Combining Evidence Across Filtrations Using Adjusters

https://arxiv.org/abs/2402.09698 August 2024



Yo Joong "YJ" Choe Data Science Institute University of Chicago yjchoe.github.io



Aaditya Ramdas Statistics & Data Science; Machine Learning Carnegie Mellon University <u>stat.cmu.edu/~aramdas</u>

Outline

- 1. Motivation: Combining e-processes for testing exchangeability
- 2. p-lifting and e-lifting: Lifting evidence into finer filtrations via adjusters
- 3. Experiments & further illustrative examples
- 4. Implications & new characterizations of adjusters for e-processes
- 5. Randomized adjustments for e-lifting

Evidence Measures for Anytime-Valid Inference

- $\mathbb{G} = (\mathcal{G}_t)_{t>0}$: filtration
- τ : any G-stopping time
- P : point null
- \mathcal{P} : composite null

Test Supermartingale $(M_t)_{t>0}$ for P w.r.t. G

- 1. $(M_t)_{t>0}$ is **adapted** to \mathbb{G} .
- 2. $M_0 = 1$ and $M_t \ge 0$, $\forall t$.
- 3. $\mathbb{E}_{\mathsf{P}}[\mathsf{M}_{\mathsf{t}} \mid \mathscr{G}_{\mathsf{t}-1}] \leq \mathsf{M}_{\mathsf{t}-1}, \forall \mathsf{t}.$



Optional Stopping

 $p_t = 1/e_t^*$ (Ville's Inequality)

What goes wrong when combining e-processes across filtrations?

Example: Testing Exchangeability "Is your data stream actually random?"



- it is a nontrivial e-process for testing exchangeability (Ramdas et al., IJAR 2022).

• We want to sequentially test whether a binary stream of data X_1, X_2, \ldots is **exchangeable**:

 \mathscr{P}^{exch} : X₁, X₂,... is exchangeable.



• This is a composite null for which no nontrivial test martingales exist in the data filtration.

• $(e_t)_{t>0}$ is a nontrivial e-process for testing **randomness** ("Is the data i.i.d.?") if and only if

Example: Testing Exchangeability "Is your data stream actually random?"

It turns out that there are two different methods to construct an e-process for $\mathscr{P}^{\mathsf{exch}}$:

- 1. Universal inference (UI) e-process (Ramdas et al., 2022): $e_t^{UI} = \frac{\text{mixture over Markov alternatives}}{\text{maximum likelihood under null}}$.
 - Powerful against <u>Markovian</u> alternatives.
 - Anytime-valid w.r.t. the data ("full") filtratic
- 2. Conformal test martingale (Vovk, 2021): $e_t^{conf} =$
 - Powerful against <u>changepoint</u> alternatives.
 - This e-process is ONLY anytime-valid w.r.t. a coarse filtration \mathbb{G} , $\mathcal{G}_t = \sigma(p_1, \dots, p_t)!$

on
$$\mathbb{F}$$
, $\mathcal{F}_t = \sigma(X_1, \dots, X_t)$.

$$=\prod_{i=1}^{t} \left[1 + \lambda \left(p_{i} - \frac{1}{2}\right)\right], \text{ where } p_{i} \text{ are conformal } p\text{-values.}$$

E-process w.r.t. a coarse filtration is NOT anytime-valid in the data filtration (in general)

Data: from i.i.d. Bernoulli. $\tau^{\mathbb{F}}$ = first time we observe five consecutive 0's.



Over 10,000 repeated trials, $\hat{\mathbb{E}}_{P}[e_{\tau^{\mathbb{F}}}^{conf}] \approx 1.33 \pm 0.02$.

The conformal test martingale only has "restricted" anytime-validity, as it does NOT allow stopping w.r.t. the data filtration.





We can't just average the two to obtain an e-process...

• What happens if we just try to take the average anyway?

- In \mathbb{F} , $(m_t)_{t>0}$ is **not** an e-process because $(e_t^{conf})_{t>0}$ is not \mathbb{F} -anytime-valid, as we just saw.
- In G, $(m_t)_{t>0}$ is also **not** an e-process because $(e_t^{UI})_{t>0}$ is not G-adapted.
- So, $(m_t)_{t>0}$ is not an e-process w.r.t. either filtration.

