

Confidence sequences

1 Definition of a confidence sequence

It has become standard practice for organizations with online presence to run large-scale randomized experiments, or A/B tests, to improve product performance and user experience. Such experiments are inherently sequential: visitors arrive in a stream and outcomes are typically observed quickly relative to the duration of the test. Results are often monitored continuously using inferential methods that assume a fixed sample, despite the well-known problem that such monitoring can inflate Type I error substantially (Armitage et al. 1969, Berman et al. 2018). Furthermore, most A/B tests are run with little formal planning and very fluid decision-making, as compared with clinical trials or industrial quality control, the traditional applications of sequential analysis.

In this mini, we present methods for deriving *confidence sequences* as a flexible tool for inference in sequential experiments (Darling & Robbins 1967a, Lai 1984, Jennison & Turnbull 1989). A confidence sequence is a sequence of confidence sets $(\text{CI}_t)_{t=1}^\infty$, typically intervals $\text{CI}_t = (L_t, U_t) \subseteq \mathbb{R}$, satisfying a uniform coverage guarantee: after observing the t^{th} unit, we calculate an updated confidence set CI_t for the unknown quantity of interest θ_t , with the coverage property

$$\mathcal{P}(\forall t \geq 1 : \theta_t \in \text{CI}_t) \geq 1 - \alpha. \quad (1)$$

With only a uniform lower bound (L_t) on $\theta_t \in \mathbb{R}$, i.e., if $U_t \equiv \infty$, we have a *lower confidence sequence*. Likewise, if $L_t \equiv -\infty$ we have an *upper confidence sequence* given by the uniform upper bound (U_t) . We will build upon the general framework for uniform exponential concentration introduced in the previous mini (Howard et al. 2018), which means our techniques apply to a wide variety of situations: scalar, matrix and Banach-space-valued observations, with possibly unbounded support; self-normalized bounds applicable to observations satisfying very weak moment or symmetry conditions; and continuous-time scalar martingales. Some bounds will yield closed-form confidence sequences, while others give a method for numerical computation of tighter intervals. Both methods allow for flexible control of the “shape” of the confidence sequence, that is, how the sequence of intervals shrink in width over time. As a simple example, given a sequence of observations from a 1-sub-Gaussian distribution whose mean we would like to track, we may choose any $\eta > 1$ and an increasing function $h : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}$ with $\sum_{k=0}^\infty 1/h(k) = 1$, to obtain a confidence sequence of the form

$$\frac{S_t}{t} \pm \frac{\eta^{1/4} + \eta^{-1/4}}{\sqrt{2}} \sqrt{\frac{\log h(\log_\eta t) + \log(2/\alpha)}{t}}. \quad (2)$$

Our theorems will generalize and sharpen related methods from Darling & Robbins (1967b, 1968), Jamieson et al. (2014), Kaufmann et al. (2014), Balsubramani (2014), Zhao et al. (2016). Our confidence sequences possess the following properties:

- (P1) **Non-asymptotic and nonparametric:** our confidence sequences offer provable coverage for all sample sizes, without exact distributional assumptions or asymptotic approximations.
- (P2) **Unbounded sample size:** our methods do not require a final sample size to be chosen ahead of time. They may be tuned for a planned sample size, but always permit additional sampling.
- (P3) **Arbitrary stopping rules:** we make no assumptions on the stopping rule used by an experimenter to decide when to end the experiment, or when to act on certain inferences.

These properties give us strong guarantees and broad applicability. An experimenter may always choose to gather more samples, and may stop at any time according to any rule, even one not formally defined, and the resulting inferential guarantees hold under the stated assumptions without any approximations. Of course, this flexibility comes with a cost: our intervals are wider than those that rely on asymptotics, and without assuming a rigid stopping rule, we cannot explicitly correct for selective bias introduced by adaptive stopping. The typical, fixed-sample confidence intervals derived from the central limit theorem do not satisfy any of these properties, and accommodating any one property necessitates wider intervals. It is remarkable that we can accommodate all three and incur a cost of less than doubling the interval width—the discrete mixture bound stays within a factor of two of the fixed-sample central limit theorem bounds over five orders of magnitude in time. Our work gives another example of gaining flexibility and robustness by “doubling” uncertainty estimates, an observation made recently in multiple testing by Katsevich & Ramdas (2018), and a theme more broadly explored by Meng (2018). It may seem that the definition (1) of a confidence sequence is stronger than necessary to achieve these properties, but as we show below, it is equivalent to a definition in terms of arbitrary, unbounded stopping times. It is therefore reasonable to say that any procedure satisfying these three properties will satisfy a guarantee similar to (1).

