

Learning Parameter Distributions to Detect Concept Drift in Data Streams

Johannes Haug
University of Tuebingen
Tuebingen, Germany

johannes-christian.haug@uni-tuebingen.de

Gjergji Kasneci
University of Tuebingen
Tuebingen, Germany

gjergji.kasneci@uni-tuebingen.de

Abstract—Data distributions in streaming environments are usually not stationary. In order to maintain a high predictive quality at all times, online learning models need to adapt to distributional changes, which are known as concept drift. The timely and robust identification of concept drift can be difficult, as we never have access to the true distribution of streaming data. In this work, we propose a novel framework for the detection of real concept drift, called *ERICs*. By treating the parameters of a predictive model as random variables, we show that concept drift corresponds to a change in the distribution of optimal parameters. To this end, we adopt common measures from information theory. The proposed framework is completely model-agnostic. By choosing an appropriate base model, *ERICs* is also capable to detect concept drift at the input level, which is a significant advantage over existing approaches. An evaluation on several synthetic and real-world data sets suggests that the proposed framework identifies concept drift more effectively and precisely than various existing works.

I. INTRODUCTION

Data streams are a potentially unbounded sequence of observations. As such, data streams are subject to a number of external factors, e.g. seasonal or catastrophic events. Hence, the distributions of a data stream are usually not stationary, but change over time, which is known as concept drift.

Concept drift can seriously affect the quality of predictions, if it goes unnoticed. Concept drift detection models help identify and handle distributional changes, allowing us to maintain a high predictive performance over time. Ideally, concept drift detection models are sensitive enough to detect drift with only a short delay. However, concept drift detection should also be robust against small perturbations of the input in order to avoid false positives and thus be reliable.

Let X and Y be random variables that correspond to the streaming observations and the associated labels. According to [1], concept drift resembles a difference in the joint probability $P(Y, X)$ at different time steps $t, u \in \{1, \dots, T\}$, i.e.

$$\begin{aligned} P_t(Y, X) &\neq P_u(Y, X) \\ \Leftrightarrow P_t(Y|X)P_t(X) &\neq P_u(Y|X)P_u(X). \end{aligned}$$

We call $P_t(Y, X)$ the active concept at time step t . Moreover, we distinguish between real and virtual concept drift. Virtual concept drift describes a change in $P(X)$, i.e. $P_t(X) \neq P_u(X)$. Hence, virtual concept drift is independent from the target distribution and does not change the decision boundary

[2]. On the other hand, real concept drift, sometimes called concept shift, corresponds to a change in the conditional target distribution, i.e. $P_t(Y|X) \neq P_u(Y|X)$. Real concept drift shifts the decision boundary, which may influence subsequent predictions [2]. It is therefore crucial to detect changes of $P(Y|X)$ in time to avoid dramatic drops in predictive performance. In this paper, we investigate the effective and robust identification of real concept drift.

Unfortunately, concept drift does not follow a clear pattern in practice. Instead, we might observe large differences in the duration and magnitude of concept drift. To this end, we distinguish between different types of concept drift [1]–[3]: Sudden drift describes an abrupt change from one concept to another. Incremental drift is a steady transition of concepts over some time period. In a gradual drift, the concepts alternate temporarily, until a new concept ultimately replaces the old one. Sometimes we also observe mixtures of different concept drift types and recurring or cyclic concepts. For further information, we refer the fellow reader to [1]. In general, concept drift detection models should allow timely and accurate detection of all types of concept drift.

In a data stream, we can only access a fraction of the data at every time step t . To detect real concept drift, we thus need to approximate $P_t(Y|X)$, by using a predictive model f_{θ_t} . Accordingly, we get $P_t(Y|X) \approx P(Y|X, \theta_t)$, with parameters $\theta_t = (\theta_{tk})_{k=1}^K$. We optimize the model parameters, given the new observations in every time step. Consequently, θ_t represents our most current information about the active concept at time step t . A concept drift detection model should therefore adhere to changes of the model parameters through the following two properties:

Property 1. Model-Aware Concept Drift Detection. Let θ_t, θ_u be the parameters of a predictive model f_{θ} at two time steps t and u . Let further \mathcal{D} be a statistical divergence measure (e.g., Kullback–Leibler, Jensen–Shannon, etc.). Concept drift detection is *model-aware*, if for a detected drift between any two time steps t and u , we observe $\mathcal{D}(\theta_t, \theta_u) > 0$.

Accordingly, we associate concept drift with updates of the predictive model f_{θ} . Given that f_{θ} is robust, model-awareness reduces the sensitivity of a concept drift detection scheme to random input perturbations, which in turn reduces the risk of false alarms.

arXiv:2010.09388v1 [cs.LG] 19 Oct 2020

Property 2. Explainable Concept Drift Detection. Concept drift detection at time step t is *explainable* with respect to the predictive model f_{θ_t} , if the concept drift can be associated with individual model parameters, i.e. each dimension of θ_t .

If we associate concept drift with individual parameters, we can make more targeted model updates. Hence, we may avoid unnecessary and costly adaptations of the predictive model. Moreover, some parameter distributions even allow us to relate concept drift to specific input features. In this way, concept drift becomes much more transparent.

In this paper, we propose a novel framework for *Effective and Robust Identification of Concept Shift (ERICS)*. *ERICS* complies with the Properties 1 and 2. We use the probabilistic framework introduced in [4] to model the distribution of the parameters θ at every time step. Specifically, we express real concept drift in terms of the marginal likelihood and the parameter distribution $P(\theta; \psi)$, which is itself parameterized by ψ . Unlike many existing models, *ERICS* does not need to access the streaming data directly [5]. Instead, we detect concept drift by investigating the differential entropy and Kullback-Leibler (KL) divergence of $P(\theta; \psi)$ at different time steps. In this context, we show that concept drift corresponds to changes in the distributional uncertainty of model parameters. In other words, real concept drift can be measured as a change in the average number of bits required to encode the parameters of the predictive model. By specifying an adequate parameter distribution, we can identify concept drift at the input level, which offers a significant advantage over existing approaches in terms of explainability. In fact, the proposed framework can be applied to almost any parameter distribution and online predictive model. For illustration, we apply *ERICS* to a Probit model. In experiments on both synthetic and real-world data sets, we show that the proposed framework can detect different types of concept drift, while having a lower average delay than state-of-the-art methods. Indeed, *ERICS* outperforms existing approaches with respect to the recall and precision of concept drift alerts.