- $m_t = \frac{1}{2} \left(e_t^{UI} + e_t^{conf} \right), \quad \forall t.$

Vladimir Vovk Alexander Gammerman **Glenn Shafer**

Algorithmic Learning in a Random World

Second Edition



In fact, the average of randomized exchangeability martingales is not guaranteed to be a randomized exchangeability martingale. Martingales in a fixed filtration form a linear space, but different runs of the simplified Bayes-Kelly martingale are martingales in different filtrations, as discussed earlier. Figure 9.14 suggests that the averaged simplified Bayes–Kelly martingale is not an exchangeability martingale any more.

From Section 9.3 (p. 296)

Example: Testing a Scale-Invariant Gaussian Mean From Pérez-Ortiz et al. (2022)

 $\mathcal{H}_{0}: \delta = \delta_{0}$

• Let F be the "full" data filtration, and let G denote the scale-invariant **coarsening** of F: $\mathcal{G}_{t} = \sigma \left(\frac{X_{1}}{|X_{1}|} \right)$

• In G, a GROW e-process $(e_t)_{t>0}$ for $(\mathcal{H}_0, \mathcal{H}_1)$ can be derived. However, it is also shown that this e-process is not anytime-valid w.r.t. F:

If $\tau^{\mathbb{F}} = 1 + 1(|X_1| \in [0.44, 1.70])$, then $\mathbb{E}[e_{\tau^{\mathbb{F}}}] \approx 1.19 > 1$.

Suppose the data X_1, X_2, \ldots is sampled from $\mathcal{N}(\mu, \sigma^2)$, and let $\delta = \mu/\sigma$. Consider testing

vs.
$$\mathcal{H}_1: \delta = \delta_1$$
.

$$\frac{X_{t}}{|X_{1}|}, \dots, \frac{|X_{t}|}{|X_{1}|}, \quad \forall t.$$



How can we combine e-processes across different filtrations?

(Especially, if the e-process in the coarser filtration isn't valid in the finer one.)

p-lifting & e-lifting: Lifting evidence across filtrations

p-lifting: P-processes can be lifted "freely"

- This result follows from the so-called **lifting lemma**, which is an extension of the "equivalence lemma" (Ramdas et al., 2020; Howard et al., 2021) to pairs of filtrations.
- Essentially, any "probability statements" of anytime-validity translate across filtrations.
 - In contrast, analogous "expectation statements" of anytime-validity do not translate.

Theorem (p-lifting). Suppose $\mathbb{G} \subseteq \mathbb{F}$, and let $(\mathfrak{p}_t)_{t>0}$ be a p-process for \mathscr{P} w.r.t. \mathbb{G} . Then, $(\mathfrak{p}_t)_{t>0}$ is a p-process for \mathscr{P} w.r.t. F.



The Equivalence Lemma Ramdas et al. (2020); Howard et al. (2021)

- Let $(\xi_t)_{t\geq 1}$ be a sequence of events adapted to a filtration G. (E.g., $\xi_t = \{p_t \leq \alpha\}$.)
- Given any probability P and any $\alpha \in (0, 1)$, the following statements are <u>equivalent</u>:
- (a) Time-uniform validity: $P(\bigcup_{t\geq 1} \xi_t) \leq \alpha$.
- (b) Random time validity: for any (possibly infinite) random time T, $P(\xi_T) \leq \alpha$.
- (c) G-anytime-validity: for any (possibly infinite) G-stopping time τ^{G} , $P(\xi_{\tau^{G}}) \leq \alpha$.

The Lifting Lemma Adaptation of the equivalence lemma to two filtrations

Let $(\xi_t)_{t>1}$ be a sequence of events adapted to a sub-filtration $\mathbb{G} \subseteq \mathbb{F}$. Given any probability P and any $\alpha \in (0, 1)$, the following statements are <u>equivalent</u>: (a) G-anytime-validity: for any (possibly infinite) G-stopping time τ^{G} , $P(\xi_{\tau^{G}}) \leq \alpha$. (b) F-anytime-validity: for any (possibly infinite) F-stopping time $\tau^{\mathbb{F}}$, $P(\xi_{\tau^{\mathbb{F}}}) \leq \alpha$.

