many applications require classification at every time. For example, in the Intensive Care Unit (ICU), critical patients' vital signs develop dynamically, the status perception and disease diagnosis are needed at any time. Timely diagnosis provides more opportunities to rescue lives. In response to the current demand, we propose a new task – Continuous Classification of Time Series (CCTS). It aims to classify as accurately as possible at every time in time series.

A.2 Concept Difference

A.2.1 One-shot Classification, Continuous Classification

The existing (early) classification of time series is the one-shot classification, where the classification is performed only once at the final or an early time. However, many real-world applications require continuous classification. For example, intensive care patients should be detected and diagnosed at all times to facilitate timely life-saving.

- **One-shot Classification.** A time series $X = \{x_1, \dots, x_T\}$ is labeled with classes C . OC classifies X at a time $t, t \leq T$ with the minimum loss $\mathcal{L}(f(X_{1:t}), C)$.
- **Continuous Classification.** A time series is $X = \{x_1, \dots, x_T\}$. At time $t, x_{1:t}$ is labeled with class c_t . Continuous Classification classifies $x_{1:t}$ at every time $t = 1, \dots, T$ with the minimum loss $\sum_{t=1}^T \mathcal{L}(f(x_{1:t}), c^t)$.

A.2.2 Classification of Time Series (CTS), Early Classification of Time Series (ECTS), Continuous Classification of Time Series (CCTS)

The popularity of time series classification has attracted increasing attention in many practical fields. The foundation is Classification of Time Series (CTS). It makes classification based on the full-length data. But in time-sensitive applications, Early Classification of Time Series (ECTS) is more critical, classifying at an early time. For example, early diagnosis helps for sepsis outcomes. But both of them classify only once and lean a single data distribution. In fact, CCTS is composed of multiple ECTS and the continuous classification is composed of multiple one-shot classification.

- **Classification of Time Series (CTS).** A dataset of time series $\mathcal{D} = \{(X^n, C^n)\}_{n=1}^N$ has N samples. Each time series X^n is labeled with a class C^n , CTS classifies time series using the full-length data by model $f : f(X) \rightarrow C$.
- **Early Classification of Time Series (ECTS).** A dataset of time series $\mathcal{D} = \{(X^n, C^n)\}_{n=1}^N$ has N samples. Each time series $X^n = X_{t=1}^T$ is labeled with a class C^n . ECTS classifies time series in an advanced time t by model $f : f(\{x_1, x_2, \dots, x_t\}) \rightarrow C$, where $t < T$.
- **Continuous Classification of Time Series (CCTS).** A dataset of time series $\mathcal{D} = \{(X^n, C^n)\}_{n=1}^N$ has N samples. Each time series $X^n = X_{t=1}^T$ is labeled with a class C^n . CCTS classifies time series in every time t by model $f : f(\{x_1, x_2, \dots, x_t\}) \rightarrow C$, where $t = 1, \dots, T$.

A.2.3 Curriculum Learning (CL), Continual Learning (CL), Online Learning (OL), Continuous Classification of Time Series (CCTS)