We will later demonstrate two applications in sequential estimation. First, under a randomization inference model in the Neyman-Rubin potential outcomes framework, we give a tight *empirical variance* confidence sequence for Bernoulli treatment assignment. This method sequentially estimates the variance of the underlying process and uses it to generate a valid confidence sequence, giving a non-asymptotic, sequential analogue of the *t*-test. Such a confidence sequence follows from a general empirical variance confidence sequences for bounded observations. Second, we give asymptotic and non-asymptotic iterated logarithm bounds for the operator norm of a matrix martingale and demonstrate their application to sequential covariance matrix estimation.

Lemma 1 Let $(A_t)_{t=1}^{\infty}$ be an adapted sequence of events in some filtered probability space and let $A_{\infty} := \limsup_{t \rightarrow \infty} A_t$. The following are equivalent:

1. $\mathcal{P}(\bigcup_{t=1}^{\infty} A_t) \leq \alpha$.
2. $\mathcal{P}(A_T) \leq \alpha$ for all random times T , possibly infinite and not necessarily stopping times.
3. $\mathcal{P}(A_{\tau}) \leq \alpha$ for all stopping times τ , possibly infinite.

Proof: The implication (a) \implies (b) follows from

$$A_T = \left(\bigcup_{t=1}^{\infty} A_t \cap \{T = t\} \right) \cup [A_{\infty} \cap \{T = \infty\}] \subseteq \bigcup_{t=1}^{\infty} A_t. \quad (3)$$

It is clear that (b) \implies (c). For (c) \implies (a), take $\tau = \inf\{t \in \mathcal{N} : A_t \text{ occurs}\}$, so that $A_{\tau} = \bigcup_{t=1}^{\infty} A_t$. ■

References

- Armitage, P., McPherson, C. K. & Rowe, B. C. (1969), ‘Repeated Significance Tests on Accumulating Data’, *Journal of the Royal Statistical Society. Series A (General)* **132**(2), 235–244.
- Balsubramani, A. (2014), ‘Sharp Finite-Time Iterated-Logarithm Martingale Concentration’, *arXiv:1405.2639 [cs, math, stat]*.
- Berman, R., Pekelis, L., Scott, A. & Van den Bulte, C. (2018), p-Hacking and False Discovery in A/B Testing, SSRN Scholarly Paper ID 3204791, Social Science Research Network, Rochester, NY.
- Darling, D. A. & Robbins, H. (1967a), ‘Confidence Sequences for Mean, Variance, and Median’, *Proceedings of the National Academy of Sciences* **58**(1), 66–68.
- Darling, D. A. & Robbins, H. (1967b), ‘Iterated Logarithm Inequalities’, *Proceedings of the National Academy of Sciences* **57**(5), 1188–1192.
- Darling, D. A. & Robbins, H. (1968), ‘Some Further Remarks on Inequalities for Sample Sums’, *Proceedings of the National Academy of Sciences* **60**(4), 1175–1182.
- Howard, S. R., Ramdas, A., McAuliffe, J. & Sekhon, J. (2018), ‘Exponential line-crossing inequalities’, *arXiv:1808.03204 [math]*.
- Jamieson, K., Malloy, M., Nowak, R. & Bubeck, S. (2014), lil’ UCB: An Optimal Exploration Algorithm for Multi-Armed Bandits, in ‘Proceedings of The 27th Conference on Learning Theory’, Vol. 35 of *Proceedings of Machine Learning Research*, pp. 423–439.

- Jennison, C. & Turnbull, B. W. (1989), ‘Interim Analyses: The Repeated Confidence Interval Approach’, *Journal of the Royal Statistical Society. Series B (Methodological)* **51**(3), 305–361.
- Katsevich, E. & Ramdas, A. (2018), ‘Towards “simultaneous selective inference”: post-hoc bounds on the false discovery proportion’, *arXiv:1803.06790 [math, stat]* .
- Kaufmann, E., Cappé, O. & Garivier, A. (2014), ‘On the Complexity of Best Arm Identification in Multi-Armed Bandit Models’, *arXiv:1407.4443 [cs, stat]* .
- Lai, T. L. (1984), ‘Incorporating scientific, ethical and economic considerations into the design of clinical trials in the pharmaceutical industry: a sequential approach’, *Communications in Statistics - Theory and Methods* **13**(19), 2355–2368.
- Meng, X.-L. (2018), ‘Double Your Variance, Dirtify Your Bayes, Devour Your Pufferfish, and Draw Your Kidstogram’.
- Zhao, S., Zhou, E., Sabharwal, A. & Ermon, S. (2016), Adaptive Concentration Inequalities for Sequential Decision Problems, *in* ‘30th Conference on Neural Information Processing Systems (NIPS 2016)’, Barcelona, Spain.