
Accurate Inference in Adaptive Linear Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Estimators computed from adaptively collected data do not behave like their non-
2 adaptive brethren. Rather, the sequential dependence of the collection policy
3 can lead to severe distributional biases that persist even in the infinite data limit.
4 We develop a general method – *W-decorrelation* – for transforming the bias of
5 adaptive linear regression estimators into variance. The method uses only coarse-
6 grained information about the data collection policy and does not need access to
7 propensity scores or exact knowledge of the policy. We bound the finite-sample bias
8 and variance of the *W*-estimator and develop asymptotically correct confidence
9 intervals based on a novel martingale central limit theorem. We then demonstrate
10 the empirical benefits of the generic *W*-decorrelation procedure in two different
11 adaptive data settings: the multi-armed bandit and the autoregressive time series.

12 1 Introduction

13 Consider a dataset of n sample points $(y_i, \mathbf{x}_i)_{i \leq n}$ where y_i represents an observed outcome and
14 $\mathbf{x}_i \in \mathbb{R}^p$ an associated vector of covariates. In the standard linear model, the outcomes and covariates
15 are related through a parameter β :

$$y_i = \langle \mathbf{x}_i, \beta \rangle + \varepsilon_i. \quad (1)$$

16 In this model, the ‘noise’ term ε_i represents inherent variation in the sample, or the variation that is
17 not captured in the model. Parametric models of the type (1) are a fundamental building block in
18 many machine learning problems. A common additional assumption is that the covariate vector \mathbf{x}_i
19 for a given datapoint i is independent of the other sample point outcomes $(y_j)_{j \neq i}$ and the inherent
20 variation $(\varepsilon_j)_{j \in [n]}$. This paper is motivated by experiments where the sample $(y_i, \mathbf{x}_i)_{i \leq n}$ is not
21 completely randomized but rather *adaptively* chosen. By adaptive, we mean that the choice of the data
22 point (y_i, \mathbf{x}_i) is guided from inferences on past data $(y_j, \mathbf{x}_j)_{j < i}$. Consider the following sequential
23 paradigms:

- 24 1. Multi-armed bandits: This class of sequential decision making problems captures the
25 classical ‘exploration versus exploitation’ tradeoff. At each time i , the experimenter chooses
26 an ‘action’ \mathbf{x}_i from a set of available actions \mathcal{X} and accrues a reward $R(y_i)$ where (y_i, \mathbf{x}_i)
27 follow the model (1). Here the experimenter must balance the conflicting goals of learning
28 about the underlying model (i.e., β) for better future rewards, while still accruing reward in
29 the current time step.
- 30 2. Active learning: Acquiring labels y_i is potentially costly, and the experimenter aims to
31 learn with as few outcomes as possible. At time i , based on prior data $(y_j, \mathbf{x}_j)_{j \leq i-1}$ the
32 experimenter chooses a new data point \mathbf{x}_i to label based on its value in learning.
- 33 3. Time series analysis: Here, the data points (y_i, \mathbf{x}_i) are naturally ordered in time, with
34 $(y_i)_{i \leq n}$ denoting a time series and the covariates \mathbf{x}_i include observations from the prior time
35 points.

36 Here, time induces a natural sequential dependence across the samples. In the first two instances, the
 37 actions or policy of the experimenter are responsible for creating such dependence. In the case of time
 38 series data, this dependence is endogenous and a consequence of the modeling. A common feature,
 39 however, is that the choice of the design or sequence $(\mathbf{x}_i)_{i \leq n}$ is typically not made for inference
 40 on the model after the data collection is completed. This does not, of course, imply that accurate
 41 estimates on the parameters β cannot be made from the data. Indeed, it is often the case that the
 42 sample is informative enough to extract consistent estimators of the underlying parameters. Indeed,
 43 this is often crucial to the success of the experimenter’s policy. For instance, ‘regret’ in sequential
 44 decision-making or risk in active learning are intimately connected with the accurate estimation of
 45 the underlying parameters [Castro and Nowak, 2008, Audibert and Bubeck, 2009, Bubeck et al.,
 46 2012, Rusmevichientong and Tsitsiklis, 2010] . Our motivation is the natural follow-up question of
 47 accurate *ex post* inference in the standard statistical sense:

48 Can adaptive data be used to compute accurate confidence regions and p -values?

49 As we will see, the key challenge is that even in the simple linear model of (1), the distribution of
 50 classical estimators can differ from the predicted central limit behavior of non-adaptive designs. In
 51 this context we make the following contributions:

- 52 • **Decorrelated estimators:** We present a general method to decorrelate arbitrary estimators
 53 $\hat{\beta}(\mathbf{y}, \mathbf{X}_n)$ constructed from the data. This construction admits a simple decomposition
 54 into a ‘bias’ and ‘variance’ term. In comparison with competing methods, like propensity
 55 weighting, our proposal requires little explicit information about the data-collection policy.
- 56 • **Bias and variance control:** Under a natural exploration condition on the data collection
 57 policy, we establish that the bias and variance can be controlled at nearly optimal levels. In
 58 the multi-armed bandit setting, we prove this under an especially weak averaged exploration
 59 condition.
- 60 • **Asymptotic normality and inference:** We establish a martingale central limit theorem
 61 (CLT) under a moment stability assumption. Applied to our decorrelated estimators, this
 62 allows us to construct confidence intervals and conduct hypothesis tests in the usual fashion.
- 63 • **Validation:** We demonstrate the usefulness of the decorrelating construction in two different
 64 scenarios: multi-armed bandits (MAB) and autoregressive (AR) time series. We observe
 65 that our decorrelated estimators retain expected central limit behavior in regimes where the
 66 standard estimators do not, thereby facilitating accurate inference.

67 The rest of the paper is organized with our main results in Section 2, discussion of related work in
 68 Section 3, and experiments in Section 4. An earlier version of this paper was published in ICML
 69 2018 (citation retracted). This version contains a new ‘limited information’ martingale central limit
 70 theorem, as well as new results on for the special case of multi-armed bandits.

71 2 Main results: W -decorrelation

72 We focus on the linear model and assume that the data pairs (y_i, \mathbf{x}_i) satisfy:

$$y_i = \langle \mathbf{x}_i, \beta \rangle + \varepsilon_i, \quad (2)$$

73 where ε_i are independent and identically distributed random variables with $\mathbb{E}\{\varepsilon_i\} = 0$, $\mathbb{E}\{\varepsilon_i^2\} = \sigma^2$
 74 and bounded third moment. We assume that the samples are ordered naturally in time and let $\{\mathcal{F}_i\}_{i \geq 0}$
 75 denote the filtration representing the sample. Formally, we let data points (y_i, \mathbf{x}_i) be adapted to this
 76 filtration, i.e. (y_i, \mathbf{x}_i) are measurable with respect to \mathcal{F}_j for all $j \geq i$.

77 Our goal in this paper is to use the available data to construct *ex post* confidence intervals and p -values
 78 for individual parameters, i.e. entries of β . A natural starting point is to consider is the standard least
 79 squares estimate:

$$\hat{\beta}_{\text{OLS}} = (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{y}_n,$$

80 where $\mathbf{X}_n = [\mathbf{x}_1^\top, \dots, \mathbf{x}_n^\top] \in \mathbb{R}^{n \times p}$ is the design matrix and $\mathbf{y}_n = [y_1, \dots, y_n] \in \mathbb{R}^n$. When data
 81 collection is non-adaptive, classical results imply that the standard least squares estimate $\hat{\beta}_{\text{OLS}}$
 82 is distributed asymptotically as $N(\beta, \sigma^2 (\mathbf{X}_n^\top \mathbf{X}_n)^{-1})$, where $N(\mu, \Sigma)$ denotes the Gaussian distribution
 83 with mean μ and covariance Σ . Lai and Wei [1982] extend these results to the current scenario:

84 **Theorem 1** (Theorems 1, 3 [Lai and Wei, 1982]). Let $\lambda_{\min}(n)$ ($\lambda_{\max}(n)$) denote the minimum (resp.
85 maximum) eigenvalue of $\mathbf{X}_n^\top \mathbf{X}_n$. Under the model (2), assume that (i) ε_i have finite third moment
86 and (ii) almost surely, $\lambda_{\min}(n) \rightarrow \infty$ with $\lambda_{\min} = \Omega(\log \lambda_{\max})$ and (iii) $\log \lambda_{\max} = o(n)$. Then
87 the following limits hold almost surely:

$$\|\widehat{\beta}_{\text{OLS}} - \beta\|_2^2 \leq C \frac{\sigma^2 p \log \lambda_{\max}}{\lambda_{\min}}$$

$$\left| \frac{1}{n\sigma^2} \|\mathbf{y}_n - \mathbf{X}_n \widehat{\beta}_{\text{OLS}}\|_2^2 - 1 \right| \leq C(p) \frac{1 + \log \lambda_{\max}}{n}.$$