- Curriculum Learning (CL). A curriculum is a sequence of training criteria over T training steps $\mathcal{C} = \{Q_1, \dots, Q_T\}$. Each criterion Q_t is a reweighting of the target training distribution $P(z)$, $Q_t(z) \propto W_t(z)P(z)$, \forall example $x \in$ training set D , such that the following three conditions are satisfied: (1) The entropy of distributions gradually increases, i.e., $H(Q_t) < H(Q_{t+1})$. (2) The weight for any example increases, i.e., $W_t(z) \leq W_{t+1}(z)$, $\forall z \in D$. (3) $Q_T(z) = P(z)$.
- Continual Learning (CL). A CL issue $\mathcal{T} = \{T^1, T^2, \dots, T^N\}$ has a sequence of N tasks. Each task $T^n = (X^n, C^n)$ is represented by the training sample X^n with classes C^n . CL learns a new task at every moment. The goal is to control the statistical risk of all seen tasks $\sum_{n=1}^N \mathbb{E}_{(X^n, C^n)}[\mathcal{L}(f_n((X^n; \theta), C^n)]$ with loss \mathcal{L} , network function f_n and parameters θ .
- Online Learning (OL). A OL issue has a sequence of dataset $\mathcal{X} = \{X^1, X^2, \dots, X^N\}$ for one task \mathcal{T} . Each dataset X^t has a distribution D^t . CL learns a new D^t at every time t . The goal is to find the optimal solution of \mathcal{T} after N iterations by minimize the regret $\mathcal{R} := \sum_{t=1}^N (f^t(X^t) - \min f^t(X^t))$.
- Continuous Classification of Time Series (CCTS). A dataset \mathcal{S} contains N time series. Each time series $X = \{x_1, \dots, x_T\}$ is labeled with a class $C \in \mathcal{C}$ at the final time T . As time series varies among time, it has a subsequence series with M different distributions $\mathcal{D} = \{\mathcal{D}^1, \dots, \mathcal{D}^M\}$, each \mathcal{D}^m has subsequence $X_{1:t^m}$. CCTS learns every \mathcal{D}^m and introduces a task sequence $\mathcal{M} = \{\mathcal{M}^1, \dots, \mathcal{M}^M\}$ to minimize the additive risk $\sum_{m=1}^M \mathbb{E}_{\mathcal{M}^m}[\mathcal{L}(f^m(\mathcal{D}^m; \theta), C)]$ with model f and parameter θ . f^m is the model f after being trained for \mathcal{M}^m . When the model is trained for \mathcal{M}^m , its performance on all observed data cannot degrade: $\frac{1}{m} \sum_{i=1}^m \mathcal{L}(f^i, \mathcal{M}^i) \leq \frac{1}{m-1} \sum_{i=1}^{m-1} \mathcal{L}(f^i, \mathcal{M}^i)$.

A.2.4 Data Stream Classification, Multi-step Prediction, Anytime Classification

- Data stream (big data, focus on the current) and time series (a whole data, focus on the overall) are two different data form. Data stream classification aims to the rapid feature extraction and model optimization for the arriving data, while our CCTS can classify new time series anytime without training again.
- Multi-step prediction predicts values in multiple future time steps, similar to the early classification of time series (ECTS). But our CCTS aims to predict the final label, like patient outcome, at every time. CCTS has to solve the problem of catastrophic forgetting, but ECTS doesn't.
- 'Anytime' is another expression of 'continuous'.

B Experiments

B.1 Datasets

- One data corresponds to one label, e.g. the label 'sepsis' for a Sepsis data and 'mortality' for a COVID-19 data. We copy the final label of the time series to its each time point.
 - Sepsis dataset [41] has 30,336 records with 2,359 diagnosed sepsis. Early diagnose is critical to improve sepsis outcome [42]. In this dataset, the time series are the changes of 40 related patient features, the label at each time is sepsis or non-sepsis.
 - COVID-19 dataset [39] has 6,877 blood samples of 485 COVID-19 patients from Tongji Hospital, Wuhan, China. Mortality prediction helps for treatment and rational resource allocation [40]. In this dataset, the time series are the changes of blood samples, the label at each time is mortality or survival.
 - UCR-EQ dataset has 471 earthquake records from UCR time series database archive. It is the univariate time series of seismic feature value. Natural disaster early warning, like earthquake warning, helps to reduce casualties and property losses.
 - USHCNrain dataset [37] has the daily meteorological data of 48 states in U.S. from 1887 to 2014. It is the multivariate time series of 5 weather features. Rainfall warning is not only the demand of daily life, but also can help prevent natural disasters.

Not that for each time series in the above four datasets, every time point is tagged with a class label, which is the same as its outcome label, such as 'mortality', 'sepsis', 'earthquake' and 'rain'.

