

CONTROL VARIATE SELECTION FOR MONTE CARLO INTEGRATION

Rémi Leluc;
François Portier and Johan Segers

DISCUSSION PAPER | 2019 / 15



Control variate selection for Monte Carlo integration

Rémi Leluc* François Portier* Johan Segers†

Abstract

Monte Carlo integration with variance reduction by means of control variates can be implemented by the ordinary least squares estimator for the intercept in a multiple linear regression model with the integrand as response and the control variates as covariates. Even without special knowledge on the integrand, significant efficiency gains can be obtained if the control variate space is sufficiently large. Incorporating a large number of control variates in the ordinary least squares procedure may however result in (i) a certain instability of the ordinary least squares estimator and (ii) a possibly prohibitive computation time. Regularizing the ordinary least squares estimator by preselecting appropriate control variates via the Lasso turns out to increase the accuracy without additional computational cost. The findings in the numerical experiment are confirmed by concentration inequalities for the integration error.

1 Introduction

Whereas the basic Monte Carlo (MC) estimate is given by $(1/n)\sum_i f_i$, for independent and identically distributed random variables f_i , the control variates method is based on $(1/n)\sum_i(f_i + h_i)$, where the h_i variables, called control variates, are constructed to have zero expectation. When the controls h_i have been selected or estimated properly (based on the samples f_i), the use of control variates might reduce the variance of the basic MC estimate significantly. The method of control variates, already used frequently to compute prices of financial derivatives [4], has been employed recently in many different fields of Machine Learning. Examples include (i) *reinforcement learning* and more particularly *policy gradient* methods [8, 9] where the score function permits to define many control variates ; (ii) inference in complex probabilistic models [15] where the Stein method allows to define accurate control variates [10]; and (iii) gradient based *optimization* [20, 6].

Suppose that $m \geq 1$ control variates are available and $n \geq 1$ samples have been generated. Any linear combination of control variates can be used as a particular control variate. In terms of the variance of the estimation error, the optimal linear combination can be estimated based on the empirical risk minimization principle applied to an ordinary least squares (OLS) regression problem [see Eq. (2.3) below]. This approach, referred to as OLSMC, is the most common implementation of the control variates method as detailed for instance in [12, Section 8.3] or [14, 17], although other implementations are possible, see Remark 2 below.

Asymptotically, the OLSMC error is bounded by the MC error and is proportional to the L_2 approximation error of the integrand in the linear span of control variates [5]. In combination with well-known approximation results in L_p -spaces [16], this representation of the OLSMC error suggests to use an increasing number of control variates. Indeed, in [14] it is shown that when m grows with n , the OLSMC error rate can be faster than $1/\sqrt{n}$.

However, when based on a large number of control variates, the OLSMC suffers from two classical problems common for least squares methods: (i) numerical instabilities when the control variates are nearly collinear, and (ii) a computational complexity in $m^3 + nm^2$, which might be prohibitive.

To deal with these two issues, it has been proposed in [17] to regularize the OLSMC estimate by adding a ℓ_1 -penalty term in the minimization problem, just as in the LASSO [18]. Simulation results in [17] show that

*Télécom Paris, Institut Polytechnique de Paris.

†Institut de Statistique, Biostatistique et Sciences Actuarielles, Université catholique de Louvain.

this approach, referred to as LASSOMC, provides great improvements in practice. However, those practical findings are not supported by an asymptotic error rate nor by a non-asymptotic error bound.

The main objective of the paper is to provide a non-asymptotic theory for the use of control variates in Monte Carlo simulations. The contributions are as follows.

1. A *new method* called LSLASSOMC is proposed. In the spirit of [1], it consists in selecting the best control variates via the LASSO, using subsampling to decrease the computation time, and then to apply OLSMC with the selected controls.
2. *Support recovery*: the LASSO is shown to select the correct control variates with large probability.
3. *Concentration inequalities* are derived for the OLSMC and LASSOMC errors. The one for the OLSMC highlights a compromise between the approximation error of the integrand in the linear span of control variates and the multicollinearities between the control variates. The one for LASSOMC shows significant improvements regarding the effects of multicollinearity.