88 Further assume the following stability condition: there exists a deterministic sequence of matrices
89 \mathbf{A}_n such that (iii) $\mathbf{A}_n^{-1} (\mathbf{X}_n^\top \mathbf{X}_n)^{1/2} \rightarrow \mathbf{I}_p$ and (iv) $\max_i \|\mathbf{A}_n^{-1} \mathbf{x}_i\|_2 \rightarrow 0$ in probability. Then,

$$(\mathbf{X}_n^\top \mathbf{X}_n)^{1/2} (\widehat{\beta}_{\text{OLS}} - \beta) \xrightarrow{d} \mathbf{N}(0, \sigma^2 \mathbf{I}_p).$$

90 At first blush, this allows to construct confidence regions in the usual way. More precisely, the result
91 implies that $\widehat{\sigma}^2 = \|\mathbf{y}_n - \mathbf{X}_n \widehat{\beta}_{\text{OLS}}\|_2^2 / n$ is a consistent estimate of the noise variance. Therefore, the
92 interval $[\widehat{\beta}_{\text{OLS},1} - 1.96\widehat{\sigma}(\mathbf{X}_n^\top \mathbf{X}_n)^{-1/2}_{11}, \widehat{\beta}_{\text{OLS},1} + 1.96\widehat{\sigma}(\mathbf{X}_n^\top \mathbf{X}_n)^{-1/2}_{11}]$ is a 95% two-sided confidence
93 interval for the first coordinate β_1 . Indeed, this result is sufficient for a variety of scenarios with
94 weak dependence across samples, such as when the (y_i, \mathbf{x}_i) form a Markov chain that mixes rapidly.
95 However, while the assumptions for consistency are minimal, the additional stability assumption
96 required for asymptotic normality poses some challenges. In particular:

- 97 1. The stability condition can provably fail to hold for scenarios where the dependence across
98 samples is non-negligible. This is not a weakness of Theorem 1: the CLT need not hold for
99 the OLS estimator [Lai and Wei, 1982, Lai and Siegmund, 1983].
- 100 2. The rate of convergence to the asymptotic CLT depends on the *quantitative rate* of the
101 stability condition. In other words, variability in the inverse covariance $\mathbf{X}_n^\top \mathbf{X}_n$ can cause
102 deviations from normality of OLS estimator [Dvoretzky, 1972]. In finite samples, this can
103 manifest itself in the bias of the OLS estimator as well as in higher moments.

104 An example of this phenomenon is the standard multi-armed bandit problem [Lai and Robbins,
105 1985]. At each time point $i \leq n$, the experimenter (or data collecting policy) chooses an arm
106 $k \in \{1, 2, \dots, p\}$ and observes a reward y_i with mean β_k . With $\beta \in \mathbb{R}^p$ denoting the mean rewards,
107 this falls within the scope of model (2), where the vector \mathbf{x}_i takes the value \mathbf{e}_k (the k^{th} basis vector),
108 if the k^{th} arm or option is chosen at time i .¹ Other stochastic bandit problems with covariates such as
109 contextual or linear bandits [Rusmevichientong and Tsitsiklis, 2010, Li et al., 2010, Deshpande and
110 Montanari, 2012] can also be incorporated fairly naturally into our framework. For the purposes of
111 this paper, however, we restrict ourselves to the simple case of multi-armed bandits without covariates.
112 In this setting, ordinary least squares estimates correspond to computing sample means for each arm.
113 The stability condition of Theorem 1 requires that $N_k(n)$, the number of times a specific arm $k \in [p]$
114 is sampled is asymptotically deterministic as n grows large. This is true for certain regret-optimal
115 algorithms [Russo, 2016, Garivier and Cappé, 2011]. Indeed, for such algorithms, as the sample
116 size n grows large, the suboptimal arm is sampled $N_k(n) \sim C_k(\beta) \log n$ for a constant $C_k(\beta)$ that
117 depends on β and the distribution of noise ε_i . However, in finite samples, the dependence on $C_k(\beta)$
118 and the slow convergence rate of $(\log n)^{-1/2}$ lead to significant deviation from the expected central
119 limit behavior.

120 Villar et al. [2015] studied a variety of multi-armed bandit algorithms in the context of clinical trials.
121 They empirically demonstrate that sample mean estimates from data collected using many standard
122 multi-armed bandit algorithms are biased. Recently, Nie et al. [2017] proved that this bias is negative
123 for Thompson sampling and UCB. The presence of bias in sample means demonstrates that standard
124 methods for inference, as advocated by Theorem 1, can be misleading when the same data is now
125 used for inference. As a pertinent example, testing the hypotheses “the mean reward of arm 1 exceeds
126 that of 2” based on classical theory can be significantly affected by adaptive data collection.

127 The papers [Villar et al., 2015, Nie et al., 2017] focus on the finite sample effect of the data collection
128 policy on the bias and suggest methods to reduce the bias. It is not hard to find examples where

¹Strictly speaking, the model (2) assumes that the errors have the same variance, which need not be true for the multi-armed bandit as discussed. We focus on the homoscedastic case where the errors have the same variance in this paper.

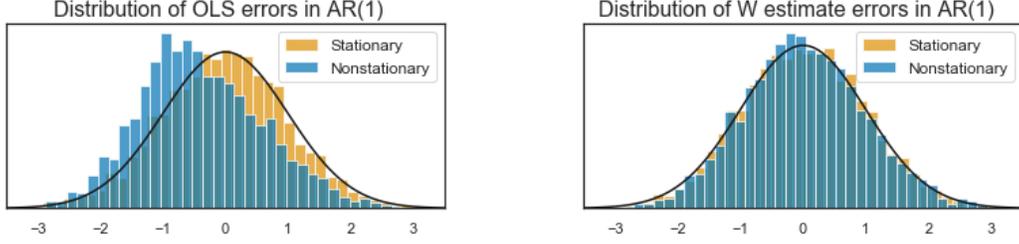


Figure 1: The distribution of normalized errors for (left) the OLS estimator for stationary and (nearly) nonstationary AR(1) time series and (right) error distribution for both models after decorrelation.

129 higher moments or tails of the distribution can be influenced by the data collecting policy. A simple,
 130 yet striking, example is the standard autoregressive model (AR) for time series data. In its simplest
 131 form, the AR model has one covariate, i.e. $p = 1$ with $\mathbf{x}_i = y_{i-1}$. In this case:

$$y_i = \beta y_{i-1} + \varepsilon_i.$$

132 Here the least squares estimate is given by $\hat{\beta}_{\text{OLS}} = \sum_{i \leq n-1} y_{i+1} y_i / \sum_{i \leq n-1} y_{i-1}^2$. When $|\beta|$ is
 133 bounded away from 1, the series is asymptotically stationary and the OLS estimate has Gaussian tails.
 134 On the other hand, when $\beta - 1$ is on the order of $1/n$ the limiting distribution of the least squares
 135 estimate is non-Gaussian and dependent on the gap $\beta - 1$ (cf. Chan and Wei [1987]). A histogram
 136 for the normalized OLS errors in two cases: (i) stationary with $\beta = 0.02$ and (ii) nonstationary with
 137 $\beta = 1.0$ is shown on the left in Figure 1. The OLS estimate yields clearly non-Gaussian errors when
 138 nonstationary, i.e. when β is close to 1.

139 On the other hand, *using the same data* our decorrelating procedure is able to obtain estimates
 140 admitting Gaussian limit distributions, as evidenced in the right panel of Figure 1. We show a similar
 141 phenomenon in the MAB setting where our decorrelating procedure corrects for the unstable behavior
 142 of the OLS estimator (see Section 4 for details on the empirics). Delegating discussion of further
 143 related work to 3, we now describe this procedure and its motivation.

