

DDG-DA: Data Distribution Generation for Predictable Concept Drift Adaptation

Wendi Li^{1,2*}, Xiao Yang^{2*}, Weiqing Liu², Yingce Xia², Jiang Bian²

¹ University of Wisconsin–Madison

² Microsoft Research

Wendi.Li@wisc.edu, {Xiao.Yang, Weiqing.Liu, Yingce.Xia, Jiang.Bian}@microsoft.com

Abstract

In many real-world scenarios, we often deal with streaming data that is sequentially collected over time. Due to the non-stationary nature of the environment, the streaming data distribution may change in unpredictable ways, which is known as concept drift. To handle concept drift, previous methods first detect when/where the concept drift happens and then adapt models to fit the distribution of the latest data. However, there are still many cases that some underlying factors of environment evolution are predictable, making it possible to model the future concept drift trend of the streaming data, while such cases are not fully explored in previous work. In this paper, we propose a novel method DDG-DA, that can effectively forecast the evolution of data distribution and improve the performance of models. Specifically, we first train a predictor to estimate the future data distribution, then leverage it to generate training samples, and finally train models on the generated data. We conduct experiments on three real-world tasks (forecasting on stock price trend, electricity load and solar irradiance) and obtain significant improvement on multiple widely-used models.

Introduction

Machine learning algorithms have been widely applied to many real-world applications. One of the foundations of its success lies in assuming that the data is independent and identically distributed (briefly, the i.i.d. assumption). However, such an assumption is not always true for those learning tasks under broadly existing scenarios with streaming data. In a typical learning task over streaming data, as shown in Figure 1, the data comes in a sequential mode, meaning that, at any timestamp t , the learning task can only observe the new coming information, i.e., $Data^{(t)}$, in addition to historical ones, i.e., $Data^{(t-1)}$, $Data^{(t-2)}$, ..., $Data^{(1)}$. The goal of the task is usually to predict a certain target related to $Data^{(t+1)}$. Due to the non-stationary nature of the real-world environment, the data distribution could keep changing with continuous data streaming. Such a phenomenon/problem is called *concept drift* (Lu et al. 2018), where the basic assumption is that concept drift happens un-

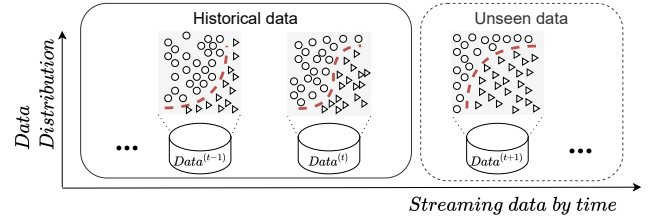


Figure 1: An example of concept drifts on streaming data. Triangle and circle represent two classes of data. Data comes like a stream. The data distribution changes over time.

expectedly and it is unpredictable (Gama et al. 2014) for streaming data.

To handle concept drift, previous studies usually leverage a two-step approach. Particularly, the first step is detecting the occurrence of concept drift in the streaming data, followed by the second step, if concept drift does occur, which adapts the model with new coming data by training a new model purely with latest data or making an ensemble between the new model with the current one, or fine-tuning the current model with latest data (Lu et al. 2018). Most of these studies share the same assumption that the latest data contains more useful information w.r.t. the upcoming streaming data (Zhao, Cai, and Zhou 2020), and thus the detection of concept drift and following model adaptation is mainly applied to the latest data.

Therefore, existing methods for handling concept drift will detect concept drift on the latest arrived data $Data^{(t)}$ at timestamp t and adapt the forecasting model accordingly. The concept drift continues, and the adapted model on $Data^{(t)}$ will be used on unseen streaming data in the future (e.g., $Data^{(t+1)}$). The previous model adaptation has a one-step delay to the concept drift of upcoming streaming data, which means a new concept drift has occurred between timestamp t and $t + 1$.

Nevertheless, apart from detecting unexpectedly occurred concept drift, most of the existing studies paid less attention to the scenarios that concept drift evolves in a gradual nonrandom way, which are in fact more common in streaming data. In Figure 2, we visualize the trace of concept drift in streaming data on three real-world tasks: forecasting on

*These authors contributed equally.

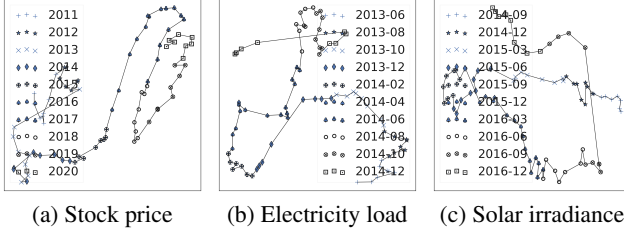


Figure 2: The visualization of concept drifts over time; though the concept drifts unexpectedly at some time, it drifts in a gradual nonrandom way most of the time. If the concept drifts are all sudden and do not follow any pattern, it may not be handled. However, in real-world scenarios, as we have shown, most concept drifts have a nonrandom trend rather than completely random.

stock price trend, electricity load and solar irradiance. Each point indicates the data distribution (a.k.a. concept) represented by several features in a specific time period and is plotted on the figure with t-SNE (Van der Maaten and Hinton 2008). To better visualize the concept drift, we connect points contiguous in time to a line. From this figure, we can find that unexpected drastic concept drifts only occur occasionally. **On the other hand, most of the points that are contiguous in time are spatially adjacent to each other and form a nonrandom trend, demonstrating a strong indication that a large portion of concept drifts do evolve in a gradual nonrandom way.** More importantly, such gradual nonrandom concept drifts are in some sense predictable since patterns exist in the concept drift trend.