In summary, we propose a generic and flexible framework that leverages the uncertainty patterns of model parameters for more effective concept drift detection in data streams. An open source version of *ERICS* is available at <https://github.com/haugjo/erics>.

II. ERICS: A CONCEPT DRIFT DETECTION FRAMEWORK

Real concept drift corresponds to a change of the conditional target distribution $P(Y|X)$ [1]. However, data streams are potentially infinite and so the true distribution $P(Y|X)$ remains unknown. Hence, we may use a predictive model f_{θ} to approximate $P(Y|X)$. Since we update the model parameters θ for every new observation, θ_t represents our most current information about the active concept at time step t . Consequently, we may identify concept drift by investigating changes in θ over time.

To this end, we adopt the general framework of [4] and treat the parameters θ as a random variable, i.e. $\theta \sim P(\theta; \psi)$. Analogously, we optimize the distribution parameters ψ at every

time step with respect to the log-likelihood. This optimization problem can be expressed in terms of the marginal likelihood $P(Y|X, \psi)$ [4]. Hence, the marginal likelihood relates to the optimal parameter distribution under the active concept. Accordingly, we may associate concept drift between two time steps t and u with a difference of the marginal likelihood for the distribution parameters ψ_t and ψ_u :

$$\begin{aligned} & P(Y|X; \psi_t) \neq P(Y|X; \psi_u) \\ \Leftrightarrow & |P(Y|X; \psi_t) - P(Y|X; \psi_u)| > 0 \\ \Leftrightarrow & \left| \int P(Y|X, \theta) [P(\theta; \psi_t) - P(\theta; \psi_u)] d\theta \right| > 0. \quad (1) \end{aligned}$$

From (1), we may obtain a general scheme for concept drift detection. To this end, we rephrase (1) in terms of the differential entropy and KL-divergence, which are common measures from information theory. The entropy of a random variable corresponds to the average degree of uncertainty of the possible outcomes. Besides, entropy is often described as the average number of bits required to encode a sample of the distribution. On the other hand, the KL-divergence measures the difference between two probability distributions. It is frequently applied in Bayesian inference models, where it describes the information gained by updating from a prior to a posterior distribution. We can derive the following proportionality:

$$\begin{aligned} & \int P(Y|X, \theta) [P(\theta; \psi_t) - P(\theta; \psi_u)] d\theta \\ \propto & \int P(\theta; \psi_t) d\theta - \int P(\theta; \psi_u) d\theta \\ \propto & \int P(\theta; \psi_t) \log P(\theta; \psi_t) d\theta - \int P(\theta; \psi_u) \log P(\theta; \psi_t) d\theta \\ = & H[P(\theta; \psi_u), P(\theta; \psi_t)] - h[P(\theta; \psi_t)] \\ = & h[P(\theta; \psi_u)] - h[P(\theta; \psi_t)] + D_{KL}[P(\theta; \psi_u) \| P(\theta; \psi_t)], \quad (2) \end{aligned}$$

where $h[P(\theta; \psi_t)]$ is the differential entropy of the parameter distribution at time step t . Note that we have rephrased the cross entropy $H[P(\theta; \psi_u), P(\theta; \psi_t)]$ by using the KL-divergence D_{KL} . We may now substitute (2) into (1) to derive a general scheme for concept drift detection:

$$\left| \underbrace{h[P(\theta; \psi_u)] - h[P(\theta; \psi_t)]}_{\Delta \text{Uncertainty}} + \underbrace{D_{KL}[P(\theta; \psi_u) \| P(\theta; \psi_t)]}_{\Delta \text{Distribution}} \right| > 0 \quad (3)$$

Intuitively, real concept drift thus corresponds to a change in the uncertainty of the optimal parameters and a divergence of the parameter distribution. On the other hand, stable concepts are characterized by a static parameter distribution and uncertainty.

Note that (3) has another interpretation in the context of Bayesian inference. As mentioned before, the KL-divergence $D_{KL}[P(\theta; \psi_u) \| P(\theta; \psi_t)]$ can be interpreted as the information gained from inferring the posterior $P(\theta; \psi_u)$ from a prior $P(\theta; \psi_t)$. According to (3), we thus find that every difference in parameter uncertainty (entropy) between time step t and

u , which can not be attributed to the inference of posterior parameters, may be traced back to a concept drift.

Finally, we show that the proposed concept drift detection scheme adheres to the Properties 1 and 2.

Proof that ERICS is model-aware (Property 1): By construction, we model the parameters θ through a distribution $P(\theta; \psi)$. According to (3), we write

$$\left| \int P(\theta; \psi_t) \log P(\theta; \psi_t) d\theta - \int P(\theta; \psi_u) \log P(\theta; \psi_t) d\theta \right|,$$

which is 0 iff $P(\theta; \psi_t) = P(\theta; \psi_u)$. Consequently, we find that Equation (3) evaluates to true, iff $P(\theta; \psi_t) \neq P(\theta; \psi_u)$. By definition, for any sensible statistical divergence measure \mathcal{D} , we know that $\mathcal{D}(P(\theta; \psi_t), P(\theta; \psi_u)) = 0 \Leftrightarrow P(\theta; \psi_t) = P(\theta; \psi_u)$. Equation (3) holds true, and thus $P(\theta; \psi_t) \neq P(\theta; \psi_u) \Leftrightarrow \mathcal{D}(P(\theta; \psi_t), P(\theta; \psi_u)) > 0$ ■

Proof that ERICS is explainable (Property 2): By construction, any parametric distribution $P(\theta; \psi)$ used in Equation (3) can be evaluated for each parameter individually, i.e. we have $P(\theta_k; \psi_k) \forall k$. ■