144 2.1 Removing the effects of adaptivity

145 We propose to decorrelate the OLS estimator by constructing:

$$\hat{\beta}^d = \hat{\beta}_{\text{OLS}} + \mathbf{W}_n (y - \mathbf{X}_n \hat{\beta}_{\text{OLS}}),$$

146 for a specific choice of a ‘decorrelating’ or ‘whitening’ matrix $\mathbf{W}_n \in \mathbb{R}^{p \times n}$. This is inspired by the
 147 high-dimensional linear regression debiasing constructions of Zhang and Zhang [2014], Javanmard
 148 and Montanari [2014b,a], Van de Geer et al. [2014]. As we will see, this construction is useful also in
 149 the present regime where we keep p fixed and $n \gtrsim p$. By rearranging:

$$\begin{aligned} \hat{\beta}^d - \beta &= (\mathbf{I}_p - \mathbf{W}_n \mathbf{X}_n) (\hat{\beta}_{\text{OLS}} - \beta) + \mathbf{W}_n \varepsilon_n \\ &\equiv \mathbf{b} + \mathbf{v}. \end{aligned}$$

150 We interpret \mathbf{b} as a ‘bias’ and \mathbf{v} as a ‘variance’. This is based on the following critical constraint on
 151 the construction of the whitening matrix \mathbf{W}_n :

152 **Definition 1** (Well-adaptedness of \mathbf{W}_n). *Without loss of generality, we assume that ε_i are adapted*
 153 *to \mathcal{F}_i . Let $\mathcal{G}_i \subset \mathcal{F}_i$ be a filtration such that \mathbf{x}_i are adapted w.r.t. \mathcal{G}_i and ε_i is independent of \mathcal{G}_i .*
 154 *We say that \mathbf{W}_n is well-adapted if the columns of \mathbf{W}_n are adapted to \mathcal{G}_i , i.e. the i^{th} column \mathbf{w}_i is*
 155 *measurable with respect to \mathcal{G}_i .*

156 With this in hand, we have the following simple lemma.

157 **Lemma 2.** *Assume \mathbf{W}_n is well-adapted. Then:*

$$\begin{aligned} \|\beta - \mathbb{E}\{\hat{\beta}^d\}\|_2 &\leq \mathbb{E}\{\|\mathbf{b}\|_2\}, \\ \text{Var}(\mathbf{v}) &= \sigma^2 \mathbb{E}\{\mathbf{W}_n \mathbf{W}_n^\top\}. \end{aligned}$$

158 A concrete proposal is to trade-off the bias, controlled by the size of $\mathbf{I}_p - \mathbf{W}_n \mathbf{X}_n$, with the the
 159 variance which appears through $\mathbf{W}_n \mathbf{W}_n^\top$. This leads to the following optimization problem:

$$\mathbf{W}_n = \arg \min_{\mathbf{W}} \|\mathbf{I}_p - \mathbf{W} \mathbf{X}_n\|_F^2 + \lambda \text{Tr}(\mathbf{W} \mathbf{W}^\top).$$

160 Solving the above in closed form yields ridge estimators for β , and by continuity, also the standard
 161 least squares estimator. Departing from Zhang and Zhang [2014], Javanmard and Montanari [2014a],
 162 we solve the above in an *online* fashion in order to obtain a well-adapted \mathbf{W}_n . We define, $\mathbf{W}_0 = 0$,
 163 $\mathbf{X}_0 = 0$, and recursively $\mathbf{W}_n = [\mathbf{W}_{n-1} \mathbf{w}_n]$ for

$$\mathbf{w}_n = \arg \min_{\mathbf{w} \in \mathbb{R}^p} \|\mathbf{I} - \mathbf{W}_{n-1} \mathbf{X}_{n-1} - \mathbf{w} \mathbf{x}_n^\top\|_F^2 + \lambda \|\mathbf{w}\|_2^2.$$

164 As in the case of the offline optimization, we may obtain closed form formulae for the columns \mathbf{w}_i
 165 (see Algorithm 1). The method as specified requires $O(np^2)$ additional computational overhead,
 166 which is typically minimal compared to computing $\hat{\beta}_{\text{OLS}}$ or a regularized version like the ridge or
 167 lasso estimate. We refer to $\hat{\beta}^d$ as a *W-estimate* or a *W-decorrelated estimate*.

168 2.2 Bias and variance

169 We now examine the bias and variance control for $\hat{\beta}^d$. We first begin with a general bound for the
 170 variance:

171 **Theorem 3** (Variance control). *For any $\lambda \geq 1$ set non-adaptively, we have that*

$$\text{Tr}\{\text{Var}(\mathbf{v})\} \leq \frac{\sigma^2}{\lambda} (p - \mathbb{E}\{\|\mathbf{I}_p - \mathbf{W}_n \mathbf{X}_n\|_F^2\}).$$

172 *In particular, $\text{Tr}\{\text{Var}(\mathbf{v})\} \leq \sigma^2 p / \lambda$. Further, if $\|\mathbf{x}_i\|_2^2 \leq C$ for all i :*

$$\text{Tr}\{\text{Var}(\mathbf{v})\} \asymp \frac{\sigma^2}{\lambda} (p - \mathbb{E}\{\|\mathbf{I}_p - \mathbf{W}_n \mathbf{X}_n\|_F^2\}).$$

173 This theorem suggests that one must set λ as large as possible to minimize the variance. While this
 174 is accurate, one must take into account the bias of $\hat{\beta}^d$ and its dependence on the regularization λ .
 175 Indeed, for large λ , one would expect that $\mathbf{I}_p - \mathbf{W}_n \mathbf{X}_n \approx \mathbf{I}_p$, which would not help control the bias.
 176 In general, one would hope to set λ , thereby determining $\hat{\beta}^d$, at a level where its bias is negligible in
 177 comparison to the variance. The following theorem formalizes this:

178 **Theorem 4** (Variance dominates MSE). *Recall that the matrix \mathbf{W}_n is a function of λ . Suppose that*
 179 *there exists a deterministic sequence $\lambda(n)$ such that:*

$$\begin{aligned} \mathbb{E}\{\|\mathbf{I}_p - \mathbf{W}_n \mathbf{X}_n\|_{\text{op}}^2\} &= o(1/\log n), \\ \mathbb{P}\{\lambda_{\min}(\mathbf{X}_n^\top \mathbf{X}_n) \leq \lambda(n) \log \log n\} &\leq 1/n, \end{aligned}$$

180 *Then we have*

$$\frac{\|\mathbb{E}\{\mathbf{b}\}\|_2^2}{\text{Tr}\{\text{Var}(\mathbf{v})\}} = o(1).$$

181 The conditions of Theorem 4, in particular the bias condition on $\mathbf{I}_p - \mathbf{W}_n \mathbf{X}_n$ are quite general. In
 182 the following proposition, we verify some sufficient conditions under which the premise of Theorem
 183 4 hold.

184 **Proposition 5.** *Either of the following conditions suffices for the requirements of Theorem 4.*

185 1. *The data collection policy satisfies for some sequence $\mu_n(i)$ and for all $\lambda \geq 1$:*

$$\begin{aligned} \mathbb{E}\left\{\frac{\mathbf{x}_i \mathbf{x}_i^\top}{\lambda + \|\mathbf{x}_i\|_2^2} \middle| \mathcal{G}_{i-1}\right\} &\succcurlyeq \frac{\mu_n(i)}{\lambda} \mathbf{I}_p, \\ \sum_i \mu_n(i) &\equiv n \bar{\mu}_n \geq K \sqrt{n}, \end{aligned} \quad (3)$$

186 *for a large enough constant K . Here we keep $\lambda(n) \asymp n \bar{\mu}_n / \log(p \log n)$.*

187 2. *The matrices $(\mathbf{x}_i \mathbf{x}_i^\top)_{i \leq n}$ commute and $\lambda(n) \log \log n$ is (at most) the $1/n^{\text{th}}$ percentile of*
 188 *$\lambda_{\min}(\mathbf{X}_n^\top \mathbf{X}_n)$.*

Algorithm 1: W -Decorrelation Method

Input: sample $(y_i, \mathbf{x}_i)_{i \leq n}$, regularization λ , unit vector $\mathbf{v} \in \mathbb{R}^p$, confidence level $\alpha \in (0, 1)$, noise estimate $\hat{\sigma}^2$.

Compute: $\hat{\beta}_{\text{OLS}} = (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{y}_n$.

Setting $\mathbf{W}_0 = 0$, compute $\mathbf{W}_i = [\mathbf{W}_{i-1} \mathbf{w}_i]$ with $\mathbf{w}_i = (\mathbf{I}_p - \mathbf{W}_{i-1} \mathbf{X}_i^\top) \mathbf{x}_i / (\lambda + \|\mathbf{x}_i\|_2^2)$, for $i = 1, 2, \dots, n$.

Compute $\hat{\beta}^d = \hat{\beta}_{\text{OLS}} + \mathbf{W}_n (y - \mathbf{X}_n \hat{\beta}_{\text{OLS}})$ and $\hat{\sigma}(\mathbf{v}) = \hat{\sigma} \langle \mathbf{v}, \mathbf{W}_n \mathbf{W}_n^\top \mathbf{v} \rangle^{1/2}$

Output: decorrelated estimate $\hat{\beta}^d$ and CI interval

$$I(\mathbf{v}, \alpha) = [\langle \mathbf{v}, \hat{\beta}^d \rangle - \hat{\sigma}(\mathbf{v}) \Phi^{-1}(1 - \alpha), \langle \mathbf{v}, \hat{\beta}^d \rangle + \hat{\sigma}(\mathbf{v}) \Phi^{-1}(1 - \alpha)].$$

189 It is useful to consider the intuition for the sufficient conditions given in Proposition 5. By Lemma 2,
 190 note that the bias is controlled by $\|\mathbf{I} - \mathbf{W}_n \mathbf{X}_n\|_{\text{op}}$, which increases with λ . Consider a case in which
 191 the samples \mathbf{x}_i lie in a strict subspace of \mathbb{R}^p . In this case, controlling the bias uniformly over $\beta \in \mathbb{R}^p$
 192 is now impossible regardless of the choice of \mathbf{W}_n . For example, in a multi-armed bandit problem,
 193 if the policy does not sample a specific arm, there is no information available about the reward
 194 distribution of that arm. Proposition 5 the intuition that the data collecting policy should explore
 195 the full parameter space. For multi-armed bandits, policies such as epsilon-greedy and Thompson
 196 sampling satisfy this assumption with appropriate $\mu_n(i)$.