In this paper, we focus on predictable concept drift by forecasting the future data distribution. We propose a novel method DDG-DA¹ to predict the data distribution of the next time-step sequentially, **such that the model of the downstream learning task can be trained on the data sample from the predicted distribution instead of catching up with the latest concept drift only.** In practice, DDG-DA is designed as a dynamic **data generator** that can create sample data from previously observed data by following predicted future data distribution. In other words, DDG-DA generates the resampling probability of each historical data sample to construct the future data distribution in estimation. However, it is quite challenging in reality to train this data generator to maximize the similarity between the predicted data distribution (represented by weighted resampling on historical data) and the ground truth future data distribution (represented by data in the future). To address this challenge, we propose to first represent a data distribution by learning a model under this data distribution and then create a differentiable distribution distance to train the data generator. To verify the effectiveness of this approach, we also conduct a thorough theoretical analysis to prove the equivalence between traditional distribution distance, e.g. KL-divergence, and the proposed differentiable distribution distance. Furthermore, it is worth not-

¹The code is available at https://github.com/Microsoft/qlib/tree/main/examples/benchmarks_dynamic/DDG-DA

ing that our proposed new approach is quite efficient since the distribution distance is differentiable and could be easy for optimization.

Extensive experiments have been conducted on a variety of popular streaming data scenarios with multiple widely-used learning models to evaluate the effectiveness of our proposed method DDG-DA. Experimental results have shown that DDG-DA could enhance the performance of learning tasks in all these scenarios. Further analysis also reveals that DDG-DA can successfully predict the future data distribution in the circumstance of predictable concept drift. The predicted data distribution could benefit the learning task by proactively adapting to drifting concepts. The code of DDG-DA is open-source on Github².

To sum up, our contributions include

- To the best of our knowledge, DDG-DA is the first method to **model the evolving of data distribution in predictable concept drift scenarios.**
- We create a differentiable distribution distance and provide theoretical analysis to prove its equivalence to KL-divergence distribution distance.
- We conduct extensive experiments on different real-world concept-drift-predictable scenarios to show that DDG-DA outperforms SOTA concept drift adaptation methods by training models on resampled dataset(predicted data distribution).

Background and Related Work

Streaming Data and Concept Drift

General Concept Drift In a large number of practical scenarios, data are collected over time sequentially. The streaming data could be formalized as $\mathbf{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(T)}\}$, where each element $\mathbf{x}^{(t)} \in \mathbb{R}^m$ pertaining to \mathbf{X} is an array of m values such that $\mathbf{x}^{(t)} = \{x_1, x_2, \dots, x_m\}$. Each one of the m scalar values corresponds to the input variable observed at time t . Given a target sequence $\mathbf{y} = \{y^{(1)}, y^{(2)}, \dots, y^{(T)}\}$ corresponding to \mathbf{X} , algorithms are designed to build the model on historical data $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^t$ and forecast \mathbf{y} on the future data $D_{test}^{(t)} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=t+1}^{t+\tau}$ at each timestamp t . Data are drawn from a joint distribution $(\mathbf{x}^{(t)}, y^{(t)}) \sim p_t(\mathbf{x}, y)$. The goal of algorithms is to build a model to make accurate predictions on unseen data $D_{test}^{(t)}$ in the future. Due to the non-stationary nature of the environment, **the joint distribution p_t is not static and changes over time**, which is a well-known problem named *concept drift*. **Formally, the concept drift between two timestamps t and $t+1$ can be defined as $\exists \mathbf{x} : p_t(\mathbf{x}, y) \neq p_{t+1}(\mathbf{x}, y)$.**

To improve the forecasting accuracy in the future, adapting models to accommodate the evolving data distribution is necessary, which is the focused problem in this paper. It can be formalized as

$$\min_{f^{(t)}, f^{(t+1)}, \dots, f^{(t+\tau)}} \sum_{i=t}^{t+\tau} l(f^{(i)}(\mathbf{x}^{(i)}), y^{(i)})$$

²<https://github.com/Microsoft/qlib>

where $f^{(t)}$ is learned from the training data in the previous stream $D_{train}^{(t)} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=t-k}^{t-1}$ with window size k and will be adapted continuously to minimize the target metric *loss* in given time period $[t, t + \tau]$. In many practical scenarios, the streaming data come continually over time and might never end. Accommodating such an amount of data in memory will be infeasible. Therefore only a limited memory size k is allowed for an online setting. Data in the memory could be used for model adaptation.

The Categorization of Concept Drifts and Our Scope

The concept drifts can be categorized into different types when based on different criteria (Ren et al. 2018). (Webb et al. 2016) provides quantitative measures to categorize concept drifts. For instance, a concept drift can be either *abrupt* or *gradual* according to its changing speed. It can also be *recurring* or *non-recurring* depending on whether similar concepts have appeared in the data stream.

Our method aims to handle scenarios with *predictable* concept drifts defined in (Minku, White, and Yao 2009), which refers to the concept drifts that follow a pattern rather than completely random due to the seasonality and reoccurring contexts. As (Ren et al. 2018; Žliobaitė, Pechenizkiy, and Gama 2016) state, the concept drifts in real-world evolving data are often caused by changes of the hidden contexts. Taking stock market, electricity load and solar irradiance as examples, because they are affected by some cyclic elements (time of the day, day of the week, month of the year, etc.), all of them have periodic, recurrent contexts. Abrupt drifts are rare due to nonrandom time-varying contexts, which is shown in Figure 2. Such reasons give these concept drifts nonrandom trends rather than completely unpredictable ones. (Khamassi et al. 2018; Webb et al. 2016) demonstrate that some characteristics of concept drift can be measured and predicted.

Related Work

Concept drift is a well-known topic in recent research works (Lu et al. 2018; Gama et al. 2014). The concept drift research works usually focus on streaming data, which come in a streaming form and are usually collected in non-stationary environments as Figure 1 shows. In the setting of most of the previous works, concept drift happens unexpectedly and is unpredictable. Without proactively predicting concept drift, these concept drift adaptation methods usually assume that the newer the data, the smaller the distribution gap between it and the future coming data. For instance, forgetting-based methods (Koychev 2000; Klinkenberg 2004) decay sample weights according to their age. When new labelled data are available, the model adaptation method detects concept drift and optionally updates the current model. However, the adapted new model will be used in the future, and the concept drift continues, which could make the model adaptation fail on unseen data in the future.