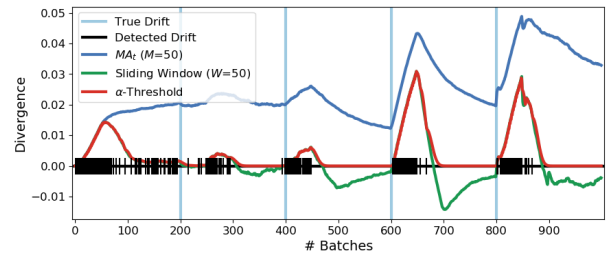
A. Continuous Concept Drift Detection

Based on the general scheme (3), we are able to identify concept drift between any two time steps t and u . In practice, we are mainly interested in concept drifts between successive time steps $t-1$ and t . However, if we were to study (3) for two time steps only, our concept drift detection model might become too sensitive to random variations of the predictive model. To be more robust, we examine the moving average of (3) instead. Specifically, we compute the moving average at time step t over M time steps as

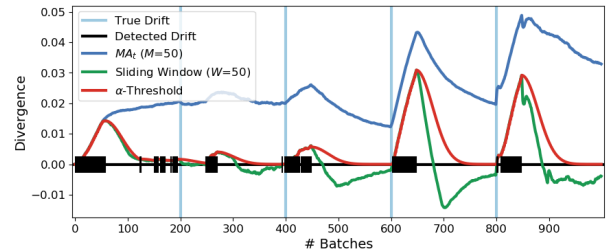
$$\text{MA}_t = \frac{1}{M} \sum_{i=t-M+1}^t \left(|h[P(\theta; \psi_i)] - h[P(\theta; \psi_{i-1})]| + D_{KL}[P(\theta; \psi_i) \| P(\theta; \psi_{i-1})] \right). \quad (4)$$

As before, the moving average contains our latest information on the model parameters and the active concept. We can adjust the sensitivity of our framework by selecting M appropriately. In general, the larger we select M , the more robust the framework becomes. However, a large M might also hide concept drifts of small magnitude or short duration.

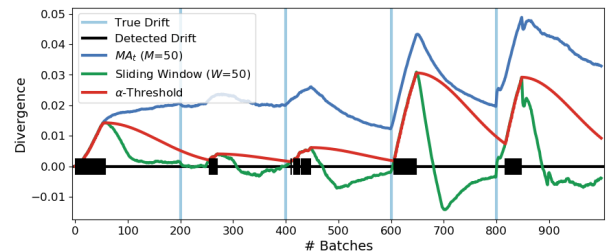
So far we have treated all changes of the parameter distribution as an indication of concept drift. Indeed, this is in line with the general definition of concept drift [1]. Still, we argue that only certain changes in the parameter distribution have practical relevance. For example, suppose that we use stochastic gradient descent (SGD) to optimize the model parameters at every time step. If we start from an arbitrary initialization, the distribution of optimal parameters usually changes significantly in early training iterations. However, given that the concept $P(Y|X)$ is stationary, SGD will almost surely converge to a local optimum. Consequently, we will ultimately minimize the entropy and KL-divergence of $P(\theta; \psi)$ in successive time steps. In other words, (4) will tend to decrease as long as we optimize the parameters ψ with respect



(a) $\beta = 0.01$



(b) $\beta = 0.001$



(c) $\beta = 0.0001$

Fig. 1. *Updating the α -Threshold.* The proposed framework uses a dynamic threshold α (red line) to detect concept drift. According to (4), we track a moving average of the divergence of the parameter distribution (dark blue line). If the total divergence in a sliding window (green line) exceeds the threshold, *ERICS* detects a concept drift (black vertical lines). By adjusting the hyperparameter β , we can control the iterative updates of α and thus regulate the sensitivity of *ERICS* after a drift is detected. Generally, the larger we choose β , the more sensitive *ERICS* becomes to changes of the parameter distribution. Here, we depict different β for the KDD data set [6]. We artificially generated four sudden concept drifts (blue vertical lines). In this example, small update steps (i.e. small β) are preferable to give the predictive model enough time to adapt to the new concept. Note that the early alerts correspond to the initial training phase of the predictive model. Hence, we would ignore them in practice.

to the active concept. However, if the decision boundary changes due to a real concept drift, SGD-updates will aim for a different optimum. This change of the objective will temporarily lead to more uncertainty in the model and thus increase the entropy of the parameter distribution.

We exploit this temporal pattern for concept drift detection. To this end, we measure the total change of (4) in a sliding window of size W :

$$\sum_{j=t-W+1}^t (\text{MA}_j - \text{MA}_{j-1}) > \alpha_t \Leftrightarrow \text{Drift at } t, \quad (5)$$

where $\alpha_t \geq 0$ is an adaptive threshold. As before, we may control the robustness of the concept drift detection with the

sliding window size W . Whenever we detect concept drift, i.e. (5) evaluates to true, we redefine α_t as

$$\alpha_t = \sum_{j=t-W+1}^t (\text{MA}_j - \text{MA}_{j-1}). \quad (6)$$

In this way, we temporarily tolerate all changes to the predictive model up to a magnitude of (6). We consider these changes to be the after-effects of the concept drift. We then update α_t in an iterative fashion. Let β be a user-defined hyperparameter in the interval $[0, 1]$. Each update depends on the current α -value, the β -hyperparameter and the time elapsed since the last concept drift alert, which we denote by Δ_{Drift} :

$$\alpha_t = \alpha_{t-1} - (\alpha_{t-1} * \beta * \Delta_{Drift}) \quad (7)$$

Note that α_t will asymptotically approach 0 over time, if there is no concept drift. In this way, we gradually reduce the tolerance of our framework after a drift is detected.