<u>Implication</u>: anytime-validity of any p-process w.r.t. $\mathbb{G} =>$ anytime-validity w.r.t. \mathbb{F} .

e-lifting: Lifting e-processes via adjusters

let $(e_t)_{t>0}$ be an e-process for \mathscr{P} w.r.t. G. Then

<u>Proof</u>: 1. By Ville's inequality, $\mathfrak{p}_t = 1/\mathfrak{e}_t^*$ is a p-process for \mathscr{P} w.r.t. G.

2. By p-lifting, $p_t = 1/e_t^*$ is also a p-process for \mathcal{P} w.r.t. F.

4. Thus, $e_t^{adj} = C(\mathfrak{p}_t) = A(e_t^*)$ is an e-process for \mathscr{P} w.r.t. F.

- <u>Theorem (e-lifting).</u> Let A be an adjuster (to be defined soon). Suppose $\mathbb{G} \subseteq \mathbb{F}$, and
 - $(A(e_t^*))_{t>0}$ is an e-process for \mathscr{P} w.r.t. F.

 $(e_t^* = \sup_{i \le t} e_i)$

- 3. The adjuster has a corresponding **p-to-e calibrator** C, such that $A(e) = C(1/e), \forall e \ge 1$.



Combining evidence across filtrations via e-lifting

: anytime-valid X : NOT anytime-valid





What are adjusters? a.k.a. lookback adjusters & martingale calibrators

• An increasing, right-continuous function A: $[1, \infty] \rightarrow [0, \infty]$ is an **adjuster** if it satisfies:

$$\int_{1}^{\infty} \frac{A(e)}{e^2} de \le 1.$$

- It is *admissible* if the above holds with equality and $A(\infty) = \infty$.
- Adjusters allow betting on the running maximum of a test supermartingale (or a capital process).



Dawid et al. (2011a;b), Shafer et al. (2011); Koolen & Vovk (2014) 18



Adjuster **A**_{mix} A_{KV} ······ A_{sqrt} ----- *A*_{1/2} $---- A_{1/4}$



Adjusters \iff P-to-E Calibrators

• A decreasing, left-continuous function $C : [0, 1] \rightarrow [0, \infty]$ is a (p-to-e) calibrator if



- It is admissible if the above holds with equality.
- Setting A(e) = C(1/e), and by change-of-variables ($\mathfrak{p} = 1/e$),

$$\int_{1}^{\infty} \frac{A(e)}{e^2} de = \int_{1}^{\infty} \frac{C(1/e)}{e^2} de = \int_{0}^{1} C(\mathfrak{p}) d\mathfrak{p} \leq 1.$$

• There is a straightforward 1-to-1 correspondence between calibrators and adjusters.

cf. Shafer et al. (2011); Vovk & Wang (2021)



Experiments & illustrative examples

Testing Exchangeability: Null Case

The lifted e-process is \mathbb{F} -anytime-valid, so we can now combine it with the UI e-process.

Data: from i.i.d. Bernoulli.

 $\tau^{\mathbb{F}}$ = first time we observe five consecutive 0's.





Across 10,000 simulations, $\mathbb{E}_{P}[A(e_{\tau^{\mathbb{F}}}^{*})] \approx 0.47 \leq 1.$



Testing Exchangeability: Alternative Case



Alternative #1: First-order Markov

The combined e-process achieves power against both alternatives.

Combined ("eLift+Avg"): $\bar{e}_t = \frac{1}{2} \left[e_t^{UI} + A((e_t^{conf})^*) \right]$ 10¹⁰ **E**-process



Alternative #2: Two changepoints

Other topical examples in the literature

- information that they had at the time of forecasting.
 - the data filtration, say $\mathbb{G}^{[k]} \subsetneq \mathbb{F}$.
- 2. Sequential independence testing. For testing independence sequentially, there is no nontrivial test martingale w.r.t. the data filtration (Henzi & Law, 2023).

