



An adaptive transfer learning perspective on classification in non-stationary environments

— *Non-parametric Regression for Online Label Shift*



Henry W J Reeve

Lecturer in Statistical Science, University of Bristol

Presented by Yu-Yang Qian

2024.9.12

南

京

大

學

Outline

- ❑ Problem Formulation
- ❑ Background
- ❑ Method
 - Step 1: Reduce Regret Minimization to Parameter Estimation
 - Step 2: Parameter Estimation using *Lepski method*
- ❑ Theory:
 - ***Dynamic Regret***: Interval Regret + Number of Intervals
- ❑ Conclusion

Outline

□ Problem Formulation

□ Background

□ Method

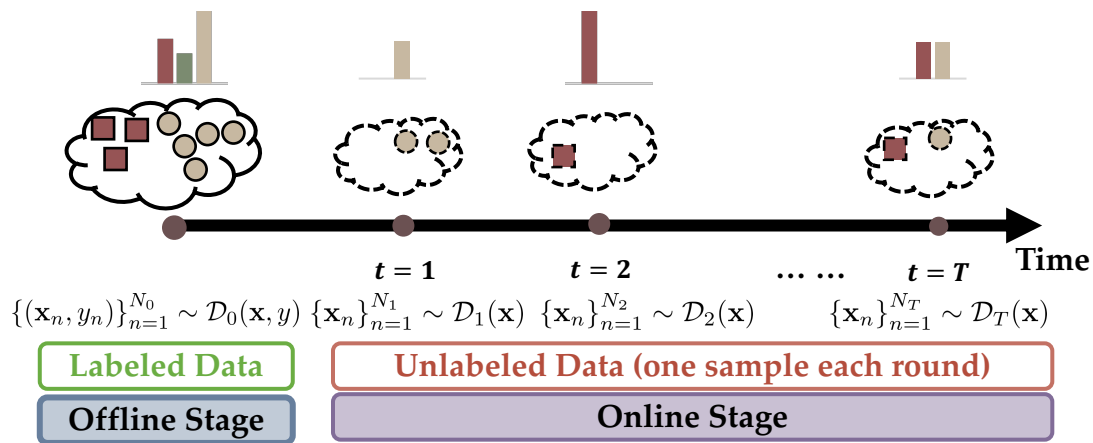
- Step 1: Reduce Regret Minimization to Parameter Estimation
- Step 2: Parameter Estimation using *Lepski method*

□ Theory:

- *Dynamic Regret*: Interval Regret + Number of Intervals

□ Conclusion

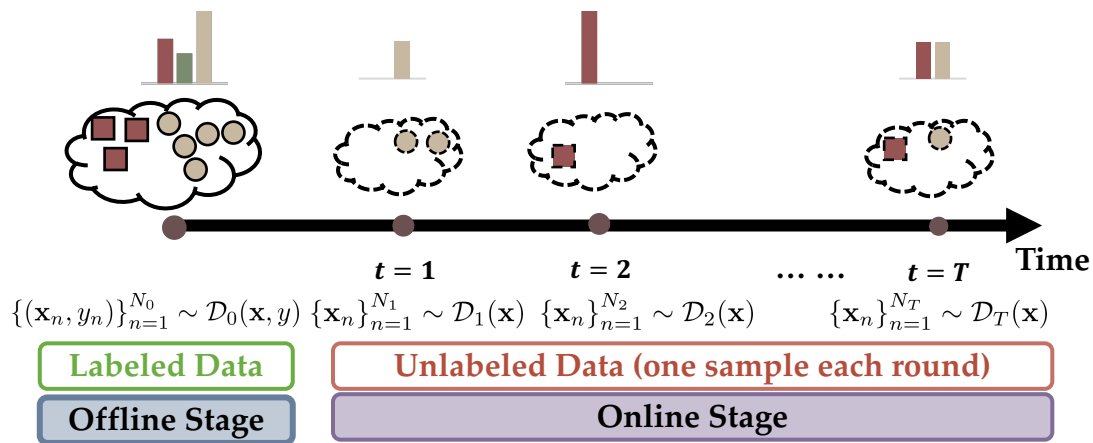
Problem Formulation



□ Two stages: offline & online

- Offline stage: a lot of labeled data
- Online stage: unlabeled data stream (one sample each round)

Problem Formulation



□ Assumption 1: Online Label shift (OLS)

- The conditional $\mathcal{D}_t(\mathbf{x} \mid y)$ is *identical* throughout the process.
- $\mathcal{D}_0(y) > 0$ for any $y \in \mathcal{Y}$.

Problem Formulation

❑ Assumption 1: Online Label shift (OLS)

- The conditional $\mathcal{D}_t(\mathbf{x} \mid y)$ is *identical* throughout the process.
- $\mathcal{D}_0(y) > 0$ for any $y \in \mathcal{Y}$.

