KDD 2024 – DELTA Workshop

On the Impact of Industrial Delays when Mitigating Distribution Drifts: an Empirical Study on Real-world Financial Systems

Thibault Simonetto¹, Maxime Cordy¹, Salah Ghamizi², Yves Le Traon¹,

Clément Lefebvre³, Andrey Boystov³, and Anne Goujon³

- ¹ University of Luxembourg, Luxembourg
- ² Luxembourg Institute of Science and Technology, Luxembourg
- ³ BGL BNP Paribas, Luxembourg





Context



Transaction system to detect and quarantine suspicious transaction

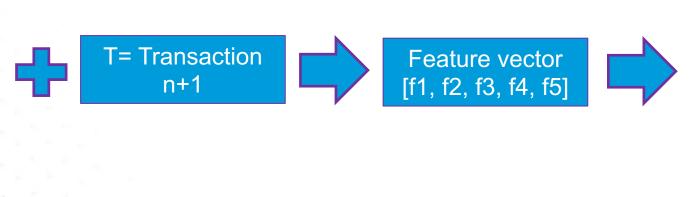
Transaction 1

Transaction 2

Transaction ...

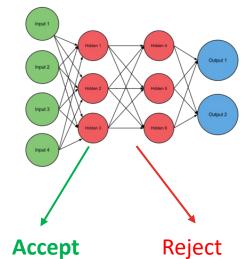
Transaction n

Client history



Incoming transaction

Feature vector



ML Model

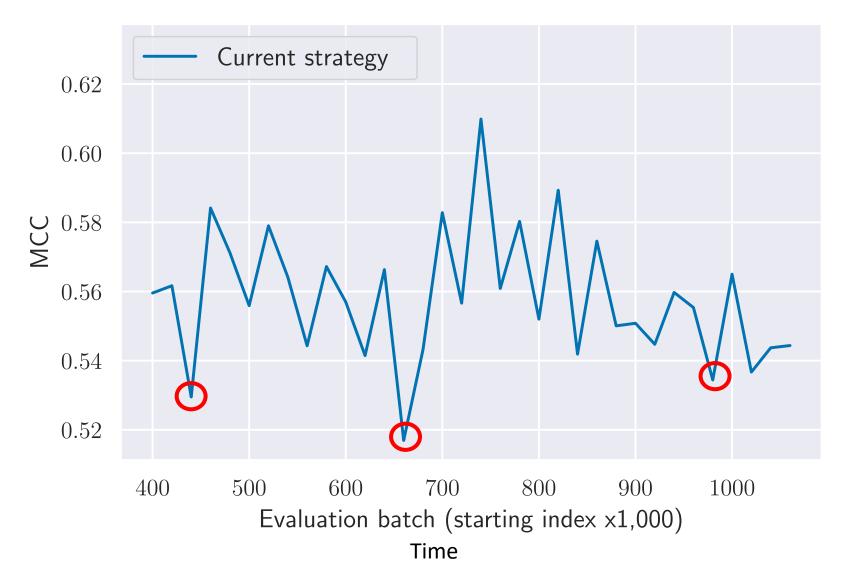
Correctly classifies up to 80% of the transactions





Problem: Performance drift







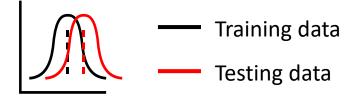
Drift detectors

Distribution
Sample stream

TTTTTT

Drift detector

Drift?