¹Henzi & Ziegel (2022); Arnold et al., (2022); Choe & Ramdas (2023) ²Balasubramani & Ramdas (2016); Shekhar & Ramdas (2023); Podkopaev et al. (2023); Henzi & Law (2023)

1. Multi-step forecast evaluation/comparison. When evaluating sequential forecasters making their forecasts h > 1 days ahead of time, we'd want to evaluate them conditioned on the

• For each offset k = 1, ..., h, there exists an e-process $(e_t^{[k]})_{t>0}$ w.r.t. different coarsenings of

To obtain an evaluation across all offsets, we would need to e-lift all h e-processes!

When combining e-processes for this null, existing e-processes have to be e-lifted.



23

Implications & new characterizations of adjusters for e-processes

Implications on coarsening the filtration

- there exist ones in a coarser filtration?
 - Unlike in the case of test supermartingales, e-lifting implies that the answer is no!
 - - admissible A), so the adjusted e-process is also powerful.
- - The original e-process is not *truly* immune to "data peeking."
 - It appears that it is necessary to sacrifice some of the evidence (via adjusters).

1. Are there testing problems for which there is no powerful **e-process** in the data filtration, but

• If there exists a powerful e-process $(e_t)_{t>0}$ for \mathscr{P} w.r.t. some $\mathbb{G} \subseteq \mathbb{F}$, then we can e-lift it to \mathbb{F} .

• If $\operatorname{limsup}_{t\to\infty} e_t = \infty$ under some $\mathbb{Q} \setminus \mathscr{P}$, then $\operatorname{limsup}_{t\to\infty} A(e_t^*) = \infty$ under $\mathbb{Q} \setminus \mathscr{P}$ (for any

2. There appears to be an **unavoidable cost** to coarsening the filtration to obtain an e-process.

A Corollary on Coarsening the Filtration

<u>Corollary.</u> Let \mathscr{P} be a composite null and let \mathscr{Q} be a composite alternative. Suppose there exists a \mathscr{Q} -powerful* e-process for \mathscr{P} in a coarsened filtration \mathbb{G} of \mathbb{F} . **Then, there exists a** \mathscr{Q} -powerful e-process for \mathscr{P} in \mathbb{F} .

*An e-process for \mathscr{P} is \mathscr{Q} -powerful if, for any $Q \in \mathscr{Q} \setminus \mathscr{P}$, $\operatorname{limsup}_{t \to \infty} e_t = \infty$, Q-almost surely.

• Interestingly, this is NOT the case if "e-process" is replaced with "test martingale"!

Is adjusting the e-process the only way?

- (The e-process can be for any null w.r.t. any filtration.)
- Is the function necessarily an adjuster?

function that maps the running maximum e_{+}^{*} to some e_{+}' for each t.

Suppose you claim to have a function that, if I give you some e-process w.r.t. a coarse filtration, then the function can transform it into an e-process w.r.t. the data filtration.

Theorem (informal): The function is **necessarily** an adjuster, as long as it is an increasing

A game-theoretic definition of adjusters How can we make betting on the running maximum a "fair game"?

 An increasing function A is an adjuster if and only if, for every test supermartingale $(M_t)_{t>0}$ for some P, there exists a test supermartingale $(M'_t)_{t>0}$ for P s.t. A is an "adjuster for test supermartingales"

$$\mathsf{A}(\mathsf{M}^*_t) \leq \mathsf{M}'_t, \quad \forall t.$$

- Game-theoretically, adjusters allow betting running maximum of the gambler's wealth
 - A is an adjuster if and only if, in Protoco Skeptic has a betting strategy to ensure

$$\mathsf{A}(\mathscr{K}^*_{\mathsf{t}}) \leq \mathscr{K}'_{\mathsf{t}}.$$

a with the	Protocol 1 Competitive scepticism
	$\mathcal{K}_0 := 1 \text{ and } \mathcal{K}'_0 := 1$
h.	for $n = 1, 2,$ do
	Forecaster announces $\mathcal{E}_n \in \mathbf{E}$
1 Rival	Sceptic announces $f_n \in [0,\infty]^{\mathcal{X}}$ such that $\mathcal{E}_n(f_n) \leq \mathcal{K}_n$.
	Rival Sceptic announces $f'_n \in [0,\infty]^{\mathcal{X}}$ such that $\mathcal{E}_n(f'_n)$
e that	Reality announces $x_n \in \mathcal{X}$
	$\mathcal{K}_n := f_n(x_n) \text{ and } \mathcal{K}'_n := f'_n(x_n)$
	end for



A Characterization Theorem for Adjusters

(a) A is an adjuster, i.e., it satisfies $\int_{1}^{\infty} \frac{A(e)}{e^{2}} de \leq 1.$

(b) A is an "adjuster for test supermartingales" (previous slide).