❑ Assumption 2: Temporal smoothness (informal)

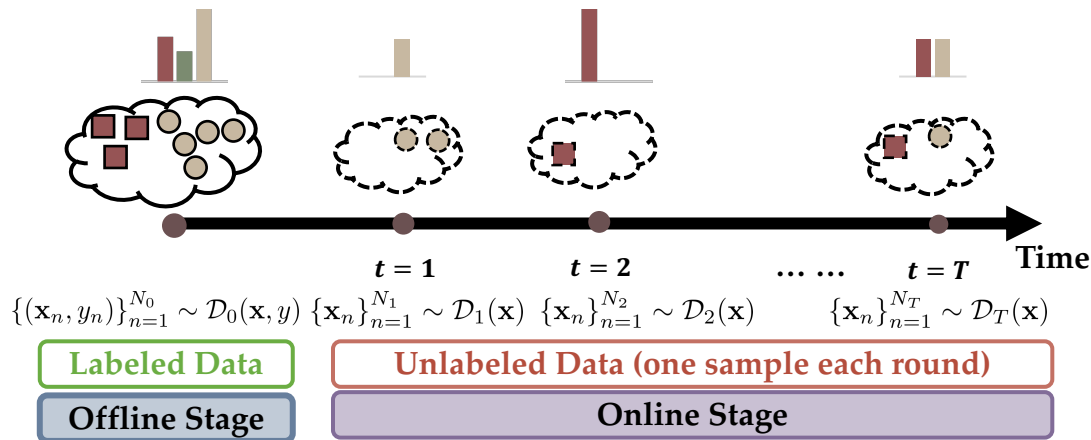
There exists $m \in [t]$, a smooth function g and a smoothness parameter λ such that for all $\ell \in \{t - m, \dots, t\}$,

$$|\pi_t - \pi_{t-\ell}| \leq \lambda \quad \text{where } \pi_t \triangleq \Pr(y_t = 1).$$

i.e., the label prior changes **smoothly** in the interval $I = [t - m, t]$.

Certainly holds for $\lambda = V_{\mathcal{I}} = \sum_{\tau=t-m}^t |\pi_{\tau} - \pi_{\tau-1}|$

Problem Formulation



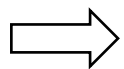
□ Goal: learn a classifier $\hat{h}_t : \mathcal{X} \mapsto \{0, 1\}$ to minimize dynamic regret:

$$\mathbf{Reg}_T^d(\{R_t, h_t^*\}_{t=1}^T) \triangleq \sum_{t=1}^T R_t(\hat{h}_t) - \sum_{t=1}^T R_t(h_t^*)$$

where $R_t(h) = \Pr(h(\mathbf{x}_t) \neq y_t)$.

Previous method for OLS

- ❑ **ATLAS** [Bai et al., 2022]: Unbiased estimator + Online Ensemble:
 - Construct unbiased gradient, then online gradient descent
- ❑ **Restart-based** [Qian et al., 2024]: Previous Online Algorithm + *Explicit Partition*:
 - Achieve a good regret in **stationary** interval, **restart** to find optimal partition
- ❑ **FLH-FTL** [Baby et al., 2023]: Adaptive Regret + *Implicit Partition* (only in analysis):
 - Achieve a good regret in **stationary** interval, **implicitly** find optimal partition



All of them achieving **optimal** dynamic regrets of $\mathcal{O}(T^{\frac{2}{3}} V_T^{\frac{1}{3}})$.

where $V_T = \sum_{t=2}^T |\pi_t - \pi_{t-1}|$ is the class-prior variation.

This paper: A totally statistical way

- ❑ All three methods somehow need “models” for classification.
 - Need to impose prior on model structure
 - e.g., linear model, logistic regression, reweighting the initial model
- ❑ This paper takes a totally *different* view:
 - *Non-parametric Regression* to solve Online Label Shift
 - A totally *statistical* method
- ❑ In the end, we will see how this method *coincident* with adaptive reg.

Outline

- Problem Formulation

- Background

- Method

 - Step 1: Reduce Regret Minimization to Parameter Estimation

 - Step 2: Parameter Estimation using *Lepski method*

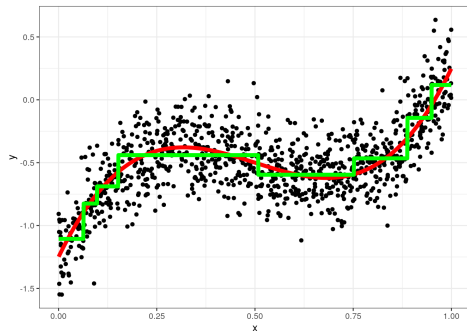
- Theory:

 - *Dynamic Regret*: Interval Regret + Number of Intervals

- Conclusion

Some History Bits: Non-parametric Regression

□ Non-parametric Regression: estimating *certain parameters*



➤ An offline setting: observe a batch of **sampled** data.