197 Given sufficient exploration, Proposition 5 recommends a reasonable value to set for the regularization
 198 parameter. In particular setting λ to a value such that $\lambda \leq \lambda_{\min} / \log \log n$ occurs with high probability
 199 suffices to ensure that the W -decorrelated estimate is approximately unbiased. Correspondingly, the
 200 MSE (or equivalently variance) of the W -decorrelated estimate need not be smaller than that of the
 201 original OLS estimate. Indeed the variance scales as $1/\lambda$, which exceeds with high probability the
 202 $1/\lambda_{\min}$ scaling for the MSE. This is the cost paid for debiasing OLS estimate.

203 Before we move to the inference results, note that the procedure requires only access to high
 204 probability lower bounds on λ_{\min} , which intuitively quantifies the exploration of the data collection
 205 policy. In comparison with methods such as propensity score weighting or conditional likelihood
 206 optimization, this represents rather coarse information about the data collection process. In particular,
 207 given access to propensity scores or conditional likelihoods one can simulate the process to extract
 208 appropriate values for the regularization $\lambda(n)$. This is the approach we take in the experiments of
 209 Section 4. Moreover, propensity scores or conditional likelihoods are ineffective when data collection
 210 policies make adaptive decisions that are deterministic given the history. A important example is that
 211 of UCB algorithms for bandits, which make deterministic choices of arms.

212 2.3 A central limit theorem and confidence intervals

213 Our final result is a central limit theorem that provides an alternative to the stability condition of
 214 Theorem 1 and standard martingale CLTs. Standard martingale CLTs [see, e.g., Lai and Wei, 1982,
 215 Dvoretzky, 1972] require convergence of $\sum_i \mathbf{w}_i \mathbf{w}_i^\top / n$ to a constant, but this convergence condition
 216 is violated in many examples of interest, including the AR examples in Section 4.

217 Let $(X_{i,n}, \mathcal{F}_{i,n}, 1 \leq i \leq n)$ be a martingale difference array, with the associated sum process
 218 $S_n = \sum_{i \leq n} X_{i,n}$ and covariance process $V_n = \sum_{i \leq n} \mathbb{E}\{X_{i,n}^2 | \mathcal{F}_{i-1,n}\}$.

219 **Assumption 1.** 1. *Moments are stable: for $a = 1, 2$, the following limit holds*

$$\lim_{n \rightarrow \infty} \mathbb{E} \left\{ \sum_{i \leq n} V_n^{-a/2} \left| \mathbb{E}\{X_{i,n}^a | \mathcal{F}_{i-1,n}, V_n\} - \mathbb{E}\{X_{i,n}^a | \mathcal{F}_{i-1,n}\} \right| \right\} = 0$$

220 2. *Martingale differences are small:*

$$\lim_{n \rightarrow \infty} \sum_{i \leq n} \mathbb{E} \left\{ \frac{|X_{i,n}|^3}{V_n^{3/2}} \right\} = 0,$$

$$\lim_{n \rightarrow \infty} \frac{\max_{i \leq n} \mathbb{E}\{X_{i,n}^2 | \mathcal{F}_{i-1,n}\}}{V_n} = 0 \text{ in probability.}$$

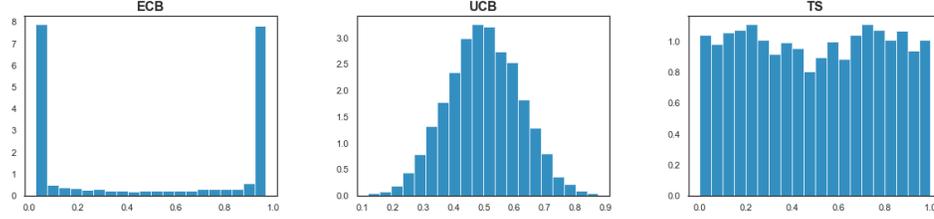


Figure 2: Histograms of the distribution of $N_1(n)/n$, the fraction of times arm 1 is picked under ε -greedy, UCB and Thompson sampling. The bandit problem has $p = 2$ arms which have i.i.d. $\text{Unif}([-0.7, 1.3])$ rewards and a time horizon of $n = 1000$. The distribution is plotted over 4000 Monte Carlo iterations.

221 **Theorem 6** (Martingale CLT). *Under Assumption 1, the rescaled process satisfies*
 222 $S_n/\sqrt{V_n} \xrightarrow{d} \mathcal{N}(0, 1)$, i.e. the following holds for any bounded, continuous test function $\varphi : \mathbb{R} \rightarrow \mathbb{R}$:

$$\lim_{n \rightarrow \infty} \mathbb{E}\{\varphi(S_n/\sqrt{V_n})\} = \mathbb{E}\{\varphi(\xi)\},$$

223 where $\xi \sim \mathcal{N}(0, 1)$.

224 The first part of Assumption 1 is an alternate form of stability. It controls the dependence of the
 225 conditional covariance of S_n on the first two conditional moments of the martingale increments $X_{i,n}$.
 226 In words, it states that the knowledge of the conditional covariance $\sum_i \mathbb{E}\{X_{i,n}^2 | \mathcal{F}_{i-1,n}\}$ does not
 227 change the first two moments of increments $X_{i,n}$ by an appreciable amount².

228 With a CLT in hand, one can now assign confidence intervals in the standard fashion, based on the
 229 assumption that the bias is negligible. For instance, we have result on two-sided confidence intervals.

230 **Proposition 7.** *Fix any $\alpha > 0$. Suppose that the data collection process satisfies the assumptions*
 231 *of Theorems 4 and 6. Set $\lambda = \lambda(n)$ as in Theorem 4, and let $\hat{\sigma}$ be a consistent estimate of σ as in*
 232 *Theorem 1. Define $\mathbf{Q} = \hat{\sigma}^2 \mathbf{W}_n \mathbf{W}_n^\top$ and the interval $I(a, \alpha) = [\hat{\beta}_a^d - \sqrt{Q_{aa}} \Phi^{-1}(1 - \alpha/2), \hat{\beta}_a^d +$
 233 $\sqrt{Q_{aa}} \Phi^{-1}(1 - \alpha/2)$. Then*

$$\limsup_{n \rightarrow \infty} \mathbb{P}\{\beta_a \notin I(a, \alpha)\} \leq \alpha.$$

234 2.4 Stability for multi-armed bandits

235 Limited information central limit theorems such as Theorem 6 (or [Hall and Heyde, 2014, Theorem
 236 3.4]), while providing insight into the problem of determining asymptotics, have assumptions that are
 237 often difficult to check in practice. Therefore, sufficient conditions such as the stability assumed in
 238 Theorem 1 are often preferred while analyzing the asymptotic behavior of martingales. In this section
 239 we circumvent this problem by proving the standard version of stability (as assumed in Theorem 1)
 240 for \mathbf{W} -estimates, assuming the matrices $\mathbf{x}_i \mathbf{x}_i^\top$ commute. While this is not a complete resolution to
 241 the problems posed by limited information martingale CLT's, it applies to important special cases
 242 like multi-armed bandits.

243 Recall that the stability assumed in Theorem 1 requires a non-random sequence of matrices \mathbf{A}_n so
 244 that

$$\mathbf{A}_n^{-1} \mathbf{X}_n \mathbf{X}_n^\top \xrightarrow{p} \mathbf{I}_p$$

245 When the vectors \mathbf{x}_i take values among $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$, a set of orthogonal vectors, we have

$$\begin{aligned} \mathbf{X}_n \mathbf{X}_n^\top &= \sum_i \mathbf{x}_i \mathbf{x}_i^\top \\ &= \sum_{a=1}^p \mathbf{v}_a \mathbf{v}_a^\top \sum_i \mathbb{I}(\text{arm } a \text{ chosen at time } i), \\ &= \sum_{a=1}^p \mathbf{v}_a \mathbf{v}_a^\top N_a(n), \end{aligned}$$

²See Hall and Heyde [2014], Theorem 3.4 for an example of a martingale central limit theorem in this flavor.

