Online multiple testing with e-values

https://arxiv.org/pdf/2311.06412

E-readers June 14th, 2024



Ziyu Xu (Neil) (CMU)



Joint work w/ Aaditya Ramdas (CMU)

> Carnegie Mellon University

Outline

- Online multiple testing
- Doubly sequential inference
- Prior work
- Our method: e-LOND (and proof of FDR control)
- Empirical results
- Stochastic rounding (randomization) for more power
- Weighted conformal selection

Online multiple testing problem

Stream of hypotheses we wish to test:



Doubly sequential inference

Imagine we have teams of data scientists doing A/B testing...

(1)

(2)



Sequential testing requires e-processes

Consider a stream of data $X_1, X_2, \ldots,$

and we construct a p-value P_i at each step *i* from X_1, \ldots, X_i

We wish to have an *anytime-valid p-value:* a fixed time (standard) p-value: $\mathbb{P}(\exists i \in \mathbb{N} : P_i \leq s) \leq s \text{ for all } s \in [0,1]$ $\mathbb{P}(P_i \leq s) \leq s \text{ for all } s \in [0,1] \text{ and } i \in \mathbb{N}$

This means that rejecting the null at step $\tau = \min\{i : P_i \le \alpha\}$, i.e., *early stopping*, is a valid test.

Anytime-valid p-values are frequently used in A/B testing [1] and traditionally has used likelihood ratios (e.g. mSPRT [2]).

An <u>e-process</u> is a nonnegative process (E_i) where E_{τ} is an e-value for any stopping time τ . All known anytime-valid p-values can be written as $P_i = 1/\max_{\substack{i \le i}} E_i$

(or is dominated by such an anytime-valid p-value).

Practically, we always have access to e-processes in the sequential regime.

- [1] Johari, Pekelis, Walsh. Always Valid Inference: Continuous Monitoring of A/B Tests. Operations Research (2022)
- [2] Robbins. Statistical Methods Related to the Law of the Iterated Logarithm. Ann. Math. Stat (1970)
- [3] Ramdas, Ruf, Larsson, Koolen. Admissible anytime-valid sequential inference must rely on nonnegative martingales. arXiv (2023)

Overview of online multiple testing



Foster and Stine (2008) Saharoni and Rosset (2014)

Tian and Ramdas (2019)

 $\mathsf{mFDR}(R_T) := \frac{\mathbb{E}[|\mathcal{N} \cap R_T|]}{\mathbb{E}[|R_T|]}$

mFDR control at stopping times (not just fixed times).

ADDIS = SAFFRON + discarding of conservative nulls

Weinstein and Ramdas (2020)

Online control FCR for selective confidence intervals

	All of the above req	uire <u>independent</u> or <u>conditional (on the past)</u> p-values (or confidence intervals).	
Zrnic et. al (2020, 2021)		Extends LORD and SAFFRON for batches, under local dependence and arbitrary dependence	
E	ao et al. (2024)	Online FCR control for selective conformal prediction interv (subsequent)	/als

See Robertson et al. (2023) for a comprehensive survey

Prior work

r-LOND algorithm [3]

$$\alpha_t^{\mathsf{r}\text{-}\mathsf{LOND}} := \alpha \gamma_t \cdot (|R_{t-1}| + 1)/\ell_t$$

where (γ_t) are nonnegative and sum to 1

$$\mathscr{C}_t := \Sigma_{i=1}^t 1/i \approx \log t$$

Theorem [3]: FDR(R_t) $\leq \alpha$ for all $t \in \mathbb{N}$ for arbitrarily dependent *p*-values (P_t) is guaranteed by *r*-LOND.

LOND algorithm [4]

$$\alpha_t^{\text{LOND}} := \alpha \gamma_t(|R_{t-1}| + 1)$$

Theorem : $FDR(R_t) \le \alpha$ for all $t \in \mathbb{N}$ for *independent* [4] *or positively dependent PRDS p-values* [3] (P_t) is guaranteed by *LOND*.

Our contribution: With e-values, we can use the more the powerful LOND, even under arbitrary dependence.