- Different time points of a time series have different labels.

- MIMIC-III dataset [48] has 19,993 admission records of 7,537 patients. We focus on 10 diagnoses (ICD-9): HIV (042), Brain Cancer (191), Diabetes(249), Hypertension (401), Heart Failure (428), Pneumonia (480-486), Gastric Ulcer (531), Hepatopathy (571), Nephropathy (580-589), SIRS (995.9). The time series are vital signs, and labels at each time are some diagnoses.
- USHCN dataset [37] has U.S. daily meteorological data from 1887 to 2014. We focus on 4 weather conditions in New York: sunny, overcast, rainfall, snowfall. The time series are records of 4 neighboring states, labels at each time are weather after a week.

B.2 Baselines

Baselines

- Early classification based methods:
 - LSTM [49] uses a base model trained by time series at every time stage.
 - SR [19] applies multiple base models trained by the full-length time series and uses the fusion result.
 - ECEC [20] has a set of base classifiers trained by time series in different time stages.
- Curriculum learning based methods:
 - DIF [44]. It arranges curriculum (data learning order/task order) by loss/difficulty getting from the teacher model.
 - UNCERT [14]. It arranges curriculum (data learning order/task order) by sentence uncertainty and model uncertainty.
- Continual learning based methods:
 - EWC [50]. It is a regularization-based method, training a model to remember the old tasks by constraining important parameters to stay close to their old values.
 - GEM [51] trains a base model to remember the old tasks by finding the new gradients which are at acute angles to the old gradients.
 - CLEAR. It is a replay-based method, using the reservoir sampling to limit the number of stored samples to a fixed budget assuming an i.i.d. data stream.
 - CLOPS [14] trains a base model by replaying old tasks with importance-guided buffer storage and uncertainty-based buffer acquisition.
- Online learning based methods:
 - OSFW [52] trains a base model by stochastic gradient estimator to achieves the stochastic regret bound.
 - ORGFW [53] trains the model by a recursive gradient estimator to achieve an optimal regret bound.
- Confidence based methods:
 - DROPOUT [29]. It uses the Monte Carlo Dropout method to approximate Bayesian inference and gets the model uncertainty.
 - TCP [45]. It designs a true class probability to estimate the confidence of the model prediction.
 - STL [46]. It gives the uncertainty about the signal temporal logic, which is more suitable for time series.

Table 8: Classification Accuracy (AUC-ROC \uparrow) of Baselines for 4 Real-world Datasets at 10 Time Steps.

¹The dataset has one task. Use online learning methods, OSFW and ORGFW, as baselines. ²The dataset has multiple tasks. Use continual learning methods, GEM and CLOPS, as baselines. ³All methods use LSTM as the base model for fairness. ⁴k% means the current classification time is k% of the total time of the full-length time series. ⁵The value is the average accuracy from 5-fold cross-validation with mean \pm std.