246 where we define $N_a(n) = \sum_{i=1}^n \mathbb{I}(\mathbf{x}_i = \mathbf{v}_a)$. Therefore, if there existed \mathbf{A}_n so that the stability
 247 condition held, then we would have, for each a , that $N_a(n)\langle \mathbf{v}_a, \mathbf{A}_n^{-1} \mathbf{v}_a \rangle \rightarrow 1$ in probability.

248 We test this assumption in a simple, but illuminating setting: a multi-armed bandit problem with
 249 $p = 2$ arms that are *statistically identical*: they each yield i.i.d. $\text{Unif}([-0.7, 1.3])$ rewards. We run
 250 ε -greedy (with a fixed value $\varepsilon = 0.1$), Thompson sampling and a variant of UCB for a time horizon of
 251 $n = 1000$ for 4000 Monte Carlo iterations. The resulting histograms of the fraction $N_1(n)/n$ of times
 252 arm 1 was picked by each of the three policies is given in Figure 2. Since the arms are statistically
 253 identical, the algorithm behavior is exchangeable with respect to switching the arm labels, viz.
 254 switching arm 1 for arm 2. In particular, the distribution of $N_1(n)$ and $N_2(n)$ is identical, for a given
 255 policy. Combining this with $N_1(n) + N_2(n) = n$, we have that $\mathbb{E}\{N_1(n)\} = \mathbb{E}\{N_2(n)\} = n/2$.
 256 Therefore, if stability a la Theorem 1 held, this would imply that the distribution of fraction $N_1(n)/n$
 257 would be close to a Dirac delta at $1/2$. However, we see that for all the three policies UCB, Thompson
 258 sampling and ε -greedy, this is not the case. Indeed, $N_1(n)/n$ has significant variance about $1/2$
 259 for all the policies; to wit, the ε -greedy indeed shows a sharp bimodal behavior. Consequently, the
 260 stability condition required by Theorem 1 *fails to hold* quite dramatically in this simple setting. As
 261 we observe in Section 4, this affects significantly the limiting distribution of the sample means, which
 262 have non-trivial bias and poor coverage of nominal confidence intervals.

263 In the following, we will prove that \mathbf{W} -estimates are indeed stable in the sense of Theorem 1, given a
 264 judicious choice of $\lambda = \lambda(n)$. Suppose that for each time i , $\mathbf{x}_i \in \{\mathbf{v}_1, \dots, \mathbf{v}_p\}$ the latter being a set
 265 of orthogonal (not necessarily unit normed) vectors \mathbf{v}_a . We also define $N_a(i) = \sum_{j \leq i} \mathbb{I}(\mathbf{x}_j = \mathbf{v}_a)$.
 266 The following proposition shows that when $\lambda = \lambda(n)$ is set appropriately, the \mathbf{W} -estimate is stable.

267 **Proposition 8.** *Suppose that the sequence $\lambda = \lambda(n)$ satisfies (i) $\lambda(n)/\lambda_{\min}(\mathbf{X}_n \mathbf{X}_n^\top) \rightarrow 0$ in
 268 probability and (ii) $\lambda(n) \rightarrow \infty$. Then the following holds:*

$$\lambda(n) \mathbf{W}_n \mathbf{W}_n^\top \xrightarrow{L_1} \frac{\mathbf{I}_p}{2}.$$

269 Along with Theorem 4 and Proposition 5, this immediately yields a simple corollary on the distribution
 270 of \mathbf{W} -estimates in the commutative setting. The key advantage here is that we are able to circumvent
 271 the assumptions of the limited information central limit Theorem 6.

272 **Corollary 9.** *Suppose that \mathbf{x}_i take values in $\{\mathbf{v}_1, \dots, \mathbf{v}_p\}$, a set of orthogonal vectors. Let $\hat{\sigma}^2$ be an
 273 estimate of the variance σ^2 as obtained from Theorem 1 and $\hat{\beta}^d$ be the \mathbf{W} -estimate obtained using
 274 $\lambda = \lambda(n)$ so that $\lambda(n) \log \log(n) \mathbb{E}\{\lambda_{\min}^{-1}(\mathbf{X}_n^\top \mathbf{X}_n)\} \rightarrow 0$. Then, with $\xi \sim \mathbf{N}(0, \mathbf{I}_p)$ and any Borel
 275 set $A \subseteq \mathbb{R}^p$:*

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{(\hat{\sigma}^2 \lambda(n) \mathbf{W}_n \mathbf{W}_n^\top)^{-1/2} (\hat{\beta}^d - \beta) \in A\right\} = \mathbb{P}\{\xi \in A\}.$$

276 3 Related work

277 There is extensive work in statistics and econometrics on stochastic regression models [Wei, 1985, Lai,
 278 1994, Chen et al., 1999, Heyde, 2008] and non-stationary time series [Shumway and Stoffer, 2006,
 279 Enders, 2008, Phillips and Perron, 1988]. This line of work is analogous to Theorem 1 or restricted
 280 to specific time series models. We instead focus on literature from sequential decision-making, policy
 281 learning and causal inference that closely resembles our work in terms of goals, techniques and
 282 applicability.

283 The seminal work of Lai and Robbins [Robbins, 1985, Lai and Robbins, 1985] has spurred a
 284 vast literature on multi-armed bandit problems and sequential experiments that propose allocation
 285 algorithms based on confidence bounds (see Bubeck et al. [2012] and references therein). A variety
 286 of confidence bounds and corresponding rules have been proposed [Auer, 2002, Dani et al., 2008,
 287 Rusmevichientong and Tsitsiklis, 2010, Abbasi-Yadkori et al., 2011, Jamieson et al., 2014] based
 288 on martingale concentration and the law of iterated logarithm. While these results can certainly be
 289 used to compute valid confidence intervals, they are conservative for a few reasons. First, they do not
 290 explicitly account for bias in OLS estimates and, correspondingly, must be wider to account for it.
 291 Second, obtaining optimal constants in the concentration inequalities can require sophisticated tools
 292 even for non-adaptive data [Ledoux, 1996, 2005]. This is evidenced in all of our experiments which
 293 show that concentration inequalities yield valid, but conservative intervals.

294 A closely-related line of work is that of learning from logged data [Li et al., 2011, Dudík et al., 2011,
 295 Swaminathan and Joachims, 2015] and policy learning [Athey and Wager, 2017, Kallus, 2017]. The
 296 focus here is efficiently estimating the reward (or value) of a certain test policy using data collected
 297 from a different policy. For linear models, this reduces to accurate prediction which is directly related
 298 to the estimation error on the parameters β . While our work shares some features, we focus on
 299 unbiased estimation of the parameters and obtaining accurate confidence intervals for linear functions
 300 of the parameters. Some of the work on learning from logged data also builds on propensity scores
 301 and their estimation [Imbens, 2000, Lunceford and Davidian, 2004].

302 Villar et al. [2015] empirically demonstrate the presence of bias for a number of multi-armed bandit
 303 algorithms. Recent work by Dimakopoulou et al. [2017] also shows a similar effect in contextual
 304 bandits. Along with a result on the sign of the bias, Nie et al. [2017] also propose conditional
 305 likelihood optimization methods to estimate parameters of the linear model. Through the lens
 306 of selective inference, they also propose methods to randomize the data collection process that
 307 simultaneously lower bias and reduce the MSE. Their techniques rely on considerable information
 308 about (and control over) the data generating process, in particular the probabilities of choosing a
 309 specific action at each point in the data selection. This can be viewed as lying on the opposite end of
 310 the spectrum from our work, which attempts to use only the data at hand, along with coarse aggregate
 311 information on the exploration inherent in the data generating process. It is an interesting, and open,
 312 direction to consider approaches that can combine the strengths of our approach and that of Nie et al.
 313 [2017].

314 4 Experiments

315 In this section we empirically validate the decorrelated estimators in two scenarios that involve
 316 sequential dependence in covariates. Our first scenario is a simple experiment of multi-armed bandits
 317 while the second scenario is autoregressive time series data. In these cases, we compare the empirical
 318 coverage and typical widths of confidence intervals for parameters obtained via three methods: (i)
 319 classical OLS theory, (ii) concentration inequalities and (iii) decorrelated estimates.