The outline of the paper is as follows. Section 2 introduces the theoretical background and the different MC estimates and provides some comments about their practical implementation and some possible alternative approaches. Section 3 contains the statements of the theoretical results. Section 4 is a simulation study to illustrate the practical behavior of the methods.

All proofs are gathered in Appendix. The approach combines well known sub-Gaussian concentration inequalities [2] with a recent concentration bound for the smallest eigenvalue of an empirical Gram matrix [21].

2 Monte Carlo integration and control variates

Background. Let $f \in L^2(P)$ be a square integrable, real-valued function on a probability space (S, \mathcal{S}, P) of which we would like to calculate the integral

$$P(f) = \int_S f(x) P(dx).$$

The MC estimator of $P(f)$ based on an independent random sample X_1, \dots, X_n from P is

$$\hat{\alpha}_n^{\text{mc}}(f) = P_n(f) = \frac{1}{n} \sum_{i=1}^n f(X_i).$$

This estimator is unbiased and has variance $n^{-1}\sigma_0^2(f)$, where $\sigma_0^2(f) = P[(f - P(f))^2]$.

The control variates are functions $h_1, \dots, h_m \in L^2(P)$ with known expectations. Without loss of generality, assume that $P(h_k) = 0$ for all $k \in \{1, \dots, m\}$. Let $h = (h_1, \dots, h_m)^T$ denote the \mathbb{R}^m -valued function with the m control variates as elements. Let $\mathcal{F}_m = \text{Span}\{h_1, \dots, h_m\} = \{\beta^T h : \beta \in \mathbb{R}^m\}$ denote the closed linear subspace of $L^2(P)$ generated by the control variates.

For any $\beta = (\beta_1, \dots, \beta_m)^T \in \mathbb{R}^m$, we have $P(f - \beta^T h) = P(f)$, so that $P_n(f - \beta^T h)$ is an unbiased estimator of $P(f)$, with variance $n^{-1}P[(f - P(f) - \beta^T h)^2]$. Any coefficient vector

$$\beta^*(f) \in \arg \min_{\beta \in \mathbb{R}^m} P[(f - P(f) - \beta^T h)^2]$$

minimizes the variance. If such a $\beta^*(f)$ would be known, the resulting oracle estimator would be

$$\hat{\alpha}_n^{\text{or}}(f) = P_n[f - \beta^*(f)^T h]. \tag{2.1}$$

By definition, the oracle estimator achieves the minimal variance $n^{-1}\sigma_m^2(f)$ where $\sigma_m^2(f)$ is the minimum value of $P[(f - P(f) - \beta^T h)^2]$ with respect to β . For any $m' = 0, 1, \dots, m$, if we use only the first m'

control variates $h_1, \dots, h_{m'}$, or even none at all in case $m' = 0$, we have $\sigma_m^2(f) \leq \sigma_{m'}^2(f)$. In particular, if $\beta^*(f)$ would be known, the use of control variates would always reduce the variance of the basic Monte Carlo estimator.

As $\beta^*(f)^T h$ is the $L_2(P)$ -projection of $f - P(f)$ on the linear vector space \mathcal{F}_m and since the control variates are centered, $\beta^*(f)$ satisfies the normal equations $P(hh^T)\beta^*(f) = P(hf)$. The integral $P(f)$ thus appears as the intercept of a linear regression model with response f and explanatory variables h_1, \dots, h_m , and it can be expressed as

$$(P(f), \beta^*(f)) \in \underset{(\alpha, \beta) \in \mathbb{R} \times \mathbb{R}^m}{\arg \min} P[(f - \alpha - \beta^T h)^2]. \quad (2.2)$$

The empirical risk minimization paradigm applied to the risk function on the right-hand side of (2.2) will lead to the OLSMC and LASSOMC estimates, to be defined further in this section. The same paradigm suggests the use of other regression methods for MC integration such as Principal Component Regression (PCR) or Ridge Regression [3], both of which will be considered in the numerical experiments.