Dataset	Method ³	Time ⁴	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
Sepsis ¹	LSTM		0.576 \pm 0.060 ⁵	0.629 \pm 0.031	0.715 \pm 0.060	0.736 \pm 0.061	0.745 \pm 0.053	0.748 \pm 0.044	0.773 \pm 0.035	0.795 \pm 0.025	0.813 \pm 0.022	0.827 \pm 0.013
	SR		0.626 \pm 0.036	0.659 \pm 0.031	0.738 \pm 0.0100	0.761 \pm 0.020	0.803 \pm 0.0010	0.807 \pm 0.030	0.815 \pm 0.010	0.835 \pm 0.010	0.835 \pm 0.020	0.850 \pm 0.020
	ECEC		0.623 \pm 0.020	0.669 \pm 0.010	0.731 \pm 0.010	0.763 \pm 0.010	0.811 \pm 0.010	0.815 \pm 0.010	0.820 \pm 0.010	0.823 \pm 0.010	0.839 \pm 0.010	0.851 \pm 0.010
	OSFW		0.670 \pm 0.010	0.712 \pm 0.020	0.735 \pm 0.010	0.765 \pm 0.020	0.821 \pm 0.030	0.825 \pm 0.020	0.839 \pm 0.010	0.850 \pm 0.010	0.851 \pm 0.010	0.861 \pm 0.010
	ORGFW		0.670 \pm 0.020	0.714 \pm 0.020	0.732 \pm 0.020	0.766 \pm 0.030	0.820 \pm 0.020	0.831 \pm 0.020	0.832 \pm 0.030	0.843 \pm 0.010	0.850 \pm 0.010	0.862 \pm 0.010
	DIF		0.666 \pm 0.025	0.710 \pm 0.021	0.732 \pm 0.024	0.755 \pm 0.027	0.825 \pm 0.025	0.832 \pm 0.027	0.839 \pm 0.028	0.844 \pm 0.012	0.847 \pm 0.010	0.848 \pm 0.016
	UNCERT		0.665 \pm 0.026	0.705 \pm 0.019	0.733 \pm 0.025	0.759 \pm 0.028	0.824 \pm 0.029	0.831 \pm 0.027	0.838 \pm 0.026	0.846 \pm 0.015	0.850 \pm 0.017	0.857 \pm 0.018
	DROPOUT		0.660 \pm 0.021	0.766 \pm 0.015	0.720 \pm 0.014	0.748 \pm 0.021	0.820 \pm 0.020	0.825 \pm 0.022	0.832 \pm 0.021	0.835 \pm 0.018	0.840 \pm 0.015	0.850 \pm 0.011
	TCP		0.662 \pm 0.021	0.705 \pm 0.016	0.722 \pm 0.020	0.758 \pm 0.025	0.826 \pm 0.027	0.827 \pm 0.027	0.839 \pm 0.025	0.840 \pm 0.017	0.843 \pm 0.011	0.862 \pm 0.016
	STL		0.660 \pm 0.023	0.709 \pm 0.016	0.727 \pm 0.021	0.760 \pm 0.024	0.824 \pm 0.026	0.833 \pm 0.028	0.836 \pm 0.015	0.840 \pm 0.016	0.845 \pm 0.016	0.855 \pm 0.017
	C³TS			0.671\pm0.025	0.715\pm0.024	0.734\pm0.021	0.768\pm0.026	0.831\pm0.023	0.840\pm0.024	0.840\pm0.018	0.851\pm0.012	0.853\pm0.012
COVID-19 ²	LSTM		0.605 \pm 0.04	0.701 \pm 0.03	0.793 \pm 0.02	0.833 \pm 0.01	0.844 \pm 0.01	0.888 \pm 0.01	0.918 \pm 0.03	0.925 \pm 0.01	0.939 \pm 0.00	0.944 \pm 0.01
	SR		0.636 \pm 0.01	0.730 \pm 0.02	0.810 \pm 0.01	0.867 \pm 0.01	0.901 \pm 0.01	0.900 \pm 0.01	0.935 \pm 0.01	0.946 \pm 0.00	0.952 \pm 0.01	0.962 \pm 0.00
	ECEC		0.639 \pm 0.01	0.732 \pm 0.02	0.829 \pm 0.01	0.870 \pm 0.01	0.901 \pm 0.02	0.904 \pm 0.01	0.937 \pm 0.00	0.948 \pm 0.01	0.952 \pm 0.00	0.955 \pm 0.01
	OSFW		0.693 \pm 0.02	0.765 \pm 0.01	0.865 \pm 0.01	0.887 \pm 0.02	0.915 \pm 0.01	0.913 \pm 0.01	0.941 \pm 0.00	0.955 \pm 0.01	0.960 \pm 0.01	0.964 \pm 0.00
	ORGFW		0.709 \pm 0.01	0.775 \pm 0.01	0.849 \pm 0.01	0.878 \pm 0.01	0.916 \pm 0.02	0.912 \pm 0.01	0.945 \pm 0.01	0.957 \pm 0.00	0.961 \pm 0.00	0.965 \pm 0.00