The choice of a suitable β usually depends on the application at hand. By way of illustration, we applied *ERICs* with different β to the KDD data set [6]. We used [7]’s method to induce sudden concept drift after every 20% of observations. For more information, see Section V. Figure 1 illustrates the components of *ERICs* for three different β -values. Notably, the larger we chose β , the more drifts we detected. Since we were dealing with a sudden concept drift in this particular example, we could be less sensitive and apply smaller update steps. For $\beta = 0.0001$, we achieved good first results in all our experiments. Therefore, this value can generally be used as a starting point for further optimization.

To conclude our general framework, we provide a pseudo code implementation in Figure 2.

B. Limitations and Advantages

The proposed framework does not access streaming observations directly, but uses the parameters of a predictive model instead. Accordingly, our approach is much more memory efficient than many related works. Yet, if the parameter distribution does not change in a drift period, concept drift may go unnoticed. In general, however, *ERICs* can detect all concept drifts that affect the predictive outcome.

One should also be aware that some predictive models are prone to adversarial attacks. Accordingly, *ERICs* can only be as robust as its underlying predictive model. This sensitivity to the predictive model is shared by most existing works. With *ERICs*, the possibility of misuse is drastically reduced, as we closely monitor the distribution of the model parameters at all times.

III. ILLUSTRATING ERICS

ERICs is model-agnostic. This means that the framework can be applied to different predictive models f_θ and parameter distributions $P(\theta; \psi)$. In this way, we enable maximum flexibility with regard to possible streaming applications. By way of illustration, we adopt a Probit model with independent normally distributed parameters. This setup has achieved state-of-the-art results in online feature selection [4]. Besides,

Require: $[\psi_t, \dots, \psi_{t-M}]; [\text{MA}_{t-1}, \dots, \text{MA}_{t-W}]; \alpha_{t-1}; \Delta_{Drift}$

```

 $\alpha_t \leftarrow Eq. (7)$ 
 $\Delta_{Drift} \leftarrow \Delta_{Drift} + 1$ 
 $\text{MA}_t \leftarrow Eq. (4)$ 
 $sumWindow \leftarrow \sum_{j=t-W+1}^t (\text{MA}_j - \text{MA}_{j-1})$ 

```

if $sumWindow > \alpha_t$ **then**

```

 $\alpha_t \leftarrow sumWindow$ 

```

```

 $\Delta_{Drift} \leftarrow 1$ 

```

end if

return $\alpha_t; \text{MA}_t; \Delta_{Drift}$

Fig. 2. Pseudo Code. Concept drift detection with *ERICs* at time step t .

it offers dramatic computational advantages due to its low complexity. In line with [4], we optimize ψ at every time step with respect to the log-likelihood for the Probit model.

The assumption of independent model parameters may appear restrictive, but in practice it often leads to good results, e.g. in the case of local feature attributions [8], [9] or feature selection [4], [10]. In fact, the independence assumption allows us to identify the parameters affected by concept drift and thus to comply with Property 2. Since the Probit model comprises one parameter per input feature, we can readily associate concept drift with individual input variables.

Accordingly, let $P(\theta; \psi_t) = \mathcal{N}(\psi_t = (\mu_t, \Sigma_t))$, where $\mu_t = (\mu_{tk})_{k=1}^K$ is a vector of mean values and Σ_t is the diagonal covariance matrix, where the diagonal entries correspond to the vector $\sigma_t^2 = (\sigma_{tk}^2)_{k=1}^K$. The differential entropy of $P(\theta; \psi_t)$ is

$$h[P(\theta; \psi_t)] = \frac{1}{2} \left(K + K \ln(2\pi) + \ln \prod_{k=1}^K \sigma_{tk}^2 \right).$$

The KL-divergence between $P(\theta; \psi_t)$ and $P(\theta; \psi_{t-1})$ is

$$D_{KL}[P(\theta; \psi_t) \| P(\theta; \psi_{t-1})] = \frac{1}{2} \left(\sum_{k=1}^K \frac{\sigma_{tk}^2 + (\mu_{t-1,k} - \mu_{tk})^2}{\sigma_{t-1,k}^2} - K + \ln \frac{\prod_{k=1}^K \sigma_{t-1,k}^2}{\prod_{k=1}^K \sigma_{tk}^2} \right).$$

According to (4), we then write the moving average as

$$\text{MA}_t = \frac{1}{2M} \sum_{i=t-M+1}^t \left| \sum_{k=1}^K \frac{\sigma_{ik}^2 + (\mu_{i-1,k} - \mu_{ik})^2}{\sigma_{i-1,k}^2} - K \right|. \quad (8)$$

Note that (8) scales linearly with the number of parameters K , i.e. it has $\mathcal{O}(K)$ time complexity.

In order to identify concept drift at individual parameters (which is equivalent to examining individual features, since we use a Probit model), we can investigate the moving average of a specific parameter θ_k :

$$\text{MA}_{tk} = \frac{1}{2M} \sum_{i=t-M+1}^t \left| \frac{\sigma_{ik}^2 + (\mu_{i-1,k} - \mu_{ik})^2}{\sigma_{i-1,k}^2} - 1 \right| \quad (9)$$

In this case, we maintain a different threshold α_k per parameter. Note that (9) has a constant time complexity.

IV. RELATED WORK

In this section, we briefly introduce some of the most prominent and recent contributions to concept drift detection.

DDM monitors changes in the classification error of a predictive model [11]. Whenever the observed error changes significantly, DDM issues a warning or an alert. We find various modifications of this general scheme, including [12] and [13]. Another well-known method for concept drift adaptation is ADWIN [14]. Here, the authors maintain a sliding window, whose size changes dynamically according to the current rate of distributional change. [15] also employ a sliding window approach and provide a feasible implementation of Fisher’s Exact test, which they use for concept drift detection. Similar to our framework, [16] use a sliding window and the entropy to detect concept drift. However, they examine entropy with regard to the predictive result and disregard the model parameters. FHDDM applies a sliding window to classification results and tracks significant differences between the current probability of correct predictions and the previously observed maximal probability [17]. To this end, FHDDM employs a threshold that is based on the Hoeffding bound. In a later approach, the same authors instead use McDiarmid’s inequality to detect concept drift [18]. EWMA is a method that monitors an increase in the probability that observations are misclassified [19]. The authors use an exponentially weighted moving average, which places greater weight on the most recent instances in order to detect changes. [20] also focus on the predictive outcome. Specifically, they investigate the distribution of the loss function via resampling. Likewise, the LFR method uses certain test statistics to detect concept drift by identifying changes through statistical hypothesis testing [21]. Finally, [22] compare the labels of close data points in successive batches to detect concept drift.