Underlying (unknown) oracle labels: $\theta_i = \mathbb{E}(Y \mid X = X_i)$. For simplicity, $X_i = i$.

We can only observe i.i.d. samples $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}, i = 1, \dots, T$,

$$Y_i \sim \mathcal{D}_i, \quad i = 1, \dots, T$$

Goal: minimize ground-truth distributions' loss: $\min \sum_{i=1}^T \mathbb{E}_{Y \sim \mathcal{D}_i} (f(X_i) - Y)^2$

Some History Bits: Non-parametric Regression?

Underlying (unknown) oracle labels: $\theta_i = \mathbb{E}(Y \mid X = X_i)$. For simplicity, $X_i = i$.

We can only observe i.i.d. samples $(X_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}, i = 1, \dots, T$,

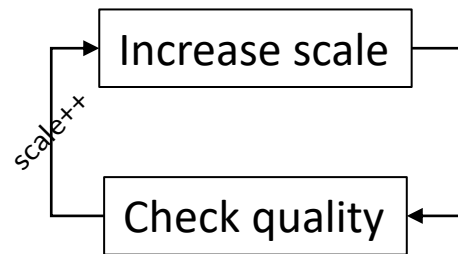
$$Y_i \sim \mathcal{D}_i, \quad i = 1, \dots, T$$

Goal: minimize ground-truth distributions' loss: $\min \sum_{i=1}^T \mathbb{E}_{Y \sim \mathcal{D}_i} (f(X_i) - Y)^2$

- ❑ f can be KNN, spline, local smoothing, historical averaging...
- ❑ Still have 'parameter' in the model, why called non-parametric?
 - 'non-parametric' does not mean we use no parameters in prediction model
 - but mean **we impose no assumption on data generation**
 - OTOH, parametric regression explicitly assume generation function on Y_i

History of Lepski method

- ❑ Lepski method is first introduced in 1997 [Lepski and Spokoiny, 1997]
- ❑ Goal: handle *non-parametric regression* tasks.
- ❑ Key steps (select the optimal trade-off between bias and variance):
 - Increasing the estimator's *scale* (e.g., use more historical data for average, more samples for KNN)
 - The selection of the *optimal scale* is by checking the error bar (e.g., check its error on current data)
 - Once the error bar is *reached*, stop and output



History of Lepski method

- ❑ Lepski method is first introduced in 1997 [Lepski and Spokoiny, 1997]
- ❑ Goal: handle *non-parametric regression* tasks.
- ❑ It is now be applied for many areas, especially *transfer learning*
 - Mostly offline (one-stage) transfer learning
 - For example, [Cai and Wei, 2019]: weighted KNN (of increasing source data size) for transfer learning
- ❑ This paper seems to be the first to apply it to online learning problem.

Outline

- Problem Formulation

- Background

- ▣ Method

 - Step 1: Reduce Regret Minimization to Parameter Estimation

 - Step 2: Parameter Estimation using ***Lepski method***

- Theory:

 - ***Dynamic Regret***: Interval Regret + Number of Intervals

- Conclusion

Method Overview

- ❑ Build entirely upon statistical method:
 - Reduce the Regret Minimization to Parameter Estimation
 - Estimate density ratio by *non-parametric regression*
 - Prove that algorithm can perform well in any *stationary intervals*
 - Find the optimal partition *implicitly* (only in analysis)

$$\mathbf{Reg}_T^{\mathbf{d}} = \sum_{i=1}^J \mathbf{Reg}_{\mathcal{I}_i}$$

intervals: $\{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_J\}$

Step 1: Reduce Regret Minimization to Parameter Estimation

□ Reduce the risk to:

Lemma 1. Suppose $h : \mathcal{X} \rightarrow \{0, 1\}$ is a classifier. Then,

$$R_t(h) = \pi_t + \int h(1 - \pi_t - \eta)d(\nu_0 + \nu_1)$$

where $\pi_t = \Pr(y_t = 1)$, $\nu_0 = \mathcal{D}(\mathbf{x} \mid y = 0)$, $\nu_1 = \mathcal{D}(\mathbf{x} \mid y = 1)$, $\eta(\mathbf{x}) = \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}$

$$\Rightarrow h_t^*(\mathbf{x}) = \mathbb{1} \{1 - \eta(\mathbf{x}) - \pi_t < 0\} = \mathbb{1} \{\eta(\mathbf{x}) + \pi_t > 1\}$$

To this end, we reduce instantaneous risk into a *parameter estimation* problem.