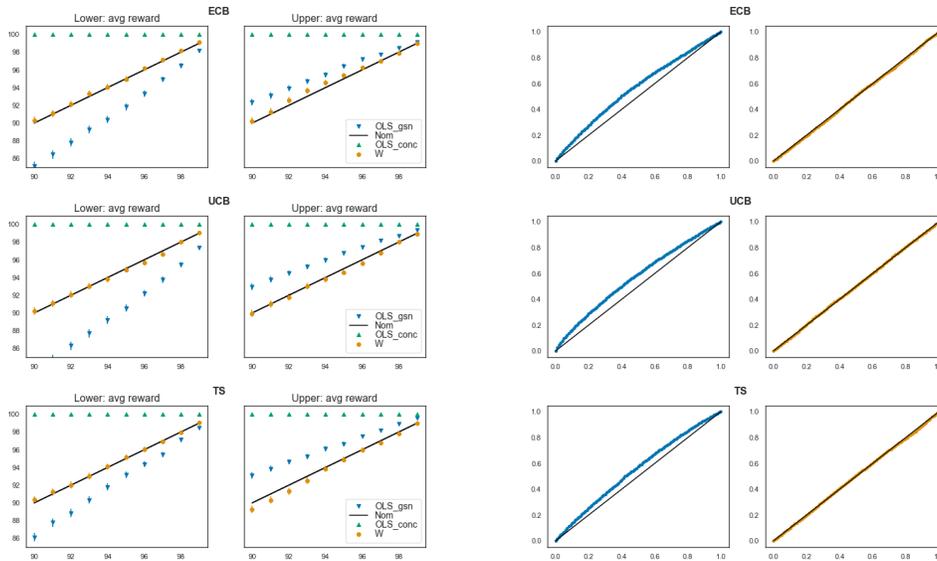


Figure 3: Multi-armed bandit results. Left: One-sided confidence region coverage for OLS and decorrelated \mathbf{W} -decorrelated estimates of the average reward $0.5\beta_1 + 0.5\beta_2$. Right: Probability (PP) plots for the OLS and \mathbf{W} -decorrelated estimate errors of the average reward.

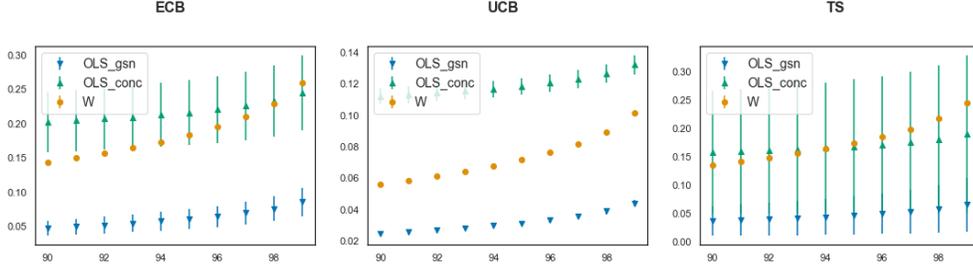


Figure 4: Multi-armed bandit results. Mean 2-sided confidence interval widths (error bars show 1 standard deviation) for the average reward $0.5\beta_1 + 0.5\beta_2$ in the MAB experiment.

320 4.1 Multi-armed bandits

321 In this section, we demonstrate the utility of the W -estimator for a stochastic multi-armed bandit
 322 setting. Villar et al. [2015] studied this problem in the context of patient allocation in clinical
 323 trials. Here the trial proceeds in a sequential fashion with the i^{th} patient given one of p treatments,
 324 encoded as $\mathbf{x}_i = \mathbf{e}_a$ with $a \in [p]$, and y_i denoting the outcome observed. We model the outcome
 325 as $y_i = \langle \mathbf{x}_i, \beta \rangle + \varepsilon_i$ where $\varepsilon_i \sim \text{Unif}([-1, 1])$ with $\beta = (0.3, 0.3)$ being the mean outcome of the
 326 treatments. Note that the two treatments are *statistically identical* in terms of outcome. As we will
 327 see, the adaptive sampling induced by the bandit strategies, however, introduces significant biases in
 328 the estimates.

329 We sequentially assign one of $p = 2$ treatments to each of $n = 1000$ patients using one of three
 330 policies (i) an ε -greedy policy (called ECB or Epsilon Current Belief), (ii) a practical UCB strategy
 331 based on the law of iterated logarithm (UCB) [Jamieson et al., 2014] and (iii) Thompson sampling
 332 [Thompson, 1933]. The ECB and TS sampling strategies are Bayesian. They place an independent
 333 Gaussian prior (with mean $\mu_0 = 0.3$ and variance $\sigma_0^2 = 0.33$) on each unknown mean outcome
 334 parameter and form an updated posterior belief concerning β following each treatment administration
 335 \mathbf{x}_i and observation y_i .

336 For ECB, the treatment administered to patient i is, with probability $1 - \varepsilon = .9$, the treatment with the
 337 largest posterior mean; with probability $1 - \varepsilon$, a uniformly random treatment is administered instead,
 338 to ensure sufficient exploration of all treatments. Note that this strategy satisfies condition (3) with
 339 $\mu_n(i) = \varepsilon/p$. For TS, at each patient i , a sample $\hat{\beta}$ of the mean treatment effect is drawn from the
 340 posterior belief. The treatment assigned to patient is the one maximizing the sampled mean treatment,
 341 i.e. $a_*(i) = \arg \max_{a \in [p]} \hat{\beta}_a$. In UCB, the algorithm maintains a score for each arm $a \in [p]$ that is a
 342 combination of the mean reward that the arm achieves and the empirical uncertainty of the reward.
 343 For each patient i , the UCB algorithm chooses the arm maximizing this score, and updates the score
 344 according to a fixed rule. For details on the specific implementation, see Jamieson et al. [2014]. Our
 345 goal is to produce confidence intervals for the β_a of each treatment based on the data adaptively
 346 collected from these standard bandit algorithms. We will compare the estimates and corresponding
 347 intervals for the *average reward* $0.5\beta_1 + 0.5\beta_2$. As the two arms/treatments are statistically identical,
 348 this isolates the effect of adaptive sampling on the obtained estimates.

349 We repeat the simulation for 5000 Monte Carlo runs. From each trial, we estimate the parameters β
 350 using both OLS and the W -estimator with $\lambda = \hat{\lambda}_{5\%, \pi}$ which is the 5th percentile of $\lambda_{\min}(n)$ achieved
 351 by the policy $\pi \in \{\text{ECB}, \text{UCB}, \text{TS}\}$. This choice is guided by Corollary 4.

352 We compare the quality of confidence regions for the average reward $0.5\beta_1 + 0.5\beta_2$ obtained from the
 353 W -decorrelated estimator, the OLS estimator with standard Gaussian theory (OLS_{gsn}), and the OLS
 354 estimator using concentration inequalities (OLS_{conc}) [Abbasi-Yadkori et al., 2011, Sec. 4]. Figure 3
 355 (left column) shows that the OLS Gaussian have have inconsistent coverage from the nominal. This
 356 is consistent with the observation that the sample means are biased negatively [Nie et al., 2017]. The
 357 concentration OLS tail bounds are all conservative, producing nearly 100% coverage, irrespective of
 358 the nominal level. This is intuitive, since they must account for the bias in sample means [Nie et al.,
 359 2017]. Meanwhile, the decorrelated intervals improves coverage uniformly over OLS intervals, often
 360 achieving the nominal coverage.

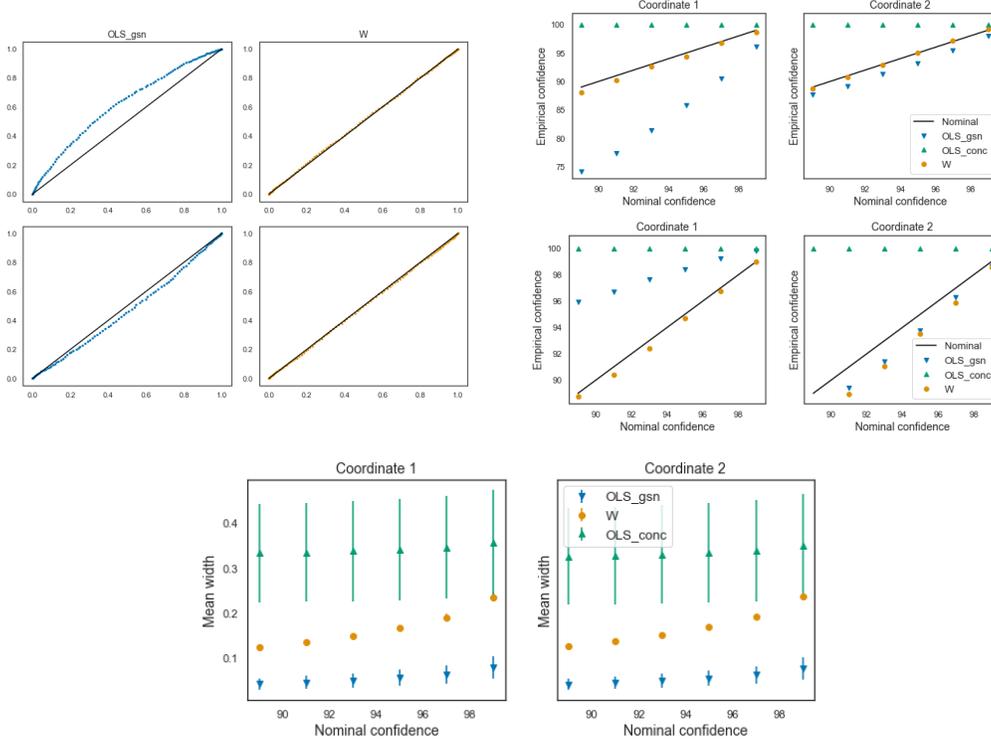


Figure 5: AR(2) time series results. Upper left: PP plot for the distribution of errors of standard OLS estimate and the W -decorrelated estimate. Upper right: Lower (top) and upper (bottom) coverage probabilities for OLS with Gaussian intervals, OLS with concentration inequality intervals, and decorrelated W -decorrelated estimate intervals. Note that ‘Conc’ has always 100% coverage. Bottom: Average 2 sided confidence interval widths obtained using the OLS estimator with standard Gaussian theory, OLS with concentration inequalities and the W -decorrelated estimator.