Previous Methods to Handle Predictable Concept Drift

Predictable concept drift is a scenario that the evolving concept drift is predictable to a certain extent **before concept drift happens**. In this scenario, methods are encouraged to

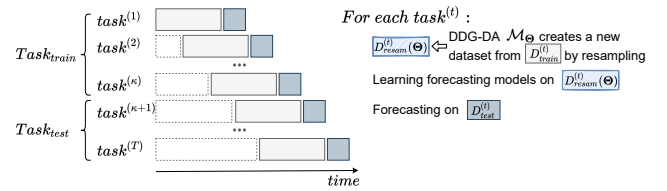


Figure 3: Training data (historical data) and test data (recent unseen data) change over time; the objective of each task is to improve the forecasting performance on test data.

predict and adapt to the new concept before drift happens. **Proactively adapting models to the future joint distribution $p(\mathbf{x}, y)$ is the key to catching up with the concept evolution.** Though most previous works assume that concept drift happens unexpectedly and is unpredictable. There are still a few of them noticing that concept drift is predictable (Minku, White, and Yao 2009; Khamassi et al. 2018) under certain additional assumptions. Some case studies in real-world scenarios demonstrate that the prediction of concept drift is possible (Gama et al. 2014). (Yang, Wu, and Zhu 2005) try to deal with recurrent concept drifts. **They use a small trigger window to detect the concept drift.** Based on the small amount of data in the window, they build a model to predict the new concept **after concept drift happens**. It is designed for the traditional concept drift setting instead of predictable concept drift. **DDG-DA tries to predict and adapt to the new concept before it happens.** As far as we know, we are the first method designed for predictable concept drift.

Types of Methods to Handle Concept Drift According to the learning mode, methods can be categorized into *re-training* (Bach and Maloof 2008) and *incremental adaptation* (Alippi and Roveri 2008). Retraining methods will discard the current model and learn a new model on a limited size of data. Incremental adaptation methods will update the current model incrementally.

From another perspective, most adaptive learning methods are designed for specific models. In lots of practical forecasting tasks, customized models will be designed specifically. These model-specific adaptive learning methods can only be carried out on specific forecasting models and are not general solutions, which narrows their practicability in real applications. Only a few adaptation methods are model-agnostic (Klinkenberg 2004; Koychev 2002, 2000; Žliobaitė 2009). DDG-DA focuses on predicting future data distribution, which is a model-agnostic solution and can benefit different customized forecasting models on streaming data.

Method Design

Overall Design

In streaming data, forecasting models are trained and adapted on historical data (training data $D_{train}^{(t)}$) and make predictions on unseen data (test data $D_{test}^{(t)}$) in the future. As shown in Figure 3, training data and test data change over time. For each timestamp t , the target of $task^{(t)}$:=

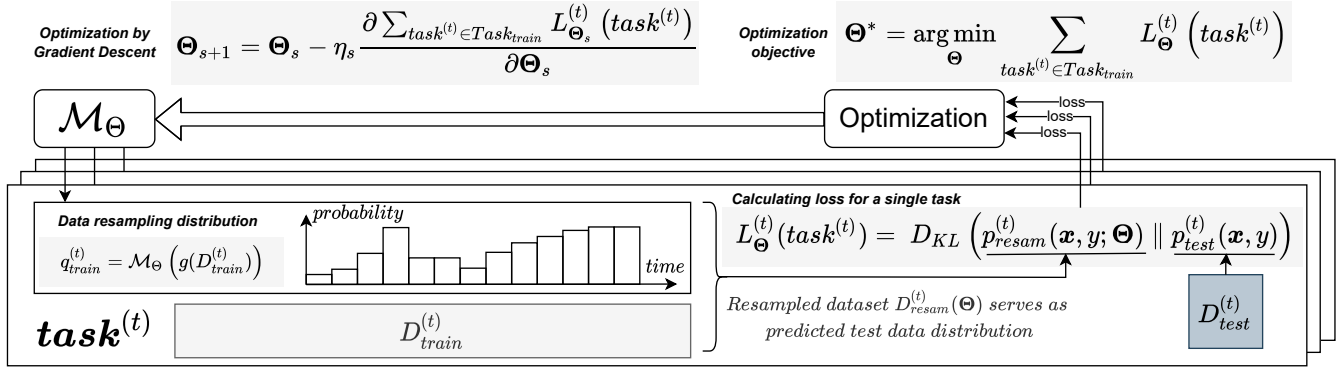


Figure 4: The learning process of DDG-DA; DDG-DA \mathcal{M}_{Θ} learns to guide the training process of forecasting model by generating dataset $D_{resam}^{(t)}(\Theta)$ resampled from $D_{train}^{(t)}$ with probability $q_{train}^{(t)}$. $q_{train}^{(t)}$ is the resampling probability given by \mathcal{M}_{Θ} at timestamp t .

$(D_{train}^{(t)}, D_{test}^{(t)})$ is to learn a new model or adapt an existing model on historical data $D_{train}^{(t)}$ and minimize the loss on $D_{test}^{(t)}$. The data come continuously and might never end. Therefore models can leverage a limited size of $D_{train}^{(t)}$ at timestamp t due to storage limitation. $D_{train}^{(t)}$ is a dataset sampled from training data distribution $p_{train}^{(t)}(\mathbf{x}, y)$ and $D_{test}^{(t)}$ from test data distribution $p_{test}^{(t)}(\mathbf{x}, y)$. The two distributions may be different (a.k.a. non-i.i.d.). This distribution gap is harmful to the forecasting accuracy on $D_{test}^{(t)}$ when learning forecasting models on $D_{train}^{(t)}$ with a distribution different from $D_{test}^{(t)}$.

DDG-DA Learning To bridge this gap, DDG-DA (annotated as \mathcal{M}_{Θ}) tries to model the concept drift and predict the **test data distribution** $p_{test}^{(t)}(\mathbf{x}, y)$. The framework of DDG-DA is demonstrated in Figure 3. DDG-DA will act like a weighted data sampler to resample on $D_{train}^{(t)}$ and create a new training dataset $D_{resam}^{(t)}(\Theta)$ whose data distribution is $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$ (the distribution of the resampled dataset serves as the prediction of test distribution). DDG-DA tries to minimize difference between the $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$ (the predicted data distribution) and the test data distribution $p_{test}^{(t)}(\mathbf{x}, y)$ (the ground truth data distribution).