Step 1: Reduce Regret Minimization to Parameter Estimation

$$R_t(h) = R_t(h_t^*) + \int_{\{h \neq h_t^*\}} |1 - \pi_t - \eta| d(\nu_0 + \nu_1)$$

$$\text{where } \pi_t = \Pr(y_t = 1), \nu_0 = \mathcal{D}(\mathbf{x} \mid y = 0), \nu_1 = \mathcal{D}(\mathbf{x} \mid y = 1), \eta(\mathbf{x}) = \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}$$

Proof of Lemma 1. By decomposing $R_t(h)$, we have

$$\begin{aligned} R_t(h) &= \pi_t \int (1 - h) d\nu_1 + (1 - \pi_t) \int h d\nu_0 = \pi_t + \int h d\{(1 - \pi_t)\nu_0 - \pi_t\nu_1\} \\ &= \pi_t + 2 \int h \{(1 - \pi_t)(1 - \eta) - \pi_t\eta\} d\nu_{1/2} = \pi_t + 2 \int h(1 - \pi_t - \eta) d\nu_{1/2} \end{aligned}$$

$$\text{where } \left(\nu_{1/2} = \frac{\nu_1}{2(\nu_0 + \nu_1)} \right)$$

Thus, with $h_t^*(\mathbf{x}) = \mathbb{1} \{ \eta(\mathbf{x}) + \pi_t > 1 \}$ we have

$$R_t(h) - R_t(h_t^*) = 2 \int (h - h_t^*)(1 - \pi_t - \eta) d\nu_{1/2} = 2 \int_{\{h \neq h_t^*\}} |1 - \pi_t - \eta| d\nu_{1/2} \quad \square$$

Step 2.1: Estimating Conditional Distribution

$$\begin{aligned}\pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})},\end{aligned}$$

$$h_t^*(\mathbf{x}) = \mathbb{1} \{ \eta(\mathbf{x}) + \pi_t > 1 \}, \quad h_t(\mathbf{x}) = \mathbb{1} \{ \hat{\eta}(\mathbf{x}) + \hat{\pi}_t > 1 \},$$

□ Next, we estimate conditional distribution $\eta(\mathbf{x})$:

➤ 1. Roadmap: build a *confidence interval*;

$$\hat{U}_b(q, \tilde{\varepsilon}) \triangleq \frac{8 \left(\sqrt{\tilde{\varepsilon} \sigma^2(q) + \tilde{\varepsilon}^2} + \{1 - 2q\} \tilde{\varepsilon} \right)}{3(1 + 2\tilde{\varepsilon})}$$

➤ 2. Use *offline data* to estimate ν_0 and ν_1 , therefore estimate η ;

$$\hat{\nu}_{y,\omega,\delta}(A) \triangleq \begin{cases} \hat{\nu}_y(A) - \hat{U}(1 - \hat{\nu}_y(A), n_y, \delta) & \text{if } \omega = -1 \\ \hat{\nu}_y(A) + \hat{U}(\hat{\nu}_y(A), n_y, \delta) & \text{if } \omega = 1 \end{cases} \quad \hat{\eta}(A) \triangleq \left(\left(\frac{\hat{\nu}_{1,\omega,\delta}(A)}{\hat{\nu}_{0,-\omega,\delta}(A) + \hat{\nu}_{1,\omega,\delta}(A)} \right) \vee 0 \right) \wedge 1$$

Step 2.1: Estimating Conditional Distribution

$$\begin{aligned}\pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})},\end{aligned}$$

$$h_t^*(\mathbf{x}) = \mathbb{1} \{ \eta(\mathbf{x}) + \pi_t > 1 \}, \quad h_t(\mathbf{x}) = \mathbb{1} \{ \hat{\eta}(\mathbf{x}) + \hat{\pi}_t > 1 \},$$

□ Next, we estimate conditional distribution $\eta(\mathbf{x})$:

➤ 3. Theoretical Guarantee: using *Dvoretzky-Kiefer-Wolfowitz-Massart* type concentration inequality:

Lemma 2. With probability at least $1 - \delta$, the estimated conditional distribution $\hat{\eta}$ has the following guarantee:

$$|\hat{\eta}(x) - \eta(x)| \leq 14 \left(\left\{ \theta^{\frac{1}{\beta}} \varepsilon_{\delta}^{\text{il}}(\check{n}) \right\}^{\frac{\beta}{2\beta+1}} \vee \sqrt{\varepsilon_{\delta}^{\text{il}}(\check{n})} \right) = \tilde{\mathcal{O}}\left(\sqrt{\frac{1}{|S_0|}}\right)$$

Details are omitted, as Step 2.1 is a standard statistical learning.