S S S S S S

Statistical test (p-value)

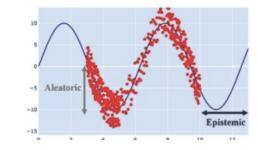
Distance metric (threshold)

Data-based detection



Requires the label

Error-based detection



Based on model behaviour for Given samples without labels

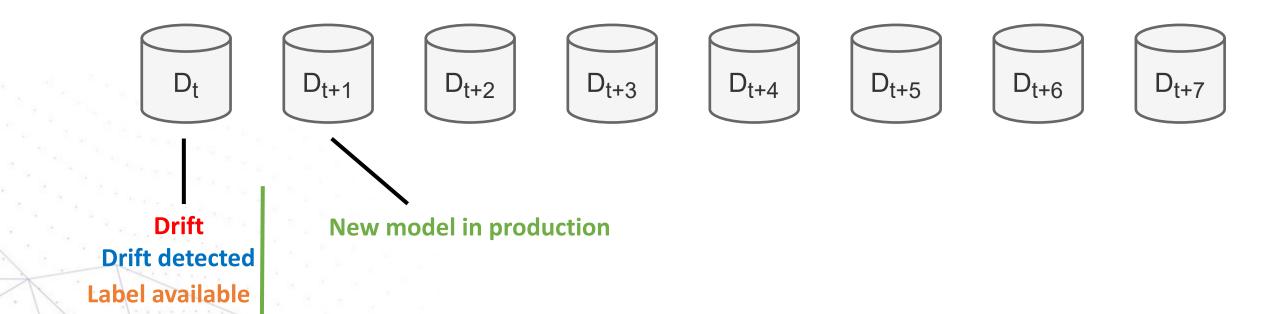
Predictive-based detection







Ideal scenario



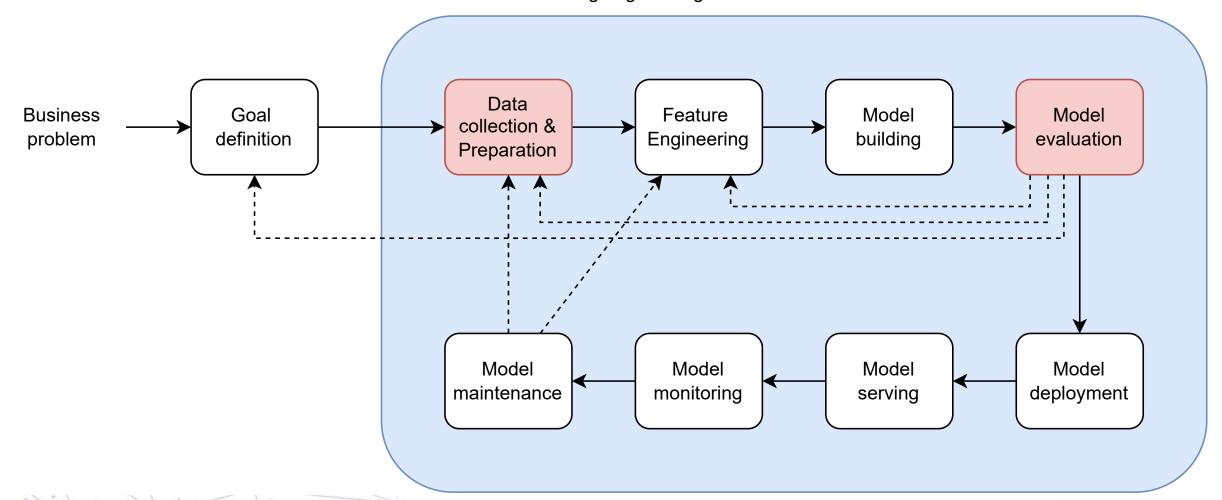
New model trained on





ML model lifecycle

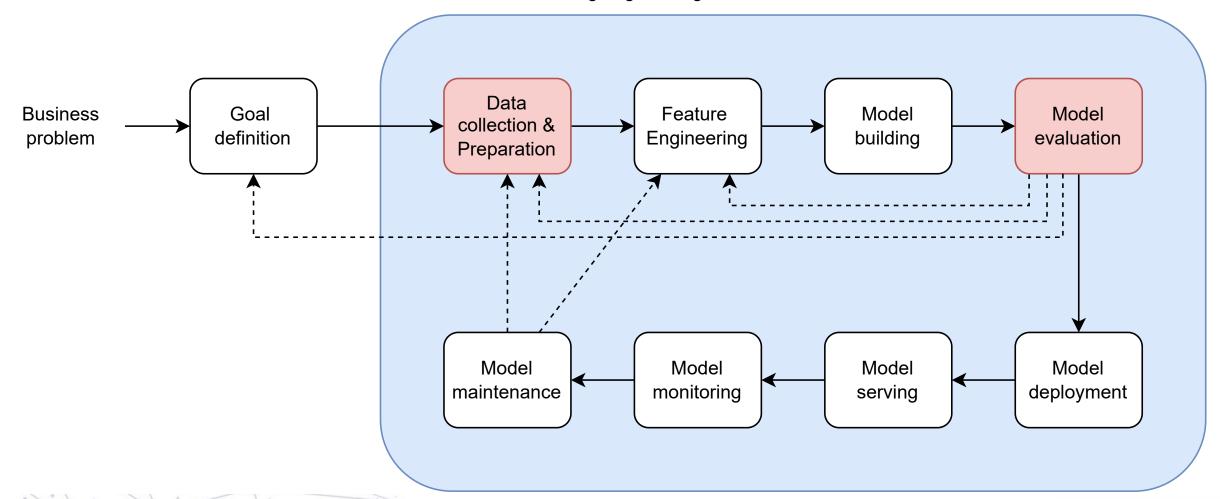
Machine learning engineering





ML model lifecycle