	DIF		0.785 \pm 0.019	0.830 \pm 0.021	0.862 \pm 0.017	0.879 \pm 0.016	0.915 \pm 0.014	0.926 \pm 0.014	0.936 \pm 0.010	0.941 \pm 0.007	0.947 \pm 0.006	0.952 \pm 0.006
	UNCERT		0.775 \pm 0.013	0.841 \pm 0.016	0.871 \pm 0.015	0.900 \pm 0.013	0.915 \pm 0.013	0.925 \pm 0.015	0.935 \pm 0.009	0.940 \pm 0.007	0.950 \pm 0.005	0.954 \pm 0.006
	DROPOUT		0.740 \pm 0.017	0.831 \pm 0.020	0.860 \pm 0.013	0.885 \pm 0.012	0.912 \pm 0.011	0.924 \pm 0.018	0.932 \pm 0.011	0.939 \pm 0.007	0.943 \pm 0.004	0.945 \pm 0.005
	TCP		0.786 \pm 0.012	0.832 \pm 0.019	0.872 \pm 0.018	0.895 \pm 0.016	0.915 \pm 0.014	0.927 \pm 0.011	0.941 \pm 0.010	0.942 \pm 0.007	0.948 \pm 0.004	0.947 \pm 0.008
	STL		0.770 \pm 0.012	0.833 \pm 0.017	0.870 \pm 0.011	0.895 \pm 0.014	0.916 \pm 0.013	0.925 \pm 0.013	0.941 \pm 0.012	0.944 \pm 0.011	0.948 \pm 0.007	0.950 \pm 0.008
	C³TS			0.790\pm0.021	0.843\pm0.019	0.877\pm0.020	0.901\pm0.015	0.919\pm0.015	0.927\pm0.012	0.945\pm0.011	0.960\pm0.011	0.967\pm0.010
MIMIC-III ²	LSTM		0.574 \pm 0.04	0.595 \pm 0.04	0.635 \pm 0.01	0.661 \pm 0.02	0.706 \pm 0.02	0.724 \pm 0.02	0.770 \pm 0.01	0.785 \pm 0.01	0.793 \pm 0.01	0.817 \pm 0.01
	SR		0.628 \pm 0.02	0.681 \pm 0.01	0.695 \pm 0.01	0.692 \pm 0.01	0.732 \pm 0.01	0.749 \pm 0.01	0.785 \pm 0.01	0.805 \pm 0.01	0.820 \pm 0.01	0.825 \pm 0.01
	ECEC		0.651 \pm 0.02	0.702 \pm 0.01	0.699 \pm 0.01	0.703 \pm 0.01	0.733 \pm 0.01	0.756 \pm 0.01	0.792 \pm 0.01	0.808 \pm 0.01	0.825 \pm 0.01	0.838 \pm 0.01
	GEM		0.685 \pm 0.01	0.714 \pm 0.02	0.732 \pm 0.01	0.747 \pm 0.01	0.759 \pm 0.01	0.772 \pm 0.01	0.805 \pm 0.01	0.816 \pm 0.01	0.826 \pm 0.01	0.827 \pm 0.01
	CLOPS		0.691 \pm 0.01	0.722 \pm 0.01	0.743 \pm 0.01	0.752 \pm 0.01	0.764 \pm 0.01	0.770 \pm 0.01	0.808 \pm 0.02	0.815 \pm 0.01	0.825 \pm 0.01	0.831 \pm 0.01
	C³TS			0.714\pm0.02	0.731\pm0.02	0.757\pm0.01	0.763\pm0.01	0.780\pm0.01	0.788\pm0.01	0.810\pm0.02	0.820\pm0.01	0.832\pm0.01
USHCN ²	LSTM		0.618 \pm 0.03	0.621 \pm 0.03	0.674 \pm 0.03	0.694 \pm 0.03	0.724 \pm 0.04	0.750 \pm 0.02	0.775 \pm 0.02	0.793 \pm 0.03	0.856 \pm 0.02	0.861 \pm 0.01
	SR		0.674 \pm 0.03	0.679 \pm 0.01	0.706 \pm 0.02	0.722 \pm 0.02	0.752 \pm 0.01	0.767 \pm 0.01	0.805 \pm 0.01	0.831 \pm 0.01	0.866 \pm 0.01	0.885 \pm 0.01
	ECEC		0.688 \pm 0.03	0.701 \pm 0.02	0.721 \pm 0.02	0.730 \pm 0.01	0.769 \pm 0.01	0.781 \pm 0.01	0.806 \pm 0.02	0.839 \pm 0.01	0.870 \pm 0.01	0.889 \pm 0.01
	GEM		0.709 \pm 0.02	0.741 \pm 0.01	0.763 \pm 0.01	0.762 \pm 0.02	0.780 \pm 0.01	0.809 \pm 0.02	0.825 \pm 0.01	0.843 \pm 0.01	0.877 \pm 0.01	0.889 \pm 0.01
	CLOPS		0.721 \pm 0.02	0.740 \pm 0.02	0.763 \pm 0.01	0.771 \pm 0.02	0.772 \pm 0.02	0.805 \pm 0.01	0.830 \pm 0.02	0.847 \pm 0.01	0.875 \pm 0.01	0.897 \pm 0.01
	C³TS			0.732\pm0.01	0.750\pm0.01	0.767\pm0.01	0.780\pm0.02	0.783\pm0.01	0.811\pm0.01	0.848\pm0.02	0.860\pm0.01	0.884\pm0.01