Step 2.2: Estimating Label Prior by *Lepski Method*

$$\begin{aligned}\pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})},\end{aligned}$$

$$h_t^*(\mathbf{x}) = \mathbb{1} \{ \eta(\mathbf{x}) + \pi_t > 1 \}, \quad h_t(\mathbf{x}) = \mathbb{1} \{ \hat{\eta}(\mathbf{x}) + \hat{\pi}_t > 1 \},$$

□ Then, we estimate the label prior $\hat{\pi}_t$:

➤ Remember in ATLAS [Bai et al., 2022], BBSE is used to estimate the label prior:

$$\hat{\pi}_t = C_0^{-1} \cdot \hat{\pi}_{\hat{y}_t}$$

where $[C_0]_{i,j} \triangleq \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_0(\mathbf{x} \mid y=j)} [\Pr(f_0(\mathbf{x}) = i)]$, $[\hat{\pi}_{\hat{y}_t}]_j = 1/n_t \cdot \sum_{\mathbf{x} \in S_t} [f_0(\mathbf{x})]_j$

$f_0 : \mathcal{X} \mapsto \Delta_2$ is the offline prediction model.

➤ However, there is no **offline model** f_0 here, we only have offline data.

Step 2.2: Estimating Label Prior by *Lepski Method*

$$h_t^*(\mathbf{x}) = \mathbb{1} \{ \eta(\mathbf{x}) + \pi_t > 1 \}, \quad h_t(\mathbf{x}) = \mathbb{1} \{ \hat{\eta}(\mathbf{x}) + \hat{\pi}_t > 1 \},$$

$$\begin{aligned} \pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}, \\ \mu_t &= \mathcal{D}_t(\mathbf{x}) \end{aligned}$$

□ Then, we estimate the label prior $\hat{\pi}_t$:

➤ We define $\mu_t = \mathcal{D}_t(\mathbf{x})$, therefore

$$\begin{aligned} \mu_t &= (1 - \pi_t)\nu_0 + \pi_t\nu_1 \\ \Rightarrow \quad \pi_t &= \frac{\mu_t - \nu_0}{\nu_1 - \nu_0} \end{aligned}$$

➤ ν_0 and ν_1 has already been estimated before, therefore we **focus on μ_t**

Step 2.2: Estimating Label Prior by *Lepski Method*

$$\Rightarrow \pi_t = \frac{\mu_t - \nu_0}{\nu_1 - \nu_0}$$

$$\begin{aligned}\pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}, \\ \mu_t &= \mathcal{D}_t(\mathbf{x})\end{aligned}$$

□ ν_0 and ν_1 has already been estimated before, therefore we **focus on μ_t**

➤ We estimate μ_t through **weighted combination** of history:

$$\hat{\mu}_t(\mathbf{x}) \triangleq \sum_{i=1}^q \underbrace{w(i, q)}_{\text{weight of history}} \cdot \underbrace{\hat{D}_{t-i}(X = \mathbf{x})}_{\text{history density}}$$

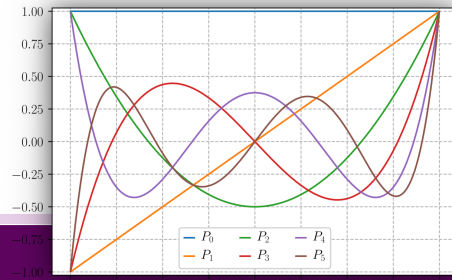
q : history length

weight of history

history density

➤ **weight of history** is defined by *Legendre polynomials* (i.e., predefined orthogonal bases):

$$\mathcal{L}_k(z) := \frac{\sqrt{2k+1}}{k!} \frac{d^k}{dz^k} \{z(z-1)\}^k$$



Step 2.2: Estimating Label Prior by *Lepski Method*

$$\Rightarrow \pi_t = \frac{\mu_t - \nu_0}{\nu_1 - \nu_0}$$

$$\begin{aligned}\pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}, \\ \mu_t &= \mathcal{D}_t(\mathbf{x})\end{aligned}$$

□ ν_0 and ν_1 has already been estimated before, therefore we **focus on μ_t**

➤ We estimate μ_t through **weighted combination** of history:

$$\hat{\mu}_t(\mathbf{x}) \triangleq \sum_{i=1}^q \underbrace{w(i, q)}_{\text{weight of history}} \cdot \underbrace{\hat{D}_{t-i}(X = \mathbf{x})}_{\text{history density}}$$

q : history length

weight of history

history density

➤ Key: How to determine the **history length q** ?

\Rightarrow **Lepski-based** method

Step 2.2: Estimating Label Prior by *Lepski Method*

$$\Rightarrow \pi_t = \frac{\mu_t - \nu_0}{\nu_1 - \nu_0} \quad \hat{\mu}_t(\mathbf{x}) \triangleq \sum_{i=1}^q w(i, q) \cdot \hat{D}_{t-i}(X = \mathbf{x})$$

$$\begin{aligned} \pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}, \\ \mu_t &= \mathcal{D}_t(\mathbf{x}) \end{aligned}$$

- Key: How to determine the **history length q** ?