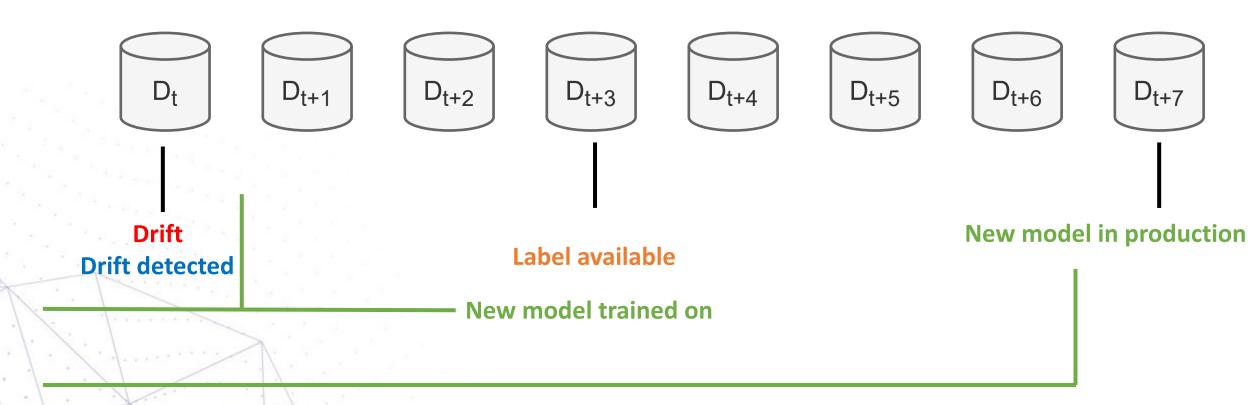
Machine learning engineering







Real-world Systems



Outdated model used to predict





Real-world Systems

 δ_d Labelling delay δ_l Deployment delay







 δ_d Labelling delay δ_l Deployment delay

 $\mathbf{h}_{t_j}: \mathcal{X} o \mathcal{Y}$ Classification model trained a time \mathbf{t}_j Using data $\{\mathbf{x}_i, y_i\}$ such that $t_i + \delta_l \leq t_j$

 \mathbf{h}_{t_j} can only make prediction on $\{\mathbf{x}_k\}$ such that $t_j+\delta_d\leq t_k$