- [1] Blanchard and Roquain. Two simple sufficient conditions for FDR control. *Elec. J. Stat. (2008)*
- [2] Benjamini and Yekutieli. The control of the false discovery rate in multiple testing under dependency. Ann. of Stat. (2001)
- [3] Zrnic, Ramdas, Jordan. Asynchronous online testing of multiple hypotheses. JMLR (2021)
- [4] Javanmard and Montanari. Online rules for control of false discovery rate and false discovery exceedance. Ann. of Stat. (2018)

e-LOND: FDR control under arbitrary dependence

e-LOND algorithm

$$\alpha_t^{\text{e-LOND}} := \alpha \gamma_t \cdot (|R_{t-1}| + 1)$$

Theorem (ours): FDR(R_t) $\leq \alpha$ for all $t \in \mathbb{N}$ for *arbitrarily dependent e-values* (E_t) is guaranteed by *e-LOND*, and $\alpha_t^{e-LOND} > \alpha_t^{r-LOND}$ for all $t \in \mathbb{N}$, i.e., strictly dominates r-LOND.



More power through randomization

Let (U_t) be a sequence of random variables that are uniformly distributed on [0,1] and independent of (E_t) or (P_t)

Ue-LOND algorithm

 $\alpha_t^{\text{Ue-LOND}} := \alpha \gamma_t \cdot (|R_{t-1}| + 1)/U_t$

Idea: Let $\hat{\alpha} \in [0,1]$ and E > 0 be dependent RVs. Define the <u>stochastically rounded e-value</u> [1]

 $S_{\widehat{\alpha}}(E) := \widehat{\alpha}^{-1} \cdot \mathbf{1}\{U \le \widehat{\alpha}E\}$

Has bounded expectation $\mathbb{E}[E] \ge \mathbb{E}[S_{\widehat{\alpha}}(E)]$

Discovery threshold property $\mathbf{1}\{S_{\widehat{\alpha}}(E) \geq \widehat{\alpha}^{-1}\} = \mathbf{1}\{E \geq U \cdot \widehat{\alpha}^{-1}\}$

Ur-LOND algorithm

$$\alpha_t^{\text{Ur-LOND}} := \alpha \gamma_t \cdot (\lfloor (|R_{t-1}| + 1)/(\ell_t \cdot U_t) \rfloor \wedge t$$

Similar ideas can be applied to improve r-LOND (by its equivalence to being calibrated p-value + e-LOND)

Theorem (ours):

- 1. FDR(R_t) $\leq \alpha$ for all $t \in \mathbb{N}$ for arbitrarily dependent e-values (E_t) is guaranteed by Ue-LOND, and $\alpha_t^{\text{Ue-LOND}} > \alpha_t^{\text{e-LOND}}$ for all $t \in \mathbb{N}$ almost surely.
- 2. FDR(R_t) $\leq \alpha$ for all $t \in \mathbb{N}$ for arbitrarily dependent *p*-values (P_t) is guaranteed by Ur-LOND, and has strictly greater power than *r*-LOND.

[1] Xu and Ramdas. More powerful multiple testing under dependence via randomization. arXiv (2023)

Empirical results



(More simulation results in paper)

Selection by conformal prediction (Jin and Candès 2023)

Assumption: exchangeability

We have a calibration dataset of size $n: (X_1, Y_1), \dots, (X_n, Y_n)$

We also receive a batch of *m* test points: X_{n+1}, \ldots, X_{n+m}

Assume that $(X_1, Y_1), \ldots, (X_{n+m}, Y_{n+m})$ are exchangeable (calibration dataset + stream are all exchangeable).

Assumption: monotonic score function

Let *V* be a scoring function that is <u>monotonic</u> in its second argument, i.e., $y \le c \Rightarrow V(x, y) \le V(x, c)$ for all *x* **Example:** $V(x, y) := y - \hat{\mu}(x)$

 $V_i := V(X_i, Y_i)$ for $i \in [n]$, $\hat{V}_{n+t} := V(X_{n+t}, c_{n+t})$ for $t \in \{1, 2, ...\}$

For each $t \in \{1,2,...\}$, our goal is to test the stochastic null hypothesis, i.e., $H_{0,t}: Y_{n+t} \leq c_{n+t}$ where c_{n+t} is a known fixed or data dependent threshold. **Goal:** Large rejection sets *R* with FDR control.

Conformalized selection (batch, no weighting)
1. Construct conformal p-values
$$P_t := \frac{\sum_{i=1}^{n} \mathbf{1}\{V_i \le \hat{V}_{n+t}\}}{n+1}$$
.
2. Apply the BH (Benjamini-Hochberg) procedure, i.e.,
 $\hat{k} := \max\left\{k \in [m] : \sum_{i=1}^{m} \mathbf{1}\{P_t \le \alpha k/m\} \ge k\right\}$, reject
 \hat{k} smallest p-values

Proof idea: Consider "oracle" p-values: $\bar{P}_{t} := \frac{\sum_{i=1}^{n} \mathbf{1}\{V_{i} \leq V_{n+t}\}}{n+1}$ $\{H_{0,t} \text{ is true }\} \Rightarrow \{\bar{P}_{t} \leq P_{t}\} \text{ (by monotonicity)}$ $P_{1}, \dots, \bar{P}_{t}, \dots, P_{m} \text{ are positively dependent, i.e., PRDS}$ (Bates et. al. 2021). BH provides FDR control for PRDS p-values.