\Rightarrow ***Lepski-based*** method

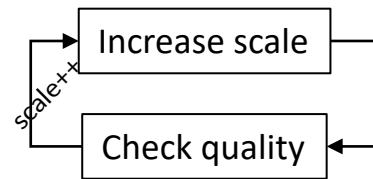
- We define the weight's complexity as:

$$\mathbb{S}_{n,\delta}(q) := \|w(q)\|_2 \sqrt{2 \log \left(\frac{\pi^2 q^2}{\delta} \right)} \quad \text{where } \|w(q)\|_2^2 \triangleq \sum_{i=1}^q w(i, q)^2.$$

We **choose** \hat{q} as the maximal value of $q \in \{8\alpha^2 (\alpha + 1)^2, \dots, n\}$ such that

$$|\hat{\mu}_n^q(f) - \hat{\mu}_n^{q_b}(f)| \leq 2 \{ \mathbb{S}_{n,\delta}(q) + \mathbb{S}_{n,\delta}(q_b) \}$$

for all $q_b \in \{8\alpha^2 (\alpha + 1)^2, \dots, q - 1\}$. (where α is a constant about temporal smoothness)



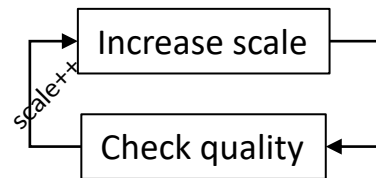
Step 2.2: Estimating Label Prior by *Lepski Method*

$$\Rightarrow \pi_t = \frac{\mu_t - \nu_0}{\nu_1 - \nu_0} \quad \hat{\mu}_t(\mathbf{x}) \triangleq \sum_{i=1}^q w(i, q) \cdot \hat{D}_{t-i}(X = \mathbf{x})$$

$$\begin{aligned} \pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}, \\ \mu_t &= \mathcal{D}_t(\mathbf{x}) \end{aligned}$$

We **choose** \hat{q} as the maximal value of $q \in \{8\alpha^2 (\alpha + 1)^2, \dots, n\}$ such that

$$|\hat{\mu}_n^q(f) - \hat{\mu}_n^{q_b}(f)| \leq 2 \{ \mathbb{S}_{n,\delta}(q) + \mathbb{S}_{n,\delta}(q_b) \}$$



for all $q_b \in \{8\alpha^2 (\alpha + 1)^2, \dots, q - 1\}$. (where α is a constant about temporal smoothness)

Corollary 6. Suppose that Assumption 1 and 2 hold. For all $\delta \in (0, 1)$ we have

$$\mathbb{P} \left\{ \frac{|\hat{\pi}_t - \pi_t| |(\nu_1 - \nu_0)|}{\mathbb{J}_\delta(\mu_t) + 3 \max_{y \in \{0,1\}} |(\hat{\nu}_y^+ - \nu_y)|} > 1 \mid \mathcal{D}_0, \mathcal{D}_1 \right\} \leq \frac{\delta}{3}$$

$$\Rightarrow |\hat{\pi}_t - \pi_t| \leq \tilde{O} \left(\frac{1}{\sqrt{t}} + \lambda \right) \text{ for } t \text{ in smooth intervals (Asp. 2)}$$

Step 2.2: Estimating Label Prior by *Lepski Method*

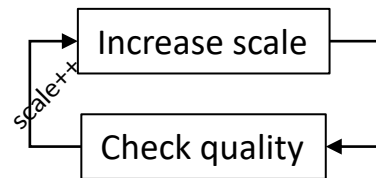
$$\Rightarrow \pi_t = \frac{\mu_t - \nu_0}{\nu_1 - \nu_0} \quad \hat{\mu}_t(\mathbf{x}) \triangleq \sum_{i=1}^{\hat{q}} w(i, q) \cdot \hat{D}_{t-i}(X = \mathbf{x})$$

$$\begin{aligned} \pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}, \\ \mu_t &= \mathcal{D}_t(\mathbf{x}) \end{aligned}$$

We **choose** \hat{q} as the maximal value of $q \in \{8\alpha^2 (\alpha + 1)^2, \dots, n\}$ such that

$$|\hat{\mu}_n^q(f) - \hat{\mu}_n^{q_b}(f)| \leq 2 \{ \mathbb{S}_{n,\delta}(q) + \mathbb{S}_{n,\delta}(q_b) \}$$

for all $q_b \in \{8\alpha^2 (\alpha + 1)^2, \dots, q - 1\}$. (where α is a constant about temporal smoothness)



Assumption 2. There exists $m \in [t]$, a smooth function g and a smoothness parameter λ such that for all $\ell \in \{t - m, \dots, t\}$,

$$|\pi_t - \pi_{t-\ell}| \leq \lambda \quad \text{where } \pi_t \triangleq \Pr(y_t = 1).$$

$$\Rightarrow |\hat{\pi}_t - \pi_t| \leq \tilde{O} \left(\frac{1}{\sqrt{t}} + \lambda \right) \text{ for } t \text{ in smooth intervals (Asp. 2)}$$