Table 9: COVID-19 Classification Accuracy with Non-uniform Training Sets and Validation Sets. \downarrow means the accuracy is greatly reduced

Subset	SR	ECEC	CLEAR	CLOPS	DIF	UNCERT	C ³ TS
Male	0.968 \pm 0.014	0.969 \pm 0.016	0.965 \pm 0.012	0.965 \pm 0.004	0.978 \pm 0.009	0.978 \pm 0.014	0.971 \pm 0.010
Female	0.915 \pm 0.004 \downarrow	0.917 \pm 0.015 \downarrow	0.919 \pm 0.018	0.928 \pm 0.003	0.919 \pm 0.008 \downarrow	0.941 \pm 0.009	0.947 \pm 0.002
Age 30-	0.965 \pm 0.014	0.967 \pm 0.015	0.967 \pm 0.013	0.964 \pm 0.009	0.977 \pm 0.008	0.979 \pm 0.012	0.972 \pm 0.010
Age 30+	0.941 \pm 0.007	0.943 \pm 0.018	0.931 \pm 0.008 \downarrow	0.923 \pm 0.040 \downarrow	0.902 \pm 0.006 \downarrow	0.914 \pm 0.007 \downarrow	0.945 \pm 0.006
Test	0.964 \pm 0.013	0.968 \pm 0.015	0.966 \pm 0.012	0.962 \pm 0.006	0.979 \pm 0.009	0.978 \pm 0.010	0.970 \pm 0.007
Valid.	0.962 \pm 0.006	0.963 \pm 0.014	0.954 \pm 0.003	0.953 \pm 0.005	0.952 \pm 0.009 \downarrow	0.954 \pm 0.004 \downarrow	0.967 \pm 0.006