In addition, we find various approaches that examine ensembles of online learners to deal with concept drift. For example, [23] compare two models; one that is trained with all streaming observations and another that is trained only with the latest observations. Likewise, [24] analyze the density of the posterior distributions of an incremental and a static estimator.

More information about the progress in concept drift detection can be found in [1]–[3], [25].

Conceptually, our work differs substantially from the remaining literature. Instead of directly examining the streaming observations or the predictive outcome, *ERICs* monitors changes in the parameters of a predictive model.

V. EXPERIMENTS

We evaluated *ERICs* in multiple experiments. All experiments were conducted on an i5-8250U CPU with 8 Gb of RAM, running 64-bit Windows 10 and Python 3.7.3. We compared our framework to the popular concept drift detection methods ADWIN [14], DDM [11], EWMA [19], FHDDM [17], MDDM [18] and RDDM [13]. We used the predefined implementations of these models as provided by the Tornado framework [26]. Besides, we applied the default

TABLE I
SYNTHETIC AND REAL WORLD DATA SETS

Name	#Samples	#Features	Data Types
SEA (synth.)	100,000	3	cont.
Agrawal (synth.)	100,000	9	cont.
Hyperplane (synth.)	100,000	20	cont.
Mixed (synth.)	100,000	9	cont.
Spambase	4,599	57	cont.
Adult	48,840	54	cont., cat.
HAR (binary)	7,450	562	cont.
KDD (sample)	100,000	41	cont., cat.
Dota	102,944	116	cat.
MNIST (binary)	10,398	784	cont.

TABLE II
HYPERPARAMETERS OF *ERICs* PER DATA SET

Data Set	M	W	β	Epochs	LR μ	LR σ
SEA	75	50	0.0001	10	0.01	0.01
Agrawal	100	50	0.001	10	0.01	0.01
Hyperplane	100	50	0.0001	10	0.01	0.01
Mixed	100	50	0.0001	10	0.1	0.01
Spambase	35	25	0.001	10	0.1	0.01
Adult	50	50	0.001	10	0.1	0.01
HAR	25	50	0.001	10	0.1	0.01
KDD	50	50	0.0001	10	0.01	0.01
Dota	75	50	0.0001	10	0.01	0.01
MNIST	25	20	0.001	50	0.1	0.01

set of parameters throughout all experiments. Note that all related models require classifications of a predictive model. To this end, we trained a Very Fast Decision Tree (VFDT) [27] in an interleaved test-then-train evaluation. The VFDT is a state-of-the-art online learner, which uses the Hoeffding bound to incrementally construct a decision tree for streaming data. We used the VFDT implementation of scikit-multiflow [28] in our experiments. Note that we consider a simple binary classification scenario in all our experiments, since it should be handled well by all models.

We optimized the hyperparameters of *ERICs* in a grid search. The search space was either chosen empirically or according to [4]. Table II lists all hyperparameters per data set. The hyperparameters “Epochs”, “LR (learning rate) μ ” and “LR σ ” control the training of the Probit model, which we adopted from [4].

A. Data Sets

In order to evaluate the timeliness and precision of a concept drift detection model, we require ground truth. Consequently, we generated multiple synthetic data sets using the scikit-multiflow package [28]. Detailed information about each generator can be obtained from the corresponding documentation. We exhibit the properties of all data sets in Table I. Note that we simulated multiple types of concept drift. Specifically, we produced sudden concept drifts with the SEA generator. To this end, we specified a drift duration (*width* parameter) of 1. We alternated between the classification functions 0-3 to produce the different concepts. With the Agrawal generator, we simulated gradual drift of different duration. Again, we

alternated between the classification functions 0-3 to shift the data distribution. With the rotating Hyperplane generator, we simulated an incremental drift over the full length of the data set. We generated 20 features with the Hyperplane generator, out of which 10 features were subject to concept drift by a magnitude of 0.5. Finally, we produced a Mixed drift using the Agrawal generator. The Mixed data contains both sudden and gradual drift, which we obtained by alternating the classification functions 0-4. All synthetic data sets contain 10% noisy data. We obtained 100,000 observations from each data stream generator.

In addition, we evaluated the proposed framework on real world data. However, since real world data usually does not provide any ground truth information, we had to artificially induce concept drift. For this reason, we applied the methodology of [7] to induce sudden concept drift in five well-known data sets from the online learning literature. First, we randomly shuffled the data to remove any natural (unknown) concept drifts. Next, we ranked all features according to their information gain. We then selected the top 50% of the ranked features and randomly permuted their values. In this way, we generated sudden drifts after every 20% of the observations. Specifically, we introduced concept drift to the real-world data sets Spambase, Adult, Human Activity Recognition (HAR), KDD 1999 and Dota2, which we took from the UCI Machine Learning repository [6]. Note that we drew a random sample of 100,000 observations from the KDD 1999 data to allow for feasible computations.

Besides, we used the MNIST data set to evaluate partial concept drift detection at the input level. We selected all observations that are either labelled 3 or 8, since these numbers are difficult to distinguish. In the first half of the observations, we treated 3 as the true class. In the second half of the observations, we switched the true class to 8. In this way, we simulated a sudden concept drift of all input features.

For all real world data sets, we normalized the continuous features to the range $[0, 1]$ and one-hot-encoded the categorical features. In the Adult data set, we imputed all NaN-values by a new category *unknown*. Moreover, we altered the labels of the HAR data set to simulate binary classification between the class *moving* (original labels *walking*, *walking_downstairs* and *walking_upstairs*) and *non-moving* (original labels *sitting*, *laying* and *standing*). We trained the online predictive models (Probit and VFDT) in batches of the following size: For Spambase and HAR we chose a batch size of 10. Adult was processed in batches of size 50. For all remaining data sets, we trained on batches of 100 observations.