 δ_d Labelling delay δ_l Deployment delay

 $\mathbf{h}_{t_j}: \mathcal{X} o \mathcal{Y}$ Classification model trained a time \mathbf{t}_j Using data $\{\mathbf{x}_i, y_i\}$ such that $t_i + \delta_l \leq t_j$

 \mathbf{h}_{t_j} can only make prediction on $\{\mathbf{x}_k\}$ $\mathrm{such}\ \mathrm{that}\ t_j + \delta_d \leq t_k$

 $\mathrm{sched} = \{ \mathrm{t}_1 \dots t_n \}$ Retraining schedule determining the sequence of models $H = \{ \mathrm{h}_{t_1} \dots h_{t_j} \dots h_{t_n} \}$



 δ_d Labelling delay δ_l Deployment delay

 $\mathrm{h}_{t_j}:\mathcal{X} o\mathcal{Y}$ Classification model trained a time t_j Using data $\{\mathrm{x}_i,y_i\}$ such that $t_i+\delta_l\leq t_j$

 \mathbf{h}_{t_j} can only make prediction on $\{\mathbf{x}_k\}$ such that $t_j+\delta_d\leq t_k$

 $\mathrm{sched} = \{ \mathrm{t}_1 \dots t_n \}$ Retraining schedule determining the sequence of models $H = \{ \mathrm{h}_{t_1} \dots h_{t_j} \dots h_{t_n} \}$

[X $_i$ is classified as $\hat{y}_i = h_{t_*}(x_i)$ where $t_* = \max\{t_k \in sched \text{ s.t. } t_k + \delta_{prod} \leq t_i\}$



 δ_d Labelling delay δ_l Deployment delay

 $\mathbf{h}_{t_j}: \mathcal{X} o \mathcal{Y}$ Classification model trained a time t_j Using data $\{\mathbf{x}_i, y_i\}$ such that $t_i + \delta_l \leq t_j$

 \mathbf{h}_{t_j} can only make prediction on $\{\mathbf{x}_k\}$ such that $t_j+\delta_d\leq t_k$

 $\mathrm{sched} = \{ \mathrm{t}_1 \dots \mathrm{t}_n \}$ Retraining schedule determining the sequence of models $H = \{ \mathrm{h}_{t_1} \dots \mathrm{h}_{t_j} \dots \mathrm{h}_{t_n} \}$

 $\{\mathbf{X}_i \;\;\; ext{is classified as} \;\;\;\; \hat{y}_i = h_{t_*}(x_i) \; ext{where} \; t_* = \max\{t_k \in sched \; ext{s.t.} \; t_k + \delta_{prod} \leq t_i\}$

We evaluate the schedule $sched = \{t_1 \dots t_n\}$ We evaluate the schedule of effectiveness and cost