Step 2.2: Estimating Label Prior by *Lepski Method*

Corollary 6. Suppose that Assumption 1 and 2 hold. For all $\delta \in (0, 1)$ we have

$$\mathbb{P} \left\{ \frac{|\hat{\pi}_t - \pi_t| |(\nu_1 - \nu_0)|}{\mathbb{J}_\delta(\mu_t) + 3 \max_{y \in \{0,1\}} |(\hat{\nu}_y^+ - \nu_y)|} > 1 \mid \mathcal{D}_0, \mathcal{D}_1 \right\} \leq \frac{\delta}{3}$$

$$\Rightarrow |\hat{\pi}_t - \pi_t| \leq \tilde{O} \left(\frac{1}{\sqrt{t}} + \lambda \right) \text{ for } t \text{ in smooth intervals (Asp. 2)}$$

Proof Sketch. (1) Decompose history distribution using *a set of orthogonal bases* (Legendre polynomials) using *Taylor expansion*

(2) The historical length is tuned to be **optimal** so that it achieve the minimal error (estimated bias + variance)

If interested, you can check the paper's Section 9...

$$\begin{aligned} \pi_t &= \Pr(y_t = 1), \\ \nu_0 &= \mathcal{D}(\mathbf{x} \mid y = 0), \\ \nu_1 &= \mathcal{D}(\mathbf{x} \mid y = 1), \\ \eta(\mathbf{x}) &= \frac{\nu_1(\mathbf{x})}{\nu_0(\mathbf{x}) + \nu_1(\mathbf{x})}, \\ \mu_t &= \mathcal{D}_t(\mathbf{x}) \end{aligned}$$

Outline

- Problem Formulation
- Background
- Method
 - Step 1: Reduce Regret Minimization to Parameter Estimation
 - Step 2: Parameter Estimation using *Lepski method*
- ▣ Theory:
 - **Dynamic Regret:** Interval Regret + Number of Intervals
- Conclusion

Dynamic Regret

Remember: $h_t^*(\mathbf{x}) = \mathbb{1}\{\eta(\mathbf{x}) + \pi_t > 1\}$, $h_t(\mathbf{x}) = \mathbb{1}\{\hat{\eta}(\mathbf{x}) + \hat{\pi}_t > 1\}$,

$$\Rightarrow R_t(h) - R_t(h_t^*) \lesssim \mathcal{O}(|\eta(\mathbf{x}_t) - \hat{\eta}(\mathbf{x}_t)| + |\pi_t - \hat{\pi}_t|)$$

$$\leq \underbrace{\frac{1}{\sqrt{|S_0|}}}_{\text{Error of Step 2.1}} + \underbrace{\left(\frac{1}{\sqrt{t}} + \lambda\right)}_{\text{Error of Step 2.2}}$$

Error of Step 2.1

Error of Step 2.2

Therefore, for an interval $\mathcal{I}_i = [s, e]$ that satisfy Assumption 2 with λ :

$$\mathbf{Reg}_{\mathcal{I}_i} = \sum_{t=s}^e R_t(h) - \sum_{t=s}^e R_t(h_t^*) = \frac{|\mathcal{I}_i|}{\sqrt{|S_0|}} + \underbrace{\sum_{t=s}^e \frac{1}{\sqrt{t-s}}}_{(\leq 2\sqrt{|\mathcal{I}_i|} - 1)} + \lambda|\mathcal{I}_i|$$

Dynamic Regret: sum of interval regret

$$\mathbf{Reg}_{\mathcal{I}_i} \leq \tilde{O} \left(\frac{|\mathcal{I}_i|}{\sqrt{|S_0|}} + \lambda |\mathcal{I}_i| + \sqrt{|\mathcal{I}_i|} \right)$$

Suppose there are J intervals: $\{\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_J\}$

$$\begin{aligned} \mathbf{Reg}_T^d &= \sum_{i=1}^J \mathbf{Reg}_{\mathcal{I}_i} = \frac{T}{\sqrt{|S_0|}} + \lambda T + \sum_{i=1}^J \sqrt{|\mathcal{I}_i|} \\ &\leq \frac{T}{\sqrt{|S_0|}} + \lambda T + \sqrt{J} \left(\sum_{i=1}^J \left(\sqrt{|\mathcal{I}_i|} \right)^2 \right)^{1/2} \\ &= \frac{T}{\sqrt{|S_0|}} + \lambda T + \sqrt{JT} \end{aligned}$$

Dynamic Regret: control number of intervals J

$$\mathbf{Reg}_T^d = \sum_{i=1}^J \mathbf{Reg}_{\mathcal{I}_i} = \frac{T}{\sqrt{|S_0|}} + \lambda T + \sqrt{JT}$$

Consider the simple case of same λ in all intervals:

Assumption 2. There exists $m \in [t]$, a smooth function g and a smoothness parameter λ such that for all $\ell \in \{t - m, \dots, t\}$,

$$|\pi_t - \pi_{t-\ell}| \leq \lambda \quad \text{where } \pi_t \triangleq \Pr(y_t = 1).$$