 $\mathbb{F} \supseteq \mathbb{G}_{t} (A(\mathfrak{e}_{t}^{*}))_{t>0}$ is an e-process for \mathscr{P} w.r.t. \mathbb{F} .

- <u>**Theorem.</u>** Let A : $[1, \infty] \rightarrow [0, \infty]$ be an increasing function. The following are <u>equivalent</u>:</u>
- (c) A is an "adjuster for e-processes": for any e-process $(e_t)_{t>0}$ for some \mathcal{P} w.r.t. G, there exists another e-process $(e'_t)_{t>0}$ for \mathcal{P} w.r.t. G such that, for all t, $A(e^*_t) \leq e'_t$.
- (d) A is an "e-lifter": for any e-process $(e_t)_{t>0}$ for some \mathscr{P} w.r.t. \mathbb{G} , and any finer filtration
- (e) For any e-process $(e_t)_{t\geq 0}$ for some \mathscr{P} w.r.t. \mathbb{G} , $(A(e_t^*))_{t\geq 0}$ is an e-process for \mathscr{P} w.r.t. \mathbb{G} .

Takeaways

Takeaways

- E-processes constructed on a coarse filtration are **not** anytime-valid in the data filtration, so they cannot be combined seamlessly with other e-processes.
 - Examples: testing exchangeability; independence; comparing multi-step forecasters
- 2. P-processes can be lifted freely across filtrations and retain their anytime-validity.
- 3. For e-processes, we can use **adjusters** to achieve validity in the data filtration.
- 4. In a sense, any function that lifts the anytime-validity of an e-process **must be** an adjuster.
- 5. U-randomization can be applied after adjustment, but not before (w/o hurting validity).

Thank You

arXiv: https://arxiv.org/abs/2402.09698 **Code:** <u>https://github.com/yjchoe/CombiningEvidenceAcrossFiltrations</u>

Questions?



YJ <u>yjchoe.github.io</u>



Aaditya stat.cmu.edu/~aramdas



Appendix

Example: Comparing Multi-Step Sequential Forecasters

$$\Delta_t^{[k]} = \frac{1}{|I_t^{[k]}|} \sum_{i \in I_t^{[k]}} \mathbb{E} \left[S(p_i, y_{i+h-1}) - S(q_i, y_{i+h-1}) \mid \mathscr{F}_{i-1} \right], \quad \forall k \in [h].$$

- If h = 2, $\Delta_{+}^{[0]}/\Delta_{+}^{[1]}$ measures the average forecast score difference on even/odd days.
- under different coarsening of the filtration \mathbb{F} for each k (updates on every even/odd days).

into the data filtration F before combining them:

Suppose we compare two sequential forecasters with lag h using some scoring rule S w.r.t. $\mathbb{F} = (\mathscr{F}_t)_{t \ge 0}$:

• When testing for the null $\mathscr{H}_0^{[k]}: \Delta_t^{[k]} \leq 0, \forall t$, for each offset k, we need to construct an e-process $(\mathfrak{e}_t^{[k]})_{t\geq 0}$

```
Each (e_t^{[k]})_{t\geq 0} is an e-process for \mathscr{H}_0^{[k]}, but only w.r.t. the sub-filtration \mathbb{G}^{[k]} \subsetneq \mathbb{F}.
```

• To test for the combined null $\mathscr{H}_0: \Delta_t^{[k]} \leq 0, \forall t, \forall k$ (an intersection), we want to e-lift all h e-processes

$$A\left((e_t^{[k]})^*\right), \forall t.$$

Henzi & Ziegel (2022) Arnold et al. (2022) Choe & Ramdas (2023)



Example: Testing Independence

distribution factorizes:

$$\mathscr{H}_0: \mathsf{P}_{\mathsf{X}\mathsf{Y}} = \mathsf{P}_{\mathsf{X}} \times \mathsf{P}_{\mathsf{Y}}$$

- the data filtration F. Two known e-processes include:
- filtrations. So we should lift both of them before taking the average.