$$s = score(Y, \hat{Y}) = MCC(Y, \hat{Y})$$
$$c = cost(sched) = |H|$$



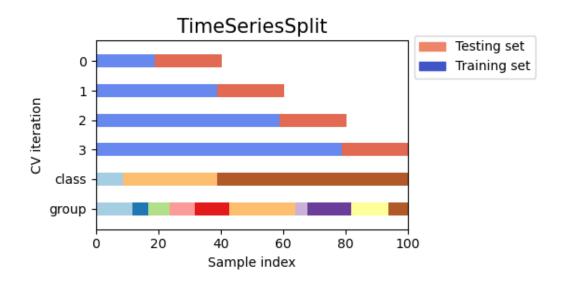


Evaluating drift detectors

Train Test

Train Drift Detectors Hyper-Pameters

25 Search iterations on with 5-fold timeseries validation





Evaluating drift detectors

Split 0 to 5

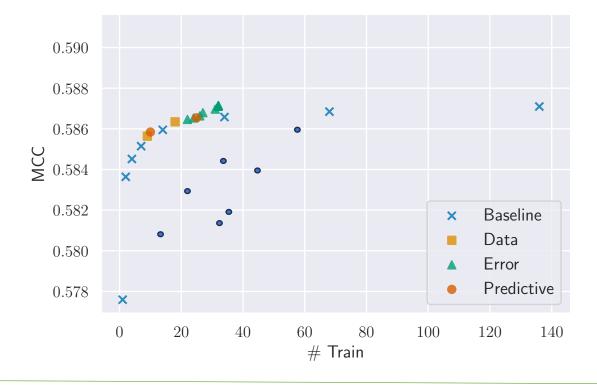
Train Test Train model on the train **Evaluate a single drift detector and parameters with realistic delays** dataset D_{t+1} D_{t+2} D_{t+3} D_{t+5} D_t D_{t+4} D_{t+6} D_{t+7} Drift New model in production Label available **Drift detected** New model trained on Outdated model used to predict



Evaluating drift detectors

Train Test

Evaluate the solution using best training parameters for each detector





Train Test 2 years of data 4 years of data Type Detector No detection Baseline Periodic Statistical test Data-based Divergence detector PC -CD DWIN (CE) DWIN (PE) DDM**EDDM** Error-based HDDM-HDDM-W detector KSWIN (CE) KSWIN (PE) Page-Hinkley (CE) Page-Hinkley (PE)