During the training process of DDG-DA, \mathcal{M}_{Θ} tries to learn patterns on $task^{(t)} \in Task_{train}$ by minimizing the data distribution distance between $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$ and $p_{test}^{(t)}(\mathbf{x}, y)$. The knowledge learned by \mathcal{M}_{Θ} from $Task_{train}$ is expected to transfer to new tasks from $Task_{test}$.

DDG-DA Forecast For a given task $task^{(t)} \in Task_{test}$, the forecasting model is trained on dataset $D_{resam}^{(t)}(\Theta)$ under distribution $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$ and forecasts on test data $D_{test}^{(t)}$. $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$ is the distribution of resampled dataset $D_{resam}^{(t)}(\Theta)$ and the decisive factor to the trained forecasting model. With the help of DDG-DA, the distri-

bution of $D_{resam}^{(t)}(\Theta)$ is more similar to $D_{test}^{(t)}$ than $D_{train}^{(t)}$. The forecasting model's performance on $D_{test}^{(t)}$ is expected to be improved.

A Practical Example To make our method more understandable, we explain it with an example in the stock price trending forecasting scenario. To handle the concept drift in data, we retrain a new model each month (the rolling time interval is 1 month) based on two years of historical data (memory size is limited in an online setting). Each chance to retrain a new model to handle concept drift is called a *task*. For example, the first $task^{(2011/01)}$ contains training data $D_{train}^{(2011/01)}$ from 2009/01 to 2010/12 and a month of test data $D_{test}^{(2011/01)}$ in 2011/01. DDG-DA creates $D_{resam}^{(2011/01)}(\Theta)$ based on $D_{train}^{(2011/01)}$ to train a forecasting model to achieve better performance on $D_{test}^{(2011/01)}$. A new task is created in each month. These tasks are arranged in chronological order and separated at the beginning of 2016 into $Task_{train}$ (all $D_{test}^{(t)}$ in $Task_{train}$ range from 2011 to 2015) and $Task_{test}$ (all $D_{test}^{(t)}$ in $Task_{test}$ range from 2016 to 2020). DDG-DA is trained on $Task_{train}$ (learning to generate $D_{resam}^{(t)}(\Theta)$ and minimize its distribution distance with $D_{test}^{(t)}$). Then DDG-DA is evaluated on $Task_{test}$.

Model Design and Learning Process

The overall model design and learning process of DDG-DA are shown in Figure 4. DDG-DA will build a set of tasks $Task_{train} = \{task^{(1)}, \dots, task^{(\tau)}\}$ for training DDG-DA. The goal of the learned DDG-DA is to improve the performance on test tasks $Task_{test} := \{task^{(\tau+1)}, \dots, task^{(T)}\}$.

Feature Design DDG-DA is expected to guide the model learning process in each $task^{(t)}$ by forecasting test data distribution. Historical data distribution information is useful to predict the target distribution of $D_{test}^{(t)}$ and is input into DDG-DA. DDG-DA will learn concept drift patterns from training tasks and help to adapt models in test tasks.

DDG-DA could be formalized as $q_{train}^{(t)} =$

$\mathcal{M}_\Theta(g(D_{train}^{(t)}))$. g is a feature extractor. It takes $D_{train}^{(t)}$ as input and outputs historical data distribution information. \mathcal{M}_Θ leverages the extracted information and output the resampling probabilities for samples in $D_{train}^{(t)}$.

Objective Function \mathcal{M}_Θ accepts the extracted feature and outputs the resampling probability on $D_{train}^{(t)}$. The resampled dataset's joint distribution $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$ serves as the distribution prediction. The learning target of DDG-DA is to minimize the difference between $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$ and $p_{test}^{(t)}(\mathbf{x}, y)$. We focus on the most important concept drift subject $p(y | \mathbf{x})$ (Kelly, Hand, and Adams 1999) and assume the difference between $p_{test}^{(t)}(\mathbf{x})$ and $p_{resam}^{(t)}(\mathbf{x}; \Theta)$ are minor. The loss of DDG-DA could be reformulated as

$$\begin{aligned} L_\Theta(task^{(t)}) &= D_{KL}(p_{test}^{(t)}(\mathbf{x}, y) \parallel p_{resam}^{(t)}(\mathbf{x}, y; \Theta)) \\ &= \mathbb{E}_{\mathbf{x} \sim p_{test}^{(t)}(\mathbf{x})} [D_{KL}(p_{test}^{(t)}(y | \mathbf{x}) \parallel p_{resam}^{(t)}(y | \mathbf{x}; \Theta))] \end{aligned} \quad (1)$$

where D_{KL} represents the Kullback–Leibler divergence.

Normal distribution assumption is reasonable for unknown variables and often used in maximum likelihood estimation (Goodfellow et al. 2016). Under this assumption, $p_{test}^{(t)}(y | \mathbf{x}) = \mathcal{N}(y_{test}^{(t)}(\mathbf{x}), \sigma)$ and $p_{resam}^{(t)}(y | \mathbf{x}) = \mathcal{N}(y_{resam}^{(t)}(\mathbf{x}; \Theta), \sigma)$ where σ is a constant.

The estimated expectation of Equation (1) on empirical sampled dataset can be reformulated as

$$L_\Theta(task^{(t)}) = \frac{1}{2} \sum_{(\mathbf{x}, y) \in D_{test}^{(t)}} \|y_{resam}^{(t)}(\mathbf{x}; \Theta) - y\|^2 \quad (2)$$

Summarizing losses of all the training tasks, the optimization target of DDG-DA could be formalized as

$$\Theta^* = \arg \min_{\Theta} \sum_{task^{(t)} \in Task_{train}} L_\Theta(task^{(t)}) \quad (3)$$

DDG-DA learns knowledge from $Task_{train}$ and transfers it to unseen test tasks $Task_{test}$. In each task, DDG-DA forecasts the future data distribution and generate dataset $D_{resam}^{(t)}(\Theta)$. Learning the forecasting models on $D_{resam}^{(t)}(\Theta)$, it adapts to upcoming streaming data better.