Given an i.i.d. stream of paired data $Z_t = (X_t, Y_t) \sim P_{XY}$, suppose we test if the joint

vs.
$$\mathscr{H}_1: \mathsf{P}_{\mathsf{X}\mathsf{Y}} \neq \mathsf{P}_{\mathsf{X}} \times \mathsf{P}_{\mathsf{Y}}$$
.

Similar to the exchangeability null, there exist no nontrivial test martingale adapted to

Pairwise betting (SR'23; PBKR'23; SR'24): adapted to the filtration w/ pairs of data.

Rank-based test martingale (HL'23): adapted to the filtration w/ rank stats of data.

In this case, BOTH e-processes are constructed w.r.t. their own, non-overlapping sub-

cf. Balasubramani & Ramdas (2016); Shekhar & Ramdas (2023); Podkopaev et al. (2023); Henzi & Law (2023)

Randomized adjusters for e-lifting

Motivation: Randomized Calibrators

- In the case of (non-sequential) **e-to-p calibration**, it is known that p = 1/e is the only admissible deterministic mapping.
- Recent papers show that there is the following "U-randomized" e-to-p calibrator can dominate the deterministic variant (almost surely):
 - $\tilde{p} = U/e$, for some $U \gtrsim Unif[0, 1]$.

Can we leverage this idea to develop a **U-randomized adjuster**?

cf. Ignatiadis et al. (2023); Ramdas & Manole (2023)



Strategy #1: Lift-then-randomize ("ltr") This is possible via UMI!

$$e_{\tau} \xrightarrow[(e-to-e)]{e-lifting} A(e_{\tau}^{*})$$

- Once we lift an e-process by adjustment, we have an e-process in F.
- τ_{i} due to "UMI" (uniformly-randomized Markov's inequality; Ramdas & Manole, 2023).
- more lenient than the p-lifted stopping rule of $p_{\tau} \leq \alpha$ (w/o adjustment).

for any $X \ge 0$ and $U \gtrsim Unif[0, 1]$, $\mathsf{P}\left(\mathsf{X} \geq \frac{\mathsf{U}}{\alpha}\right) \leq \alpha \cdot \mathbb{E}[\mathsf{X}].$

 $\begin{array}{c} U-\text{rand.} \\ \longrightarrow \\ (e-to-p) \end{array} \quad \tilde{p}_{\tau}^{\text{ltr}} = \frac{U}{A(e_{\tau}^{*})} \wedge 1 \end{array}$

• So the randomization strategy still works: \tilde{p}_{τ}^{ltr} is a valid p-value for any \mathbb{F} -stopping time

• Since we end up with a p-value, this is only practically useful if the stopping rule $\tilde{p}_{\tau}^{tr} \leq \alpha$ is



Strategy #2: Randomize-then-lift ("rtl") This does not guarantee anytime-validity (empirically, at least)

$$\mathfrak{e}_{\tau} \xrightarrow{\text{rand. } e-to-p}_{\mathsf{U}\gtrsim\mathsf{Unif}[0,1]} \qquad \tilde{\mathfrak{p}}_{\tau}^{\mathsf{rtl}} := \frac{\mathsf{U}}{\mathsf{e}_{\tau}} \wedge 1 \qquad \stackrel{\mathsf{p}-to-\mathsf{e}}{\longrightarrow} \qquad \tilde{\mathfrak{e}}_{\tau}^{\mathsf{rtl}} := \mathsf{C}\left(\frac{\mathsf{U}}{\mathsf{e}_{\tau}} \wedge 1\right)$$

- Once we add the external r.v. U, the resulting sequence is **not adapted** to \mathbb{G} !
- $\alpha \in (0, 1)$, we have, for any \mathbb{F} -stopping time τ and any $\mathsf{P} \in \mathscr{P}$,

 ≈ 0.065 for conformal mtg.)

• The U-lifting "lemma": If $(e_t)_{t>0}$ is an e-process for \mathscr{P} w.r.t. $\mathbb{G} \subseteq \mathbb{F}$, $U \gtrsim Unif[0, 1]$, and



End of Slides