361 Figure 3 (right column) shows the PP plots of OLS and W -estimator errors for the average reward
 362 $0.5\beta_1 + 0.5\beta_2$. Recall that a PP plot between two distributions on the real line with densities P and
 363 Q is the parametric curve $(P(z), Q(z))$, $z \in \mathbb{R}$ [Gibbons and Chakraborti, 2011, Chapter 4.7]. The
 364 distribution of OLS errors is clearly seen to be distinctly non-Gaussian.

365 Figure 4 summarizes the distribution of 2-sided interval widths produced by each method for the
 366 sum reward. As expected, the W -decorrelation intervals are wider than those of OLS_{gsn} but compare
 367 favorably with those provided by OLS_{conc} . For UCB, the mean OLS_{conc} widths are always largest.
 368 For TS and ECB, W -decorrelation yields smaller intervals than OLS_{conc} for moderate confidence
 369 levels and comparable for high confidence levels. From this, we see that W -decorrelation intervals
 370 can be considerably less conservative than the concentration-based confidence intervals.

371 4.2 Autoregressive time series

372 In this section, we consider the classical $AR(p)$ model where $y_i = \sum_{\ell \leq p} \beta_\ell y_{i-\ell} + \varepsilon_i$. We generate
 373 data for the model with parameters $p = 2$, $n = 50$, $\beta = (0.95, 0.2)$, $y_0 = 0$ and $\varepsilon_i \sim \text{Unif}([-1, 1])$;
 374 all estimates are computed over 4000 monte carlo iterations.

375 We plot the coverage confidences for various values of the nominal on the right panel of Figure 5.
 376 The PP plot of the error distributions on the bottom right panel of Figure 5 shows that the OLS errors
 377 are skewed downwards, while the W -estimate errors are nearly Gaussian. We obtain the following
 378 improvements over the comparison methods of OLS standard errors OLS_{gsn} and concentration
 379 inequality widths OLS_{conc} [Abbasi-Yadkori et al., 2011]

380 The Gaussian OLS confidence regions systematically give incorrect empirical coverage. Meanwhile,
381 the concentration inequalities provide very conservative intervals, with nearly 100% coverage,
382 irrespective of the nominal level. In contrast, our decorrelated intervals achieve empirical coverage
383 that closely approximates the nominal levels. These coverage improvements are enabled by an
384 increase in width over that of OLS_{gsn} , but the W -estimate widths are systematically smaller than
385 those of the concentration inequalities.

386 Acknowledgements

387 The authors would like to thank Adel Javanmard and Lucas Janson for feedback on an earlier version
388 of this paper.

389 References

- 390 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Online least squares estimation with
391 self-normalized processes: An application to bandit problems. *arXiv preprint arXiv:1102.2670*,
392 2011.
- 393 Susan Athey and Stefan Wager. Efficient policy learning. *arXiv preprint arXiv:1702.02896*, 2017.
- 394 Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits.
395 In *COLT*, pages 217–226, 2009.
- 396 Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*
397 *Learning Research*, 3(Nov):397–422, 2002.
- 398 Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic
399 multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- 400 Rui M Castro and Robert D Nowak. Minimax bounds for active learning. *IEEE Transactions on*
401 *Information Theory*, 54(5):2339–2353, 2008.
- 402 Ngai H Chan and Ching-Zong Wei. Asymptotic inference for nearly nonstationary ar (1) processes.
403 *The Annals of Statistics*, pages 1050–1063, 1987.
- 404 Kani Chen, Inchi Hu, Zhiliang Ying, et al. Strong consistency of maximum quasi-likelihood estimators
405 in generalized linear models with fixed and adaptive designs. *The Annals of Statistics*, 27(4):
406 1155–1163, 1999.
- 407 Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit
408 feedback. In *COLT*, pages 355–366, 2008.
- 409 Yash Deshpande and Andrea Montanari. Linear bandits in high dimension and recommendation
410 systems. In *Communication, Control, and Computing (Allerton), 2012 50th Annual Allerton*
411 *Conference on*, pages 1750–1754. IEEE, 2012.
- 412 Maria Dimakopoulou, Susan Athey, and Guido Imbens. Estimation considerations in contextual
413 bandits. *arXiv preprint arXiv:1711.07077*, 2017.
- 414 Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. *arXiv*
415 *preprint arXiv:1103.4601*, 2011.
- 416 Aryeh Dvoretzky. Asymptotic normality for sums of dependent random variables. In *Proc. 6th*
417 *Berkeley Symp. Math. Statist. Probab*, volume 2, pages 513–535, 1972.
- 418 Walter Enders. *Applied econometric time series*. John Wiley & Sons, 2008.
- 419 Aurélien Garivier and Olivier Cappé. The kl-ucb algorithm for bounded stochastic bandits and
420 beyond. In *Proceedings of the 24th annual Conference On Learning Theory*, pages 359–376, 2011.
- 421 Jean Dickinson Gibbons and Subhabrata Chakraborti. *Nonparametric statistical inference*. Springer,
422 2011.

- 423 Peter Hall and Christopher C Heyde. *Martingale limit theory and its application*. Academic press,
424 2014.
- 425 Christopher C Heyde. *Quasi-likelihood and its application: a general approach to optimal parameter*
426 *estimation*. Springer Science & Business Media, 2008.
- 427 Guido W Imbens. The role of the propensity score in estimating dose-response functions. *Biometrika*,
428 87(3):706–710, 2000.
- 429 Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lilucb: An optimal
430 exploration algorithm for multi-armed bandits. In *Conference on Learning Theory*, pages 423–439,
431 2014.
- 432 Adel Javanmard and Andrea Montanari. Confidence intervals and hypothesis testing for high-
433 dimensional regression. *Journal of Machine Learning Research*, 15(1):2869–2909, 2014a.
- 434 Adel Javanmard and Andrea Montanari. Hypothesis testing in high-dimensional regression under the
435 gaussian random design model: Asymptotic theory. *IEEE Transactions on Information Theory*, 60
436 (10):6522–6554, 2014b.
- 437 Nathan Kallus. Balanced policy evaluation and learning. *arXiv preprint arXiv:1705.07384*, 2017.
- 438 TseLeung Lai and David Siegmund. Fixed accuracy estimation of an autoregressive parameter. *The*
439 *Annals of Statistics*, pages 478–485, 1983.
- 440 Tze Leung Lai. Asymptotic properties of nonlinear least squares estimates in stochastic regression
441 models. *The Annals of Statistics*, pages 1917–1930, 1994.
- 442 Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in*
443 *applied mathematics*, 6(1):4–22, 1985.
- 444 Tze Leung Lai and Ching Zong Wei. Least squares estimates in stochastic regression models with
445 applications to identification and control of dynamic systems. *The Annals of Statistics*, pages
446 154–166, 1982.
- 447 M. Ledoux. *Isoperimetry and Gaussian analysis*, volume 1648. Springer, Providence, 1996.
- 448 Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical
449 Soc., 2005.
- 450 Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
451 personalized news article recommendation. In *Proceedings of the 19th international conference on*
452 *World wide web*, pages 661–670. ACM, 2010.
- 453 Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of contextual-
454 bandit-based news article recommendation algorithms. In *Proceedings of the fourth ACM interna-*
455 *tional conference on Web search and data mining*, pages 297–306. ACM, 2011.
- 456 Jared K Lunceford and Marie Davidian. Stratification and weighting via the propensity score
457 in estimation of causal treatment effects: a comparative study. *Statistics in medicine*, 23(19):
458 2937–2960, 2004.
- 459 Xinkun Nie, Tian Xiaoying, Jonathan Taylor, and James Zou. Why adaptively collected data have
460 negative bias and how to correct for it. 2017.
- 461 Peter CB Phillips and Pierre Perron. Testing for a unit root in time series regression. *Biometrika*, 75
462 (2):335–346, 1988.
- 463 Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected*
464 *Papers*, pages 169–177. Springer, 1985.
- 465 Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of*
466 *Operations Research*, 35(2):395–411, 2010.