B. Delay, Recall and Precision

In our first experiment, we applied the concept drift detection models to all synthetic and real-world data sets. Figure 3 exhibits the drift alerts of every model. The blue vertical lines and shaded areas indicate periods of concept drift. Each black vertical line corresponds to one drift alert. Most models identify concept drift in early iterations. This is due to the initial training phase of the predictive model and therefore

TABLE III
AVERAGE DELAY IN NUMBER OF BATCHES

Datasets	Drift Detection Models						
	ERICs	ADWIN	DDM	EWMA	FHDDM	MDDM	RDDM
SEA	52.75	34.55	16.45	0.26	178.61	178.58	7.42
Agrawal	71.33	16.80	0.00	0.31	0.023	0.20	137.22
Hyperplane	2.00	42.27	2.23	2.36	11.24	11.22	2.42
Mixed	34.00	27.95	0.34	11.55	75.26	75.25	227.98
Spambase	7.35	79.0	60.23	29.96	71.63	71.58	117.45
Adult	2.61	63.15	488.39	172.57	488.39	488.39	488.39
HAR	13.05	372.45	372.45	1.55	372.45	372.45	77.10
KDD	22.01	50.35	0.00	0.55	100.25	100.21	13.21
Dota	44.04	25.32	514.72	43.40	3.56	3.54	116.46
Mean	27.68	79.09	161.65	29.17	144.60	144.60	131.96
Rank	1	3	7	2	5	5	4

has no practical relevance. For the upcoming evaluations, we have therefore ignored all drift alerts in the first 80 batches.

By Figure 3, the proposed framework *ERICs* performs well in all data sets. Given the low complexity of the underlying Probit model, some concept drifts do not infer a change of the parameter distribution immediately. This can be seen in small delays, such as for the Agrawal data, for example. Still, *ERICs* achieves the smallest average delay of all concept drift detection models, which is shown in Table III.

Strikingly, *ERICs* generally seems to produce fewer false alarms than related models. We find support for this intuition by examining the average recall (Figure 4) and precision (Figure 5) over all data sets. Similar to [29], we evaluated the detected drifts for different detection ranges. The detection range corresponds to the number of batches after a known drift, during which we consider an alert as a true positive. Whenever there is no drift alert in the detection range, we count this as a false negative. Besides, all drift alerts outside of the detection range are false positives. We used these scores to compute the recall and precision values. Again, we find that *ERICs* tends to struggle in the early stages, right after a drift happens. As mentioned before, we attribute this to the slowly updating Probit model that we used for illustration. The VFDT, which is used by all related models, is much more complex and can thus adapt to changes faster. Additionally, we must treat some recall scores with care. For example, in four data sets, the DDM model detects drift in almost every time step. Hence, it achieves perfect recall, although the drift alerts are not reliable at all. Still, *ERICs* ultimately outperforms all related models in terms of both recall and precision. The superiority of our framework is even more apparent, if we look at the harmonic mean of precision and recall, which is the F1 score that we show in Figure 6.

C. Detecting Drift at the Input Level

As mentioned before, by using a Probit model and treating parameters as independently Gaussian distributed, we are able to associate concept drift with specific input features. By means of illustration, we apply *ERICs* to a sample of the MNIST data set, which we induced with concept drift. In Figure 7, we exhibit the mean of all observations corresponding to the true class before and after the concept drift (left