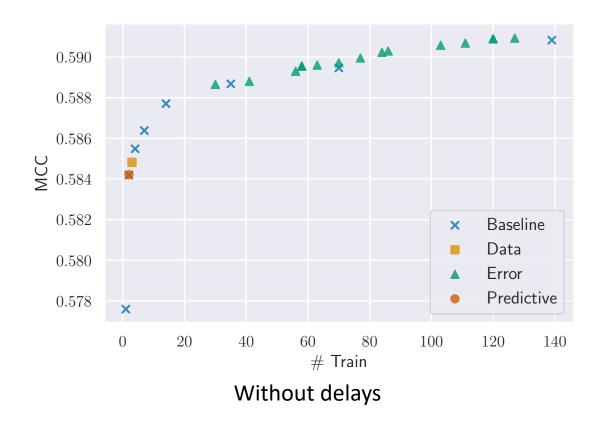
Predictive-based

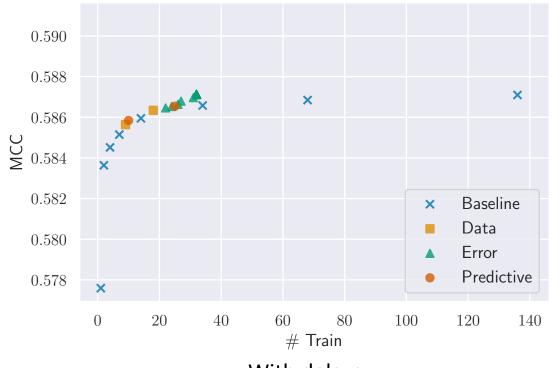
detector

Uncertainty ries DWIN





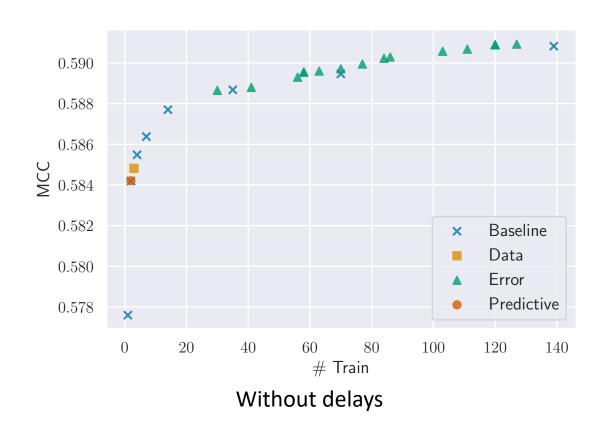


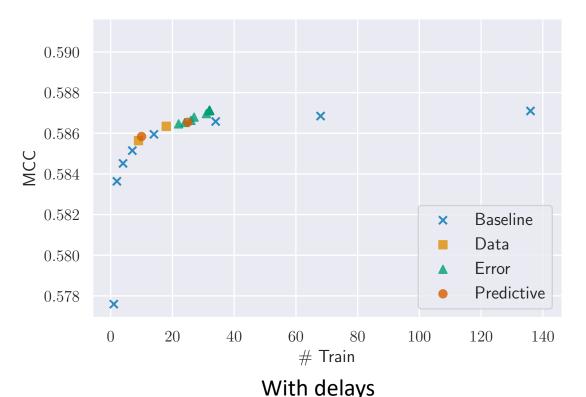


With delays Label: 10days, Validation/Deployment: 4 weeks







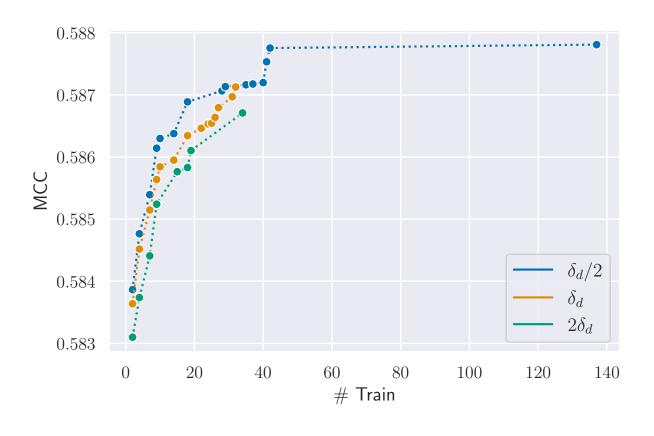


Label: 10days, Validation/Deployment: 4 weeks

Not considering delay overestimates the effectiveness/efficiency trade-off of retraining schedules

Change in delay disrupts the ranking of drift detectors



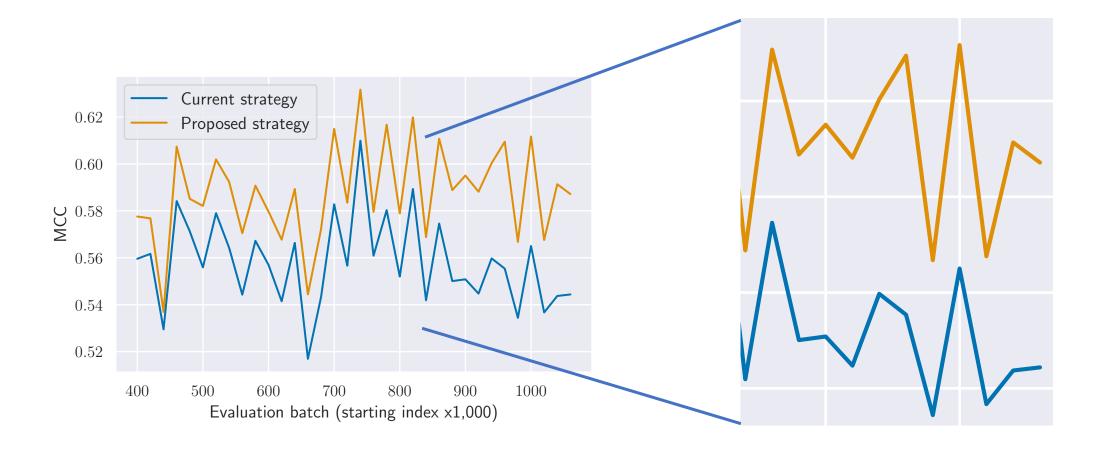


Change in delay has an inverse effect on efficiency/effectiveness.