Optimization In this section, we will introduce the optimization process of (2). $y_{test}^{(t)}(\mathbf{x})$ can be get directly from dataset $D_{test}^{(t)}$. To approximate $y_{resam}^{(t)}(\mathbf{x}; \Theta)$, DDG-DA builds a regression proxy model $y_{proxy}(\mathbf{x}; \phi^{(t)})$ on $D_{resam}^{(t)}(\Theta)$. The optimization of the proxy model can be formulated as

$$\phi^{(t)} = \arg \min_{\phi} \sum_{(\mathbf{x}, y) \in D_{resam}^{(t)}(\Theta)} \|y_{proxy}(\mathbf{x}; \phi) - y\|^2 \quad (4)$$

The learning process of DDG-DA becomes a *bi-level optimization problem* (Gould et al. 2016). The goal of the upper-level is Equation (2) and (3) ($y_{resam}^{(t)}(\mathbf{x}; \Theta)$ is replaced by

$y_{proxy}(\mathbf{x}; \phi^{(t)})$). The goal of the lower-level optimization Equation (4) can be regarded as a constraint.

The overall bi-level optimization formulation of the DDG-DA is

$$\begin{aligned} \arg \min_{\Theta} \quad & \sum_{task^{(t)} \in Task_{train}} \left(\sum_{(\mathbf{x}, y) \in D_{test}^{(t)}} \|y_{proxy}(\mathbf{x}; \phi^{(t)}) - y\|^2 \right) \\ \text{s.t.} \quad & \phi^{(t)} = \arg \min_{\phi} \sum_{(\mathbf{x}', y') \in D_{resam}^{(t)}(\Theta)} \|y_{proxy}(\mathbf{x}'; \phi) - y'\|^2 \end{aligned} \quad (5)$$

where Θ is the parameters of DDG-DA.

Many methods have been proposed to solve such a bi-level optimization problem (Gould et al. 2016). Some methods based on hyperparameter optimization (Bengio 2000; Baydin and Pearlmutter 2014; Maclaurin, Duvenaud, and Adams 2015) are proposed. But they have limitations on either the network size or optimization accuracy.

One of the key barriers that stop researchers from optimizing Equation (2) directly is that the lower-level optimization Equation (4) usually can not be solved in a closed-form. The arg min operator is usually not differentiable, which makes some popular and efficient optimization algorithms (such as gradient-descend-based methods) impossible in the optimization of the upper-level (Equation (3)). (Lee et al. 2019; Bertinetto et al. 2018) argue that (Bengio 2000; Baydin and Pearlmutter 2014) it is too costly and adopt algorithms with closed-form solutions in the lower-level optimization.

DDG-DA adopts a model with a closed-form solution as $y_{proxy}^{(t)}(\mathbf{x}; \phi^{(t)})$. We have many choices, such as logistic regression (Kleinbaum et al. 2002), kernel-based non-linear model (Liu 2003) and differentiable closed-form solvers (Bertinetto et al. 2018). We choose a linear model $h(\mathbf{x}; \phi^{(t)}) = \mathbf{x} \phi^{(t)}$ for $y_{proxy}(\mathbf{x}; \phi^{(t)})$ for simplicity. The resampling probability $q_{train}^{(t)}$ outputted by \mathcal{M}_Θ could be regarded as sample weights when learning forecasting models. The loss function in Equation (4) could be formulated as

$$\begin{aligned} l_{\phi^{(t)}}(D_{resam}^{(t)}(\Theta)) &= \frac{1}{2} \sum_{(\mathbf{x}, y) \in D_{train}^{(t)}} q_{train}^{(t)}(\mathbf{x} \phi^{(t)} - y)^2 \\ &= \frac{1}{2} (\mathbf{X}^{(t)} \phi^{(t)} - \mathbf{y}^{(t)})^\top \mathbf{Q}^{(t)} (\mathbf{X}^{(t)} \phi^{(t)} - \mathbf{y}^{(t)}) \end{aligned}$$

where $\mathbf{X}^{(t)}$, $\mathbf{y}^{(t)}$ and $\mathbf{Q}^{(t)}$ represent the concatenated features, labels and resampling probability in $D_{train}^{(t)}$.

It equals to a *weighted linear regression* problem (Ruppert and Wand 1994). $\phi^{(t)}$ has a closed-form solution formalized as

$$\begin{aligned} \phi^{(t)} &= \arg \min_{\phi} l_{\phi}(D_{resam}^{(t)}(\Theta)) \\ &= \left((\mathbf{X}^{(t)})^\top \mathbf{Q}^{(t)} \mathbf{X}^{(t)} \right)^{-1} (\mathbf{X}^{(t)})^\top \mathbf{Q}^{(t)} \mathbf{y}^{(t)} \end{aligned} \quad (6)$$

This closed-form solution to Equation (4) makes the distribution distance differentiable. It makes the optimization objective of \mathcal{M}_Θ (Equation (3)) differentiable. Therefore, simple and efficient optimization algorithms can be used to train DDG-DA (e.g. stochastic gradient descent).

Method	Stock Price Trend Forecasting					Electricity Load		Solar Irradiance		
	<i>IC</i>	<i>ICIR</i>	<i>Ann.Ret.</i>	<i>Sharpe</i>	<i>MDD</i>	<i>NMAE</i>	<i>NRMSE</i>	<i>Skill (%)</i>	<i>MAE</i>	<i>RMSE</i>
RR	0.1178	1.0658	0.1749	1.5105	-0.2907	0.1877	0.9265	7.3047	21.7704	48.0117
GF-Lin	0.1227	<u>1.0804</u>	0.1739	1.4590	-0.2690	0.1843	0.9109	<u>9.3503</u>	21.6878	<u>46.9522</u>
GF-Exp	0.1234	1.0613	0.1854	1.5906	-0.2984	0.1839	0.9084	9.2652	21.6841	46.9963
ARF	0.1240	1.0657	0.1994	1.8844	-0.1176	<u>0.1733</u>	<u>0.8901</u>	8.6267	<u>21.0962</u>	47.3270
Condor	<u>0.1273</u>	1.0635	<u>0.2157</u>	<u>2.1105</u>	-0.1624					
DDG-DA	0.1312	1.1299	0.2565	2.4063	<u>-0.1381</u>	0.1622	0.8498	12.1327	18.7997	45.5110

Table 1: Performance comparison of the concept drifts adaptation methods.