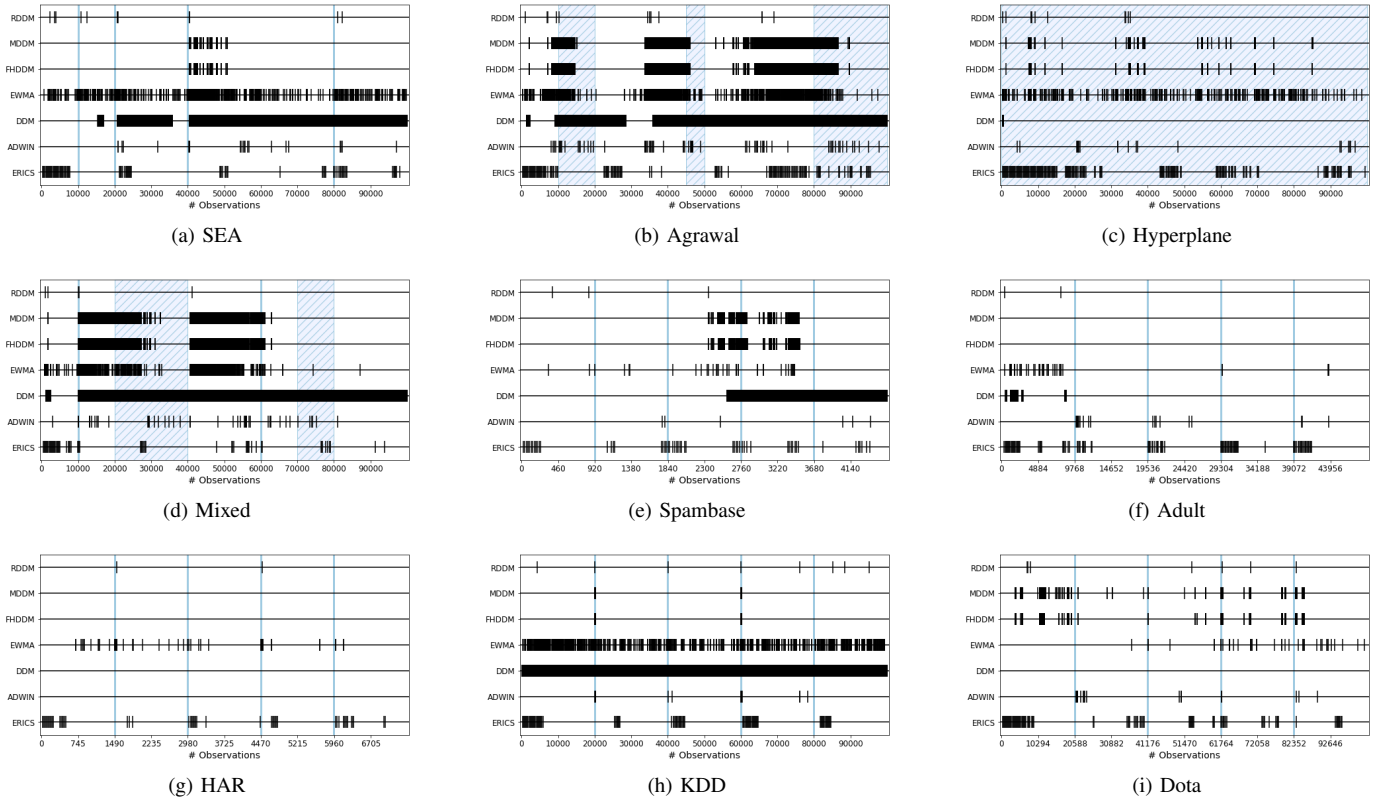


Fig. 3. *Drifts Alerts*. For all data sets, we illustrate the drift alerts obtained from each concept drift detection model. The blue vertical lines and shaded areas correspond to known concept drifts. Each black marker stands for one drift alert. Early drift alerts can be attributed to the initial training of the predictive model and were therefore ignored. Notably, *ERICS* (ours) seems to detect most concept drifts, while triggering considerably fewer false alarms than most related models. We find support for this intuition in the remaining figures.

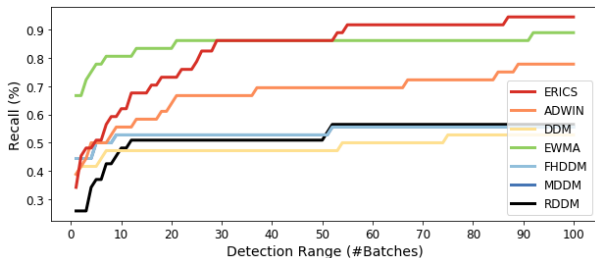


Fig. 4. *Recall*. We show the average recall over all data sets for different detection ranges. *ERICS* (ours) ultimately detects more than 90% of the known concept drifts. The apparent disadvantage of *ERICS* in early batches can be attributed to the slower update speed of the Probit model as compared to the VFDT [27], which was used by the remaining concept drift detection methods. Besides, the recall scores should be considered with care, since some methods tend to detect drift at almost every time step and are thus not reliable.

subplots). We also show the absolute difference between those mean values. In the outer most subplot on the right, we illustrate the drift alerts per input feature in the first 15 batches after the concept drift. The color intensity corresponds to the number of drift alerts (where many alerts correspond to darker patterns). Strikingly, the frequency of drift alerts closely maps the absolute difference between the two concepts. This shows that *ERICS* is generally able to identify the input features that

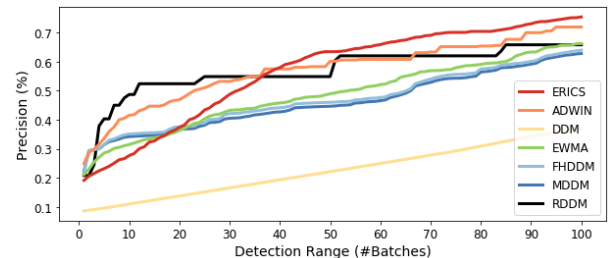


Fig. 5. *Precision*. We show the average precision over all data sets for different detection ranges. *ERICS* (ours) tends to identify concept drift later than some related models. Therefore, we count fewer true positives in small detection ranges, which leads to lower precision. However, for larger detection ranges, our framework is more precise than any other model in the evaluation.

are most affected by concept drift. We expect this pattern to become even clearer, when using more complex base models.

VI. CONCLUSION

In this work, we proposed a novel and generic framework for the detection of concept drift in streaming applications. Our framework monitors changes in the parameters of a predictive model to effectively identify distributional changes of the input. We exploit common measures from information theory, by showing that real concept drift corresponds to

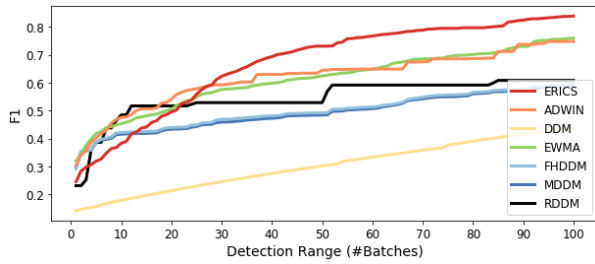


Fig. 6. *F1*: We illustrate the *F1* measure, which is the harmonic mean of the precision and recall shown in earlier plots. Here, the advantage of *ERICCS* (ours) is most apparent, since it significantly outperforms all related methods for a detection range greater than 30 batches.

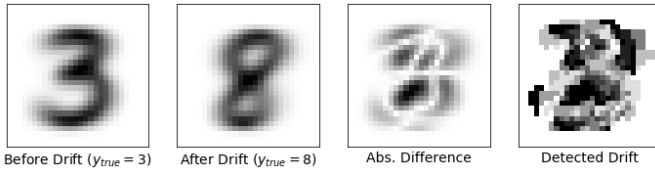


Fig. 7. *Partial Drift Detection*. By choosing an appropriate base model and parameter distribution, *ERICCS* can attribute concept drift to individual input features. We selected all observations of MNIST with the label 3 or 8 and induced concept drift by changing the true class after half of the observations. In the left subplots, we exhibit the mean of the true class before and after the concept drift. The third subplot depicts the absolute difference of these mean values. In the right subplot, we show the alerts of *ERICCS* in the first 15 batches after the concept drift. The color intensity corresponds to the frequency of drift alerts per input feature. Strikingly, the drift alerts seem to map the absolute difference between both concepts. This suggests, that *ERICCS* does indeed identify concept drift for those input features that are most affected by a distributional change.

changes of the uncertainty regarding the optimal parameters. Given an appropriate parameter distribution, the proposed framework can also attribute drift to specific input features. In experiments, we highlighted the advantages of our approach over multiple existing methods, using both synthetic and real-world data. Strikingly, *ERICCS* detects concept drift with less delay on average, while outperforming existing models in terms of both recall and precision.