Weighted conformalized selection (Jin and Candès 2023)

Assumption: known covariate shift (weighted exchangeability)

The calibration dataset $(X_1, Y_1), \ldots, (X_n, Y_n)$ are i.i.d. draws from \mathscr{P}

We also receive a batch of *m* test points: X_{n+1}, \ldots, X_{n+m} where $(X_{n+1}, Y_{n+1}), \ldots, (X_{n+m}, Y_{n+m})$ are i.i.d. draws from QCovariate shift is known, i.e., $w(x) := \frac{dQ}{d\mathcal{P}}(x, y)$



FDR control for WCS batch



Hence $E_t := \frac{K}{\alpha |\hat{R}_t|} \mathbf{1} \left\{ P_t \le \frac{\alpha |\hat{R}_t|}{K} \right\}$ satisfies $\mathbb{E}[E_t \cdot \mathbf{1}\{H_{0,t} \text{ is true}\}] \le \mathbb{E}[\bar{E}_t] \le 1$

(can deploy stochastic rounding to have power when $|\hat{R}_t| < |R|$)

Default approach to online weighted conformalized selection



 $\begin{array}{l} \begin{array}{l} \text{Online weighted conformalized selection (attempt)} \\ \text{Do the following for each } t \in \{1,2,\ldots\} \\ 1. \qquad \text{Construct multiple weighted conformal p-values (i.e., for each } \ell^{\circ} \in [t]): \\ P_{\ell}^{(t)} := \displaystyle \frac{w(X_{n+t})\mathbf{1}\{\hat{V}_{n+t} \leq \hat{V}_{n+\ell}\} + \sum\limits_{i=1}^{n} w(X_i)\mathbf{1}\{V_i \leq \hat{V}_{n+\ell}\}}{w(X_{n+t}) + \sum\limits_{i=1}^{n} w(X_i)} \\ 2. \qquad \text{Apply LOND to } (P_{\ell}^{(t)})_{\ell \in [t-1]} \text{ to derive } \hat{R}_{t-1}. \\ 3. \qquad \text{Let } E_t := \displaystyle \frac{\mathbf{1}\{P_t \leq \alpha\gamma_t(|\hat{R}_{t-1}|+1)\}}{\alpha\gamma_t(|\hat{R}_{t-1}|+1)}. \\ \text{Apply e-LOND or (Ue-LOND) to } (E_t) \text{ to get } (R_t) \end{array}$



Our approach to online weighted conformalized selection



Theorem (ours): FDR(R_t) $\leq \alpha$ for all $t \in \mathbb{N}$ when e-LOND or Ue-LOND is applied to (E_t) as defined above.

And we can also show a result for arb. dep. p-values as well:

Theorem (ours): FDR(R_t) $\leq \alpha$ for all $t \in \mathbb{N}$ when r-LOND or Ur-LOND is applied to (P_t), i.e., the standard weighted p-values.

Results on drug discovery dataset

Drug discovery dataset of X = chemical structure, Y = binds to target protein We have access to neural network predictor $\hat{\mu}$ and $\bar{\mu} = \frac{1}{n_{\text{train}}} \sum_{i=1}^{n_{\text{train}}} \hat{\mu}(X_i)$

Sample *t*th drug (i.e., run an experiment to learn Y_{n+t}) to determine if it binds with probability $p(X_{n+t}) := \operatorname{sigmoid}(\widehat{\mu}(X_{n+t}) - \overline{\mu}) \wedge 0.8$

Resulting covariate shift of $w(x) \propto 1/p(x)$.



Extensions and conclusion

Takeaways:

- Unknown dependence is commonplace in the doubly sequential framework modern data science falls into.
- With e-values, we can avoid correction when controlling the FDR under unknown or arbitrary dependence.
- We can utilize randomization to maximize power for e-LOND and r-LOND
- We can develop an online weighted conformal selection procedure using e-values to handle dependence between weighted conformal p-values.
- We can also control the false coverage rate (FCR) of e-value based CIs for any selection rule using a CI analog of e-LOND (see paper)