Experiments

In this section, we conduct experiments aiming to answer the following research questions:

- **Q1:** Can DDG-DA outperform SOTA concept drift adaptation methods in predictable concept drift scenarios?
- **Q2:** Can DDG-DA generalize to different forecasting models in different scenarios?

Experiment Setup

The experiments are conducted on multiple datasets in three real-world popular scenarios (forecasting on stock price trend, electricity load and solar irradiance (Grinold and Kahn 2000; Pedro, Larson, and Coimbra 2019)).

Experiments Results

Concept Drift Adaptation Methods Comparison (Q1)

In this part, we compare DDG-DA with concept drift adaptation methods in different streaming data scenarios to answer Q1 that DDG-DA is the SOTA in concept-drift-predictable scenarios. We compared all methods both in classification tasks (stock price trend forecasting) and regression tasks (electricity load forecasting, solar irradiance forecasting).

The following methods are compared:

- **RR:** Periodically **R**olling **R**etrain model on data in memory with equal weights. The memory only stores the most recent data in a limited window size.
- **GF-Lin** (Koychev 2000): Based on RR, **G**radual **F**orgetting by weights decaying **L**inearly by time.
- **GF-Exp** (Klinkenberg 2004): Based on RR, **G**radual **F**orgetting by weights decaying **E**xponentially by time.
- **ARF** (Gomes et al. 2017, 2018): To deal with concept drift in the streaming data, **A**daptive **R**andom **F**orest does both internal and external concept drift detecting for each newly-created tree. The final prediction will be obtained by its voting strategy.
- **Condor** (Zhao, Cai, and Zhou 2020): **C**oncept **D**rift via **R**euse is an ensemble method that handles non-stationary environments by both building new models and assigning weights for previous models.
- **DDG-DA:** Our proposed method. DDG-DA is based on the setting of RR and generates a new dataset by resampling to retrain forecasting models.

The compared methods can be categorized into two sets based on if they are model-agnostic. A model-agnostic solution can conduct concept drift adaptation without knowing the details of the forecasting model. Therefore, a model-agnostic solution can be applied to any model, which could be designed specifically to target scenarios. RR, GF-Lin, GF-Exp, DDG-DA are model-agnostic. They use the same forecasting model and historical data memory size. ARF and Condor are not model-agnostic. Condor is designed for classification, so it is not evaluated on the regression scenarios.

We use the open-source models released in Qlib’s model zoo as our candidate forecasting model. On the validation data, LightGBM (Ke et al. 2017) performs best on average. Therefore, we select LightGBM as our final forecasting model. The same forecasting model is also adopted by RR, GF-Lin and GF-Exp. For the training process of DDG-DA, we create tasks according to Figure 3. The tasks are created in a more coarse-grained frequency to reduce the retraining cost. The time interval of two adjacent tasks is 20 trading days, 7 days and 7 days in stock price trend forecasting, electricity load forecasting and solar irradiance forecasting respectively. The test data still arrives in fine granularity.

The results are shown in Table 1. RR is the most naive method and simply periodically retrains models on data in memory. It performs worst among the compared methods. GF-Lin and GF-Exp assume that the recent data distribution is more similar to the upcoming data. The weights of the most recent data are the highest and then decay by time. Such an assumption is valid in most datasets, and they outperform RR in most scenarios. ARF and Condor are the most recent SOTA solutions for concept drift adaptation. They adapt models to the most recent data to handle concept drift. The forecasting performance is improved further. However, the concept drift continues after adapting models to the most recent historical data. The adapted model could fail again when the concept drifts in the future. DDG-DA solves such a problem and performs best. It models the trend of concept drifts and generates a new dataset whose distribution is closer to that in the future. Then the forecasting model is trained on the new dataset to handle the future concept drift.

Generalization to Different Forecasting Models (Q2)

To answer Q2, we conduct experiments to demonstrate that DDG-DA can enhance the performance of different forecasting models in different streaming data scenarios.

The experiments involve different forecasting models. Therefore, we only compare model-agnostic concept drift

Model	Method	Stock Price Trend Forecasting					Electricity Load		Solar Irradiance		
		<i>IC</i>	<i>ICIR</i>	<i>Ann.Ret.</i>	<i>Sharpe</i>	<i>MDD</i>	<i>NMAE</i>	<i>NRMSE</i>	<i>Skill (%)</i>	<i>MAE</i>	<i>RMSE</i>
Linear	RR	0.0859	0.7946	0.1578	1.4211	-0.1721	0.2080	1.0207	1.4133	24.0910	51.0632
	GF-Exp	0.0863	0.8018	0.1632	1.4373	-0.1462	0.2009	0.9787	1.3660	24.1391	51.0877
	DDG-DA	0.0971	0.9193	0.1763	1.5733	-0.2130	0.1973	0.9702	3.7378	23.6901	49.8592
MLP	RR	0.1092	0.8647	0.1803	1.5200	-0.1797	0.1928	0.9682	8.2145	22.1285	47.5405
	GF-Exp	0.1091	0.8654	0.1889	1.6121	-0.1738	0.1898	0.9588	8.4668	21.6236	47.4098
	DDG-DA	0.1211	0.9921	0.2181	1.9409	-0.1864	0.1882	0.9537	10.0341	20.3422	46.5980
Light-GBM	RR	0.1178	1.0658	0.1749	1.5105	-0.2907	0.1877	0.9265	7.3047	21.7704	48.0117
	GF-Exp	0.1234	1.0613	0.1854	1.5906	-0.2984	0.1839	0.9084	9.2652	21.6841	46.9963
	DDG-DA	0.1312	1.1299	0.2565	2.4063	-0.1381	0.1622	0.8498	12.1327	18.7997	45.5110
LSTM	RR	0.1003	0.7066	0.1899	1.8919	-0.1264	0.1494	0.8386	7.5990	21.7948	49.8593
	GF-Exp	0.1008	0.7110	0.1928	1.9238	-0.1450	0.1370	0.7369	11.4428	18.5193	45.8684
	DDG-DA	0.1049	0.7456	0.2122	2.0334	-0.1745	0.1298	0.7290	12.4905	18.3295	45.3257
GRU	RR	0.1122	0.9638	0.1841	1.6645	-0.1740	0.1352	0.7090	8.5684	21.0297	47.3572
	GF-Exp	0.1182	1.0007	0.1872	1.7588	-0.1207	0.1281	0.6688	8.4578	21.0399	47.4145
	DDG-DA	0.1183	1.0091	0.1928	1.7906	-0.1182	0.1250	0.6588	10.5918	20.1891	46.3092