REFERENCES

- [1] G. I. Webb, R. Hyde, H. Cao, H. L. Nguyen, and F. Petitjean, "Characterizing concept drift," *Data Mining and Knowledge Discovery*, vol. 30, no. 4, pp. 964–994, 2016.
- [2] J. Gama, I. Žliobaitė, A. Bifet, M. Pechenizkiy, and A. Bouchachia, "A survey on concept drift adaptation," *ACM computing surveys (CSUR)*, vol. 46, no. 4, p. 44, 2014.
- [3] I. Žliobaitė, "Learning under concept drift: an overview," *arXiv preprint arXiv:1010.4784*, 2010.
- [4] J. Haug, M. Pawelczyk, K. Broelemann, and G. Kasneci, "Leveraging model inherent variable importance for stable online feature selection," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 1478–1502.
- [5] P. Zhao, L.-W. Cai, and Z.-H. Zhou, "Handling concept drift via model reuse," *Machine Learning*, vol. 109, no. 3, pp. 533–568, 2020.
- [6] D. Dua and C. Graff, "UCI machine learning repository," 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [7] T. S. Sethi and M. Kantardzic, "On the reliable detection of concept drift from streaming unlabeled data," *Expert Systems with Applications*, vol. 82, pp. 77–99, 2017.

- [8] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in neural information processing systems*, 2017, pp. 4765–4774.
- [9] G. Kasneci and T. Gottron, "Licon: A linear weighting scheme for the contribution of input variables in deep artificial neural networks," in *Proceedings of the 25th ACM International Conference on Information and Knowledge Management*, 2016, pp. 45–54.
- [10] V. Borisov, J. Haug, and G. Kasneci, "Cancelout: A layer for feature selection in deep neural networks," in *International Conference on Artificial Neural Networks*. Springer, 2019, pp. 72–83.
- [11] J. Gama, P. Medas, G. Castillo, and P. Rodrigues, "Learning with drift detection," in *Brazilian symposium on artificial intelligence*. Springer, 2004, pp. 286–295.
- [12] M. Baena-García, J. del Campo-Ávila, R. Fidalgo, A. Bifet, R. Gavalda, and R. Morales-Bueno, "Early drift detection method," in *Fourth international workshop on knowledge discovery from data streams*, vol. 6, 2006, pp. 77–86.
- [13] R. S. Barros, D. R. Cabral, P. M. Gonçalves Jr, and S. G. Santos, "Rddm: Reactive drift detection method," *Expert Systems with Applications*, vol. 90, pp. 344–355, 2017.
- [14] A. Bifet and R. Gavalda, "Learning from time-changing data with adaptive windowing," in *Proceedings of the 2007 SIAM international conference on data mining*. SIAM, 2007, pp. 443–448.
- [15] D. R. de Lima Cabral and R. S. M. de Barros, "Concept drift detection based on fisher's exact test," *Information Sciences*, vol. 442, pp. 220–234, 2018.
- [16] L. Du, Q. Song, and X. Jia, "Detecting concept drift: an information entropy based method using an adaptive sliding window," *Intelligent Data Analysis*, vol. 18, no. 3, pp. 337–364, 2014.
- [17] A. Pesaraghader and H. L. Viktor, "Fast hoeffding drift detection method for evolving data streams," in *Joint European conference on machine learning and knowledge discovery in databases*. Springer, 2016, pp. 96–111.
- [18] A. Pesaraghader, H. L. Viktor, and E. Paquet, "Mcdiarmid drift detection methods for evolving data streams," in *2018 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2018, pp. 1–9.
- [19] G. J. Ross, N. M. Adams, D. K. Tasoulis, and D. J. Hand, "Exponentially weighted moving average charts for detecting concept drift," *Pattern recognition letters*, vol. 33, no. 2, pp. 191–198, 2012.
- [20] M. Harel, S. Mannor, R. El-Yaniv, and K. Crammer, "Concept drift detection through resampling," in *International Conference on Machine Learning*, 2014, pp. 1009–1017.
- [21] H. Wang and Z. Abraham, "Concept drift detection for streaming data," in *2015 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2015, pp. 1–9.
- [22] P. Sobhani and H. Beigy, "New drift detection method for data streams," in *International conference on adaptive and intelligent systems*. Springer, 2011, pp. 88–97.
- [23] S. H. Bach and M. A. Maloof, "Paired learners for concept drift," in *2008 Eighth IEEE International Conference on Data Mining*. IEEE, 2008, pp. 23–32.
- [24] C. H. Tan, V. Lee, and M. Salehi, "Online semi-supervised concept drift detection with density estimation," *arXiv preprint arXiv:1909.11251*, 2019.
- [25] P. M. Gonçalves Jr, S. G. de Carvalho Santos, R. S. Barros, and D. C. Vieira, "A comparative study on concept drift detectors," *Expert Systems with Applications*, vol. 41, no. 18, pp. 8144–8156, 2014.
- [26] A. Pesaraghader, H. Viktor, and E. Paquet, "Reservoir of diverse adaptive learners and stacking fast hoeffding drift detection methods for evolving data streams," *Machine Learning*, vol. 107, no. 11, pp. 1711–1743, 2018.
- [27] G. Hulten, L. Spencer, and P. Domingos, "Mining time-changing data streams," in *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, 2001, pp. 97–106.
- [28] J. Montiel, J. Read, A. Bifet, and T. Abdesslem, "Scikit-multiflow: A multi-output streaming framework," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 2915–2914, 2018.
- [29] S. Yu and Z. Abraham, "Concept drift detection with hierarchical hypothesis testing," in *Proceedings of the 2017 SIAM International Conference on Data Mining*. SIAM, 2017, pp. 768–776.