Figure 9: Different Concepts

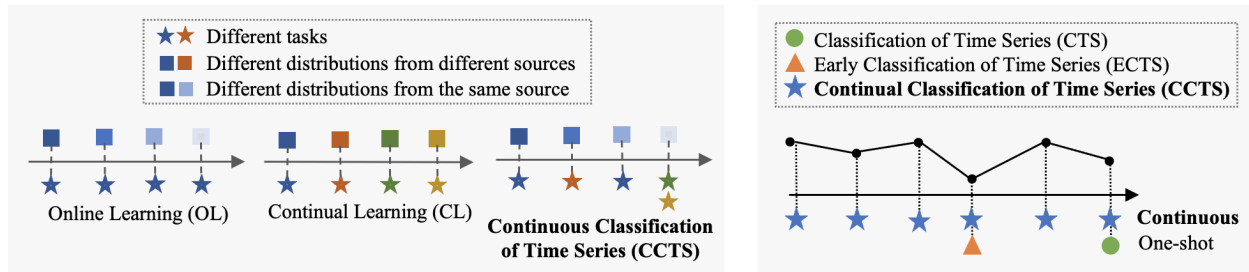


Figure 10: Multi-distribution in MIMIC-III Dataset

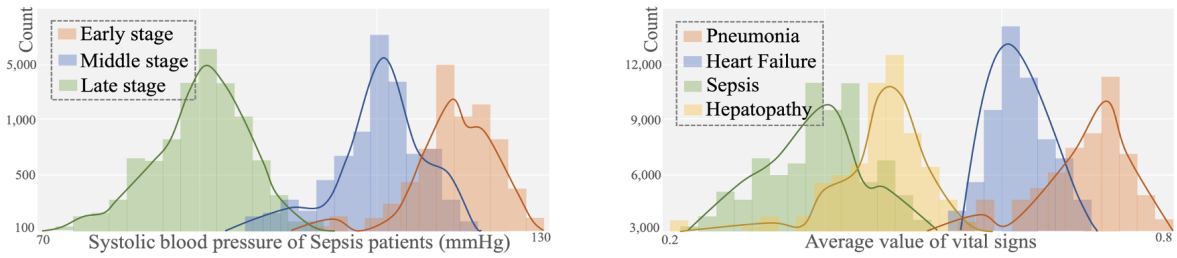


Figure 11: The Important Samples in SEPSIS Distribution Buffers

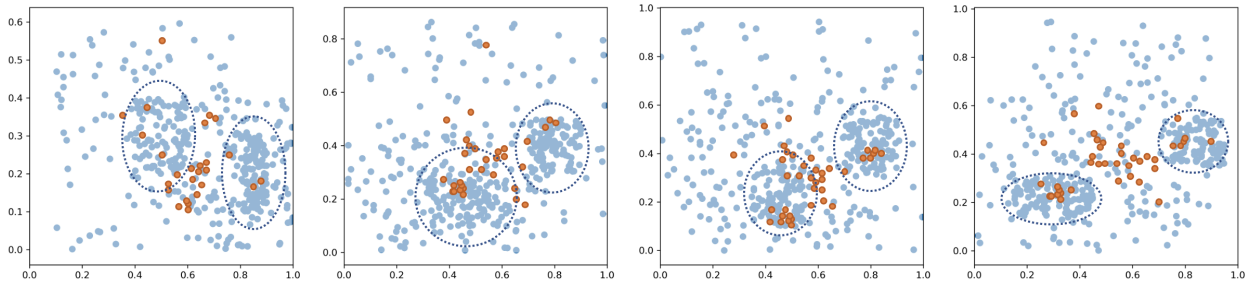


Figure 12: Task Similarity

In MIMIC-III dataset, the diagnoses with ICD-9 order are 1:HIV, 2:Brain Cancer, 3:Diabetes, 4:Hypertension, 5:Heart Failure, 6:Pneumonia, 7:Gastric Ulcer, 8:Hepatopathy, 9:Nephropathy, 10:SIRS. The new similarity order is 1, 10, 2, 4, 5, 8, 3, 9, 7, 6.