Table 2: Performance comparison of the model-agnostic solutions on various scenarios and forecasting models.

adaptation solutions. For each forecasting model in each scenario, we compare DDG-DA with the RR and GF-Exp (GF-Exp is an advanced version of GF-Lin. So GF-Lin is not included). We conduct experiments on multiple popular forecasting models, including **Linear** (Graybill 1976), **MLP** (Gardner and Dorling 1998), **LightGBM** (Ke et al. 2017), **LSTM** (Hochreiter and Schmidhuber 1997) and **GRU** (Chung et al. 2014).

As the results shown in Table 2, DDG-DA outperforms others in most cases. The results demonstrate that DDG-DA has captured the pattern of concept drift over time, which is the inherent nature of data and model-agnostic.

Comparison of Optimization Methods DDG-DA converts the distribution distance optimization into a bi-level optimization by using a proxy model to model data distribution $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$. DDG-DA adapts a proxy model with a closed-form solution and convert Equation (4) into a differentiable operator. Therefore, the distribution distance between $D_{resam}^{(t)}(\Theta)$ and $D_{test}^{(t)}$ becomes differentiable. Besides this approach, hyperparameter-optimization-based methods are eligible for the bi-level optimization. Gradient-based Hyperparameter Optimization (GHO) is one of them and adopted by a lot of research works (Maclaurin, Duvenaud, and Adams 2015; Fan et al. 2020; Baydin et al. 2017). Instead of requiring a proxy model with a closed-form solution, GHO only requires a differentiable proxy model. It initializes the proxy model’s parameters with $\phi_{(t),0}$ and then update it to $\phi_{(t),k}$ by k -steps with a gradient-based optimization method. At last, the $\phi_{(t),k}$ serves as $\phi_{(t)}$ in Equation (4) and makes Equation (2) differentiable.

To compare different bi-level optimization algorithms, we create a new solution named DDG-DA[†] based on DDG-DA with GHO for optimization and an MLP proxy model. As the results show in Table 3, DDG-DA[†] outperforms RR, GF-Lin and GF-Exp. It demonstrates the effectiveness of

Method	<i>IC</i>	<i>ICIR</i>	<i>Ann.Ret.</i>	<i>Sharpe</i>	<i>MDD</i>
RR	0.1092	0.8647	0.1803	1.5200	-0.1797
GF-Lin	0.1090	0.8641	0.1880	1.5682	-0.1656
GF-Exp	0.1091	0.8654	0.1889	1.6121	-0.1738
DDG-DA [†]	0.1204	0.9732	0.2065	1.9243	-0.1705
DDG-DA	0.1211	0.9921	0.2181	1.9409	-0.1864

Table 3: Comparison of DDG-DA with different bi-level optimization algorithms and proxy model; DDG-DA[†] is based on GHO and an MLP proxy model.

our framework. However, it slightly underperforms DDG-DA. DDG-DA[†] leverages a model with a larger capacity to model data distribution. It alleviates the problem of underfitting $p_{resam}^{(t)}(\mathbf{x}, y; \Theta)$. But the gap between the proxy model and the ground truth distribution still exists. The hyperparameter k is sensitive. A smaller k may result in underfitting. A greater k may waste computing resources and tend to overfit.

Conclusions and Future Works

In this paper, we propose a novel concept drift adaptation method DDG-DA to adapt models to future data distribution rather than latest historical data only like previous works. DDG-DA models the concept drift and generates a new dataset by resampling historical data to guide the training process. Experiments on diverse scenarios demonstrate its effectiveness in concept-drift-predictable scenarios. DDG-DA is model-agnostic, so it can be transferred to any forecasting model used for downstream tasks.

Currently, DDG-DA is a static model that assumes the pattern of concept drift is static. This may not be valid in some scenarios (e.g. the concept drift frequency of $Task_{test}$ and $Task_{train}$ may be different). For future work, we will improve DDG-DA to a dynamic one to solve this limitation.