- 467 Daniel Russo. Simple bayesian algorithms for best arm identification. In *Conference on Learning*
468 *Theory*, pages 1417–1418, 2016.
- 469 Robert H Shumway and David S Stoffer. *Time series analysis and its applications: with R examples*.
470 Springer Science & Business Media, 2006.
- 471 Adith Swaminathan and Thorsten Joachims. Batch learning from logged bandit feedback through
472 counterfactual risk minimization. *Journal of Machine Learning Research*, 16:1731–1755, 2015.
- 473 William R Thompson. On the likelihood that one unknown probability exceeds another in view of
474 the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- 475 Sara Van de Geer, Peter Bühlmann, Yaacov Ritov, Ruben Dezeure, et al. On asymptotically optimal
476 confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–
477 1202, 2014.
- 478 Sofia Villar, Jack Bowden, and James Wason. Multi-armed bandit models for the optimal design
479 of clinical trials: benefits and challenges. *Statistical science: a review journal of the Institute of*
480 *Mathematical Statistics*, 30(2):199, 2015.
- 481 Ching-Zong Wei. Asymptotic properties of least-squares estimates in stochastic regression models.
482 *The Annals of Statistics*, pages 1498–1508, 1985.
- 483 Cun-Hui Zhang and Stephanie S Zhang. Confidence intervals for low dimensional parameters in
484 high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical*
485 *Methodology)*, 76(1):217–242, 2014.

486 **References**

- 487 Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Online least squares estimation with
488 self-normalized processes: An application to bandit problems. *arXiv preprint arXiv:1102.2670*,
489 2011.
- 490 Susan Athey and Stefan Wager. Efficient policy learning. *arXiv preprint arXiv:1702.02896*, 2017.
- 491 Jean-Yves Audibert and Sébastien Bubeck. Minimax policies for adversarial and stochastic bandits.
492 In *COLT*, pages 217–226, 2009.
- 493 Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine*
494 *Learning Research*, 3(Nov):397–422, 2002.
- 495 Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic
496 multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- 497 Rui M Castro and Robert D Nowak. Minimax bounds for active learning. *IEEE Transactions on*
498 *Information Theory*, 54(5):2339–2353, 2008.
- 499 Ngai H Chan and Ching-Zong Wei. Asymptotic inference for nearly nonstationary ar (1) processes.
500 *The Annals of Statistics*, pages 1050–1063, 1987.
- 501 Kani Chen, Inchi Hu, Zhiliang Ying, et al. Strong consistency of maximum quasi-likelihood estimators
502 in generalized linear models with fixed and adaptive designs. *The Annals of Statistics*, 27(4):
503 1155–1163, 1999.
- 504 Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under bandit
505 feedback. In *COLT*, pages 355–366, 2008.
- 506 Yash Deshpande and Andrea Montanari. Linear bandits in high dimension and recommendation
507 systems. In *Communication, Control, and Computing (Allerton), 2012 50th Annual Allerton*
508 *Conference on*, pages 1750–1754. IEEE, 2012.
- 509 Maria Dimakopoulou, Susan Athey, and Guido Imbens. Estimation considerations in contextual
510 bandits. *arXiv preprint arXiv:1711.07077*, 2017.

- 511 Miroslav Dudík, John Langford, and Lihong Li. Doubly robust policy evaluation and learning. *arXiv*
512 *preprint arXiv:1103.4601*, 2011.
- 513 Aryeh Dvoretzky. Asymptotic normality for sums of dependent random variables. In *Proc. 6th*
514 *Berkeley Symp. Math. Statist. Probab*, volume 2, pages 513–535, 1972.
- 515 Walter Enders. *Applied econometric time series*. John Wiley & Sons, 2008.
- 516 Aurélien Garivier and Olivier Cappé. The kl-ucb algorithm for bounded stochastic bandits and
517 beyond. In *Proceedings of the 24th annual Conference On Learning Theory*, pages 359–376, 2011.
- 518 Jean Dickinson Gibbons and Subhabrata Chakraborti. *Nonparametric statistical inference*. Springer,
519 2011.
- 520 Peter Hall and Christopher C Heyde. *Martingale limit theory and its application*. Academic press,
521 2014.
- 522 Christopher C Heyde. *Quasi-likelihood and its application: a general approach to optimal parameter*
523 *estimation*. Springer Science & Business Media, 2008.
- 524 Guido W Imbens. The role of the propensity score in estimating dose-response functions. *Biometrika*,
525 87(3):706–710, 2000.
- 526 Kevin Jamieson, Matthew Malloy, Robert Nowak, and Sébastien Bubeck. lilucb: An optimal
527 exploration algorithm for multi-armed bandits. In *Conference on Learning Theory*, pages 423–439,
528 2014.
- 529 Adel Javanmard and Andrea Montanari. Confidence intervals and hypothesis testing for high-
530 dimensional regression. *Journal of Machine Learning Research*, 15(1):2869–2909, 2014a.
- 531 Adel Javanmard and Andrea Montanari. Hypothesis testing in high-dimensional regression under the
532 gaussian random design model: Asymptotic theory. *IEEE Transactions on Information Theory*, 60
533 (10):6522–6554, 2014b.
- 534 Nathan Kallus. Balanced policy evaluation and learning. *arXiv preprint arXiv:1705.07384*, 2017.
- 535 TseLeung Lai and David Siegmund. Fixed accuracy estimation of an autoregressive parameter. *The*
536 *Annals of Statistics*, pages 478–485, 1983.
- 537 Tze Leung Lai. Asymptotic properties of nonlinear least squares estimates in stochastic regression
538 models. *The Annals of Statistics*, pages 1917–1930, 1994.
- 539 Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Advances in*
540 *applied mathematics*, 6(1):4–22, 1985.
- 541 Tze Leung Lai and Ching Zong Wei. Least squares estimates in stochastic regression models with
542 applications to identification and control of dynamic systems. *The Annals of Statistics*, pages
543 154–166, 1982.
- 544 M. Ledoux. *Isoperimetry and Gaussian analysis*, volume 1648. Springer, Providence, 1996.
- 545 Michel Ledoux. *The concentration of measure phenomenon*. Number 89. American Mathematical
546 Soc., 2005.
- 547 Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
548 personalized news article recommendation. In *Proceedings of the 19th international conference on*
549 *World wide web*, pages 661–670. ACM, 2010.
- 550 Lihong Li, Wei Chu, John Langford, and Xuanhui Wang. Unbiased offline evaluation of contextual-
551 bandit-based news article recommendation algorithms. In *Proceedings of the fourth ACM interna-*
552 *tional conference on Web search and data mining*, pages 297–306. ACM, 2011.
- 553 Jared K Lunceford and Marie Davidian. Stratification and weighting via the propensity score
554 in estimation of causal treatment effects: a comparative study. *Statistics in medicine*, 23(19):
555 2937–2960, 2004.

- 556 Xinkun Nie, Tian Xiaoying, Jonathan Taylor, and James Zou. Why adaptively collected data have
557 negative bias and how to correct for it. 2017.
- 558 Peter CB Phillips and Pierre Perron. Testing for a unit root in time series regression. *Biometrika*, 75
559 (2):335–346, 1988.
- 560 Herbert Robbins. Some aspects of the sequential design of experiments. In *Herbert Robbins Selected
561 Papers*, pages 169–177. Springer, 1985.
- 562 Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of
563 Operations Research*, 35(2):395–411, 2010.
- 564 Daniel Russo. Simple bayesian algorithms for best arm identification. In *Conference on Learning
565 Theory*, pages 1417–1418, 2016.
- 566 Robert H Shumway and David S Stoffer. *Time series analysis and its applications: with R examples*.
567 Springer Science & Business Media, 2006.
- 568 Adith Swaminathan and Thorsten Joachims. Batch learning from logged bandit feedback through
569 counterfactual risk minimization. *Journal of Machine Learning Research*, 16:1731–1755, 2015.
- 570 William R Thompson. On the likelihood that one unknown probability exceeds another in view of
571 the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- 572 Sara Van de Geer, Peter Bühlmann, Yaacov Ritov, Ruben Dezeure, et al. On asymptotically optimal
573 confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–
574 1202, 2014.
- 575 Sofia Villar, Jack Bowden, and James Wason. Multi-armed bandit models for the optimal design
576 of clinical trials: benefits and challenges. *Statistical science: a review journal of the Institute of
577 Mathematical Statistics*, 30(2):199, 2015.
- 578 Ching-Zong Wei. Asymptotic properties of least-squares estimates in stochastic regression models.
579 *The Annals of Statistics*, pages 1498–1508, 1985.
- 580 Cun-Hui Zhang and Stephanie S Zhang. Confidence intervals for low dimensional parameters in
581 high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical
582 Methodology)*, 76(1):217–242, 2014.