Our method help to mitigate performance drift





Conclusion

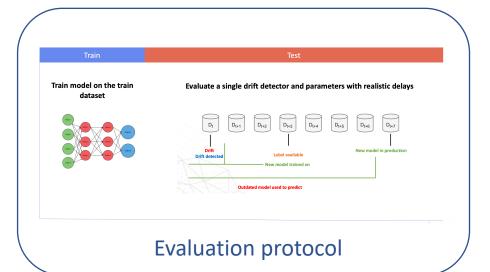
 $\begin{array}{ll} \delta_d \quad \text{Labelling delay} \quad \delta_l \quad \text{Deployment delay} \\ \mathbf{h}_{t_j}: \mathcal{X} \to \mathcal{Y} \quad \text{Classification model trained a time} \qquad \mathbf{t}_j \quad \text{Using data} \quad \{\mathbf{x}_i, y_i\} \text{ such that } t_i + \delta_l \leq t_j \\ \mathbf{h}_{t_j} \quad \text{can only make prediction on} \quad \left\{\mathbf{x}_k\right\} \text{ such that } t_j + \delta_d \leq t_k \\ \text{sched} = \left\{\mathbf{t}_1 \dots t_n\right\} \quad \text{Retraining schedule determining the sequence of models} \quad \mathbf{H} = \left\{\mathbf{h}_{t_1} \dots h_{t_j} \dots h_{t_n}\right\} \\ \mathbf{X}_i \quad \text{is classified as} \qquad \hat{y}_i = h_{t_*}(x_i) \text{ where } t_* = \max\{t_k \in sched \text{ s.t. } t_k + \delta_{prod} \leq t_i\} \\ \text{We evaluate the schedule} \quad \text{sched} = \left\{\mathbf{t}_1 \dots t_n\right\} \quad \text{We evaluate the schedule of effectiveness and cost} \\ s = score(Y, \hat{Y}) = MCC(Y, \hat{Y}) \\ c = cost(sched) = |H| \end{array}$

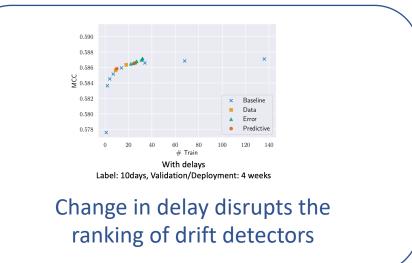
Problem definition of delays in drift mitigation





Empirical study on a real world financial system











Conclusion

 δ_d Labelling delay δ_l Deployment delay

 $\mathbf{h}_{t_i}: \mathcal{X} \to \mathcal{Y}$ Classification model trained a time \mathbf{t}_i Using data $\{\mathbf{x}_i, y_i\}$ such that $t_i + \delta_l \leq t_i$

 \mathbf{h}_{t_d} can only make prediction on $\{\mathbf{x}_k\}$ such that $t_j + \delta_d \leq t_k$

 $\mathrm{sched} = \{ t_1 \dots t_n \}$ Retraining schedule determining the sequence of models $H = \{ h_t, h_t \}$

 \hat{y}_i is classified as $\hat{y}_i = h_{t_*}(x_i) ext{ where } t_* = \max\{t_k \in sched ext{ s.}\}$

We evaluate the schedule $\operatorname{sched} = \{ t_1 \dots t_n \}$ We evaluate the schedule of effective

$$s = score(Y, \hat{Y}) = MCC(Y, \hat{Y})$$

 $c = cost(sched) = |H|$

Problem definition of delay mitigation





Empirical study on a real world financial system

Checkout our repo

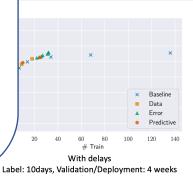
Train model on the train



D_t D_{t+1} D_{t+2} D_{t+3} D_{t+4} D_{t+5} D_{t+6} D_{t+7} Drift Drift detected New model In production New model trained on

Evaluate a single drift detector and parameters with realistic delays

valuation protocol



Change in delay disrupts the ranking of drift detectors