References

- Alippi, C.; and Roveri, M. 2008. Just-in-time adaptive classifiers—Part II: Designing the classifier. *IEEE Transactions on Neural Networks*, 19(12): 2053–2064.
- Bach, S. H.; and Maloof, M. A. 2008. Paired learners for concept drift. In *2008 Eighth IEEE International Conference on Data Mining*, 23–32. IEEE.
- Baydin, A. G.; Cornish, R.; Rubio, D. M.; Schmidt, M.; and Wood, F. 2017. Online learning rate adaptation with hyper-gradient descent. *arXiv preprint arXiv:1703.04782*.
- Baydin, A. G.; and Pearlmutter, B. A. 2014. Automatic differentiation of algorithms for machine learning. *arXiv preprint arXiv:1404.7456*.
- Bengio, Y. 2000. Gradient-based optimization of hyperparameters. *Neural computation*, 12(8): 1889–1900.
- Bertinetto, L.; Henriques, J. F.; Torr, P. H.; and Vedaldi, A. 2018. Meta-learning with differentiable closed-form solvers. *arXiv preprint arXiv:1805.08136*.
- Chung, J.; Gulcehre, C.; Cho, K.; and Bengio, Y. 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*.
- Fan, Y.; Xia, Y.; Wu, L.; Xie, S.; Liu, W.; Bian, J.; Qin, T.; Li, X.-Y.; and Liu, T.-Y. 2020. Learning to teach with deep interactions. *arXiv preprint arXiv:2007.04649*.
- Gama, J.; Žliobaitė, I.; Bifet, A.; Pechenizkiy, M.; and Bouchachia, A. 2014. A survey on concept drift adaptation. *ACM computing surveys (CSUR)*, 46(4): 1–37.
- Gardner, M. W.; and Dorling, S. 1998. Artificial neural networks (the multilayer perceptron)—a review of applications in the atmospheric sciences. *Atmospheric environment*, 32(14-15): 2627–2636.
- Gomes, H. M.; Barddal, J. P.; Ferreira, L. E. B.; and Bifet, A. 2018. Adaptive random forests for data stream regression. In *ESANN*.
- Gomes, H. M.; Bifet, A.; Read, J.; Barddal, J. P.; Enembreck, F.; Pfahringer, B.; Holmes, G.; and Abdesslem, T. 2017. Adaptive random forests for evolving data stream classification. *Machine Learning*, 106(9-10): 1469–1495.
- Goodfellow, I.; Bengio, Y.; Courville, A.; and Bengio, Y. 2016. *Deep learning*. 2. MIT press Cambridge.
- Gould, S.; Fernando, B.; Cherian, A.; Anderson, P.; Cruz, R. S.; and Guo, E. 2016. On differentiating parameterized argmin and argmax problems with application to bi-level optimization. *arXiv preprint arXiv:1607.05447*.
- Graybill, F. A. 1976. *Theory and application of the linear model*, volume 183. Duxbury press North Scituate, MA.
- Grinold, R. C.; and Kahn, R. N. 2000. *Active portfolio management*. McGraw Hill New York.
- Hochreiter, S.; and Schmidhuber, J. 1997. Long short-term memory. *Neural computation*, 9(8): 1735–1780.
- Ke, G.; Meng, Q.; Finley, T.; Wang, T.; Chen, W.; Ma, W.; Ye, Q.; and Liu, T.-Y. 2017. Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in neural information processing systems*, 3146–3154.
- Kelly, M. G.; Hand, D. J.; and Adams, N. M. 1999. The impact of changing populations on classifier performance. In *Proceedings of the fifth ACM SIGKDD international conference on Knowledge discovery and data mining*, 367–371.
- Khamassi, I.; Sayed-Mouchaweh, M.; Hammami, M.; and Ghédira, K. 2018. Discussion and review on evolving data streams and concept drift adapting. *Evolving systems*, 9(1): 1–23.
- Kleinbaum, D. G.; Dietz, K.; Gail, M.; Klein, M.; and Klein, M. 2002. *Logistic regression*. Springer.
- Klinkenberg, R. 2004. Learning drifting concepts: Example selection vs. example weighting. *Intelligent data analysis*, 8(3): 281–300.
- Koychev, I. 2000. Gradual forgetting for adaptation to concept drift. In *Proceedings of ECAI 2000 Workshop on Current Issues in Spatio-Temporal Reasoning*.
- Koychev, I. 2002. Tracking changing user interests through prior-learning of context. In *International Conference on Adaptive Hypermedia and Adaptive Web-Based Systems*, 223–232. Springer.
- Lee, K.; Maji, S.; Ravichandran, A.; and Soatto, S. 2019. Meta-learning with differentiable convex optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 10657–10665.
- Liu, B. 2003. Kernel-based nonlinear discriminator with closed-form solution. In *International Conference on Neural Networks and Signal Processing, 2003. Proceedings of the 2003*, volume 1, 41–44. IEEE.
- Lu, J.; Liu, A.; Dong, F.; Gu, F.; Gama, J.; and Zhang, G. 2018. Learning under concept drift: A review. *IEEE Transactions on Knowledge and Data Engineering*, 31(12): 2346–2363.
- MacLaurin, D.; Duvenaud, D.; and Adams, R. 2015. Gradient-based hyperparameter optimization through reversible learning. In *International conference on machine learning*, 2113–2122. PMLR.
- Minku, L. L.; White, A. P.; and Yao, X. 2009. The impact of diversity on online ensemble learning in the presence of concept drift. *IEEE Transactions on knowledge and Data Engineering*, 22(5): 730–742.
- Pedro, H. T.; Larson, D. P.; and Coimbra, C. F. 2019. A comprehensive dataset for the accelerated development and benchmarking of solar forecasting methods. *Journal of Renewable and Sustainable Energy*, 11(3): 036102.
- Ren, S.; Liao, B.; Zhu, W.; and Li, K. 2018. Knowledge-maximized ensemble algorithm for different types of concept drift. *Information Sciences*, 430: 261–281.
- Ruppert, D.; and Wand, M. P. 1994. Multivariate locally weighted least squares regression. *The annals of statistics*, 1346–1370.
- Van der Maaten, L.; and Hinton, G. 2008. Visualizing data using t-SNE. *Journal of machine learning research*, 9(11).
- Webb, G. I.; Hyde, R.; Cao, H.; Nguyen, H. L.; and Petitjean, F. 2016. Characterizing concept drift. *Data Mining and Knowledge Discovery*, 30(4): 964–994.

- Yang, Y.; Wu, X.; and Zhu, X. 2005. Combining proactive and reactive predictions for data streams. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, 710–715.
- Zhao, P.; Cai, L.-W.; and Zhou, Z.-H. 2020. Handling concept drift via model reuse. *Machine Learning*, 109(3): 533–568.
- Žliobaitė, I. 2009. Combining time and space similarity for small size learning under concept drift. In *International Symposium on Methodologies for Intelligent Systems*, 412–421. Springer.
- Žliobaitė, I.; Pechenizkiy, M.; and Gama, J. 2016. An overview of concept drift applications. *Big data analysis: new algorithms for a new society*, 91–114.