Asp. 2 certainly holds for $\lambda = V_{\mathcal{I}_i} = \sum_{t=s}^e |\pi_t - \pi_{t-1}|$, therefore we choose $\lambda \leq V_{\mathcal{I}_i}$

$$\lambda J \leq \sum_{j=1}^J V_{\mathcal{I}_j} = V_T \quad \Longrightarrow \quad J \leq \frac{V_T}{\lambda}$$

Dynamic Regret: finish the proof

$$\mathbf{Reg}_T^d = \sum_{i=1}^J \mathbf{Reg}_{\mathcal{I}_i} \leq \frac{T}{\sqrt{|S_0|}} + \lambda T + \sqrt{\frac{V_T T}{\lambda}}$$

Finally, tune λ to get the regret bound:

$$\begin{aligned} \mathbf{Reg}_T^d &= \sum_{i=1}^J \mathbf{Reg}_{\mathcal{I}_i} \leq \frac{T}{\sqrt{|S_0|}} + \lambda T + \frac{1}{2} \sqrt{\frac{V_T T}{\lambda}} + \frac{1}{2} \sqrt{\frac{V_T T}{\lambda}} \\ &\leq \mathcal{O} \left(\frac{T}{\sqrt{|S_0|}} + (V_T T^2)^{1/3} \right) \quad \left(\lambda = V_T^{1/3} T^{-1/3}, \quad J = T^{1/3} V_T^{2/3} \right) \\ &\leq \mathcal{O} \left(\max \left(\frac{T}{\sqrt{|S_0|}} + (V_T)^{1/3} T^{2/3}, \sqrt{T} \right) \right) \quad (J \geq 1) \quad \square \end{aligned}$$

Compared with **FLH-FTL** [Baby et al., 2023]:

- ❑ Get a good regret in any interval (by Lepski method)
- ❑ Implicit find the optimal partition (only in analysis)
- ❑ Get a good regret in any interval (by Strongly Adaptive Methods)
- ❑ Implicit find the optimal partition (only in analysis)

Getting the adaptive regret using non-parametric method!

If fact: if check partition in this paper:

$$\left\{ \begin{array}{l} \text{number of intervals: } J = T^{1/3} V_T^{2/3} \\ \text{Regret in each interval: } \mathcal{O}(\sqrt{\mathcal{I}_i} + 1) \end{array} \right.$$

Matches with Baby's best partition!

Lemma 5 (key partition) Initialize $\mathcal{Q} \leftarrow \Phi$. Starting from time 1, spawn a new bin $[i_s, i_t]$ whenever $\sum_{j=i_s+1}^{i_t+1} |u_j - u_{j-1}| > B/\sqrt{n_i}$, where $n_i = i_t - i_s + 2$. Add the spawned bin $[i_s, i_t]$ to \mathcal{Q} . Consider the following post processing routine.

1. Initialize $\mathcal{P} \leftarrow \Phi$.
2. For $i \in [|\mathcal{Q}|]$:
 - if $u_{i_t} = u_{i_t+1}$:
 - (a) Let p be the largest time point with $u_{p:i_t}$ being constant and let q be the smallest time point with $u_{i_t+1:q}$ being constant.
 - (b) Add bin $[i_s, p-1]$ to \mathcal{P} .
 - (c) If $(i+1)_s > q$ then add $[p, q]$ to \mathcal{P} and set $(i+1)_s \leftarrow q+1$.
 - (d) Goto Step 2.
 - Add $[i_s, i_t]$ to \mathcal{P} . Goto Step 2.

Let $M := |\mathcal{P}|$. We have $M = O\left(1 \vee n^{1/3} C_n^{2/3} B^{-2/3}\right)$. Further for any bin $[i_s, i_t] \in \mathcal{P}$, it holds that $\sum_{j=i_s+1}^{i_t} |u_j - u_{j-1}| \leq B/\sqrt{n_i}$ where $n_i = i_t - i_s + 1$.

Compared with *FLH-FTL* [Baby et al., 2023]:

□ Pros

- No diameter assumption
- Truly non-parametric, parameter-free

□ Cons

- Very costly: computational and storage is $O(T)$, can hardly be runnable in practice
- Need to derive specific method for specific problem structure
(可能一旦不是OLS/密度估计问题, 就很难再有良好保障)

Take home message

- ❑ A new view of handling non-stationary environments:
 - A ***non-parametric regression*** method (Lepski method) can also get an *strongly adaptive regret*
 - Use ***implicit partition*** (only in analysis) to get the optimal dynamic regret
- ❑ May have a chance to accelerate:
 - Lazy update the summation of history margin distribution
 - A segment tree to get the sum of history intervals
 - May achieve $O(\log T)$ complexity each round

Thanks!