Remark 1 (Choice of control variates). Which control variates work well depends on the problem. In the Black–Scholes model, for instance, an effective control variate for the price of an option is the geometric average of the price series [4, Example 4.1.2]). Two generic ways to construct control variates are to be noted. Whenever $P(dx) = w(x)Q(dx)$, where $w : S \rightarrow [0, \infty)$ and Q is a probability measure on (S, \mathcal{S}) , the quantity of interest is $P(f) = Q(wf)$, so that we can use control variates for wf with respect to Q . This trick can be useful in combination with importance sampling [11]. If P has density p with respect to the Lebesgue measure and if we have access to the derivatives of p , Stein’s method might be used to build infinitely many control functions [10].

Ordinary Least Squares Monte Carlo. Replacing the distribution P by the sample measure P_n in (2.2), we obtain the OLSMC estimator $\hat{\alpha}_n^{\text{ols}}(f)$ of $P(f)$ as a minimizer of the empirical risk:

$$(\hat{\alpha}_n^{\text{ols}}(f), \hat{\beta}_n^{\text{ols}}(f)) \in \underset{(\alpha, \beta) \in \mathbb{R} \times \mathbb{R}^m}{\arg \min} \|f^{(n)} - \alpha \mathbf{1}_n - H\beta\|_2^2, \quad (2.3)$$

where $\|\cdot\|_2$ denotes the Euclidean norm, $f^{(n)} = (f(X_1), \dots, f(X_n))^T \in \mathbb{R}^n$, $\mathbf{1}_n = (1, \dots, 1)^T \in \mathbb{R}^n$, and H is the random $n \times m$ matrix defined by

$$H = (h_j(X_i))_{\substack{i=1, \dots, n \\ j=1, \dots, m}}.$$

The minimization problem in (2.3) can be expressed using an OLS estimate with centered variables as

$$\begin{aligned} \hat{\alpha}_n^{\text{ols}}(f) &= P_n[f - \hat{\beta}_n^{\text{ols}}(f)^T h], \\ \hat{\beta}_n^{\text{ols}}(f) &= \underset{\beta \in \mathbb{R}^m}{\arg \min} \|f_c^{(n)} - H_c \beta\|_2^2, \end{aligned} \quad (2.4)$$

where $f_c^{(n)} = f^{(n)} - \mathbf{1}_n(\mathbf{1}_n^T f^{(n)})/n$ and $H_c = H - \mathbf{1}_n(\mathbf{1}_n^T H)/n$. Indeed, for fixed $\beta \in \mathbb{R}^m$, the minimizer over $\alpha \in \mathbb{R}$ of the objective function in Eq. (2.3) is just $P_n(f - \beta^T h) = P_n(f) - \beta^T P_n(h)$, and since $P_n(f) = (\mathbf{1}_n^T f^{(n)})/n$ and $P_n(h) = (\mathbf{1}_n^T H)/n$, the equivalence of (2.3) and (2.4) follows.

Remark 2 (Variations). The objective function in (2.4) involves the empirical covariance matrix $n^{-1}H_c^T H_c = P_n(hh^T) - P_n(h)P_n(h^T)$. Using different estimates of the Gram matrix $P(hh^T)$ leads to alternative control variate MC estimates for $P(f)$ [5, 14]. For fixed m and as $n \rightarrow \infty$, all these estimators are consistent and asymptotically normal. The OLSMC, however, is the only one that can integrate both the constant functions and the control functions without error. In [14], the alternative estimators have been shown to perform poorly compared to the OLSMC.

Remark 3 (Invariance). The OLSMC estimator does not change if we replace the control variate vector h by Ah , where A is an arbitrary invertible $m \times m$ matrix. Provided the control functions are linearly independent, the property of isotropy, i.e., $P(hh^T) = I_m$, can always be enforced by an appropriate linear transformation of the vector of control variates.

Remark 4 (Computation time). The reliance on least squares makes the OLSMC computing time to be in $nm^2 + m^3 + nT$, where T stands for the time needed to evaluate f . Computational benefits occur when there are multiple integrands, since the OLSMC estimate can be represented as $w^T f^{(n)}$, where the weight vector $w \in \mathbb{R}^n$ does not depend on the integrands [14]. If q integrals need to be evaluated, the computing time then becomes $nm^2 + m^3 + qnT$.

LASSO Monte Carlo. The LASSO, introduced in [18], is a regression technique that consists in minimizing the usual least squares loss plus an ℓ_1 -penalty term on the vector of regression coefficients. In contrast with OLS, the LASSO usually produces a vector with many zero coefficients, meaning that the corresponding variables are no longer included in the predictive model. The LASSO thus achieves estimation and variable selection at the same time. As the use of control variates in MC integration is linked with regression, the LASSO can be used to take advantage from situations where many control variates are present but not all of them are useful.

The LASSOMC estimator $\hat{\alpha}_n^{\text{lasso}}(f)$ of $P(f)$ follows from (2.3), adding a penalization to the regression coefficient: we have

$$(\hat{\alpha}_n^{\text{lasso}}(f), \hat{\beta}_n^{\text{lasso}}(f)) = \arg \min_{(\alpha, \beta) \in \mathbb{R} \times \mathbb{R}^m} \frac{1}{2n} \|f^{(n)} - \alpha \mathbb{1}_n - H\beta\|_2^2 + \lambda \|\beta\|_1,$$

where $\|\cdot\|_1$ denotes the ℓ_1 -norm. By the same argument used to justify the equivalence of (2.3) and (2.4), the LASSOMC can be based on centered variables via

$$\begin{aligned} \hat{\alpha}_n^{\text{lasso}}(f) &= P_n[f - \hat{\beta}_n^{\text{lasso}}(f)^T h], \\ \hat{\beta}_n^{\text{lasso}}(f) &= \arg \min_{\beta \in \mathbb{R}^m} \frac{1}{2n} \|f_c^{(n)} - H_c \beta\|_2^2 + \lambda \|\beta\|_1. \end{aligned} \quad (2.5)$$

LSLASSO Monte Carlo. Another approach is to use the LASSO to select the active variables among a large number of control variates and then to compute OLSMC using only the variables selected at the previous stage. We refer to this approach as the LSLASSOMC. To decrease the computation time when the dimensions involved in the problem, either n or m , are large, we recommend to use sub-sampling of a smaller size $N \leq n$ when conducting the first step.

Let $\hat{S} = \{k \in \{1, \dots, m\} : \hat{\beta}_{N,k}^{\text{lasso}}(f) > 0\}$ denote the estimated active set of control variates based on the subsample of size N . The LSLASSOMC estimate $\hat{\alpha}_n^{\text{lslasso}}(f)$ of $P(f)$ is defined by

$$(\hat{\alpha}_n^{\text{lslasso}}(f), \hat{\beta}_n^{\text{lslasso}}(f)) \in \arg \min_{(\alpha, \beta) \in \mathbb{R} \times \mathbb{R}^{\hat{\ell}}} \|f^{(n)} - \alpha \mathbb{1}_n - H_{\hat{S}} \beta\|_2^2, \quad (2.6)$$

where $\hat{\ell} = |\hat{S}|$ and $H_{\hat{S}}$ is the $n \times \hat{\ell}$ matrix made of the columns of H with index k in \hat{S} .

Remark 5 (Computation). As for the LASSO, the LASSOMC and LSLASSOMC can be computed by *cyclical coordinate descent* using at each step the *soft-thresholding operator* [18, Section 2.4]. The LASSOMC then requires nD operations, where D stands for the number of updated coordinates in $\hat{\beta}_n^{\text{lasso}}(f)$. This kind of optimization strategy allows to compute approximate solutions in a reduced time by, for instance, choosing at random the coordinates to update and thus reducing D . For the LSLASSOMC, the number of operations would be in $ND + n\hat{\ell}^2 + \hat{\ell}^3$, combining the cost of selecting the control variates on the subsample of size N and running the OLSMC estimate based on the selected control variates for the full sample of size n .

3 Non-asymptotic bounds

To derive concentration inequalities for the errors of the estimators proposed in Section 2, we use the notion of sub-Gaussianity as defined for instance in [2, Section 2.3]. Recall that the moment generating function of a centered Gaussian random variable with variance σ^2 is equal to $\lambda \mapsto \exp(\lambda^2 \sigma^2 / 2)$.

Definition 1. A centered random variable Y is sub-Gaussian with variance factor $\tau^2 > 0$, notation $Y \in \mathcal{G}(\tau^2)$, if $\log \mathbb{E}[\exp(\lambda Y)] \leq \lambda^2 \tau^2 / 2$ for all $\lambda \in \mathbb{R}$.

If $Y \in \mathcal{G}(\tau^2)$, then necessarily $\text{Var}(Y) \leq \tau^2$ [2, Exercise 2.16]. Chernoff's inequality provides exponential bounds on the tails of sub-Gaussian random variables. Moreover, the sum of independent sub-Gaussian variables is again sub-Gaussian. Centered, bounded random variables taking values in an interval $[a, b]$ are sub-Gaussian with variance factor at most $(b - a)^2 / 4$ [2, Lemma 2.2].

The concentration inequalities for the various Monte Carlo methods with control variates will be largely due to the following assumption that requires the residuals to be sub-Gaussian.

Assumption 1 (sub-Gaussian residuals). *The residual function $\epsilon = f - P(f) - \beta^*(f)^T h$ satisfies $\epsilon \in \mathcal{G}(\tau^2)$ for some $\tau > 0$, that is, $\int_{\mathcal{S}} \exp(\lambda x) \epsilon(x) P(dx) \leq \exp(\lambda^2 \tau^2 / 2)$ for all $\lambda \in \mathbb{R}$.*

The estimation error of the oracle estimator in (2.1) is just $\hat{\alpha}_n^{\text{or}} - P(f) = P_n(\epsilon) = n^{-1} \sum_{i=1}^n \epsilon(X_i)$. Under Assumption 1, this is a sub-Gaussian variable with variance factor τ_m^2 / n . Chernoff's inequality [2, p. 25] then implies that for all $\delta \in (0, 1)$ and all $n = 1, 2, \dots$, with probability at least $1 - \delta$,

$$|\hat{\alpha}_n^{\text{or}}(f) - P(f)| \leq \sqrt{2 \log(2/\delta)} \frac{\tau}{\sqrt{n}} \quad (3.1)$$

This concentration inequality provides a baseline when the best possible control variate in the space \mathcal{F}_m is selected. The case $m = 0$ also covers the basic MC method: in that case: τ^2 is the variance factor of the sub-Gaussian variable $f - P(f)$ on (S, \mathcal{S}, P) .

From now on, consider $m \geq 1$ control variates $h_1, \dots, h_m \in L^2(P)$, all of which are centered. To analyze the OLSMC method we need some additional assumptions on the control functions. For a function $v : S \rightarrow \mathbb{R}$, write $\|v\|_{\infty} = \sup_{x \in S} |v(x)|$.

Assumption 2 (Bounded control variates). *The control variates are uniformly bounded. Put $U := \max_{j=1, \dots, m} \|h_j\|_{\infty}$.*

Assumption 3 (Linearly independent control variates). *The control variates $h_1, \dots, h_m \in L^2(P)$ are linearly independent. As a consequence, the Gram matrix $G := P(hh^T) \in \mathbb{R}^m$ is positive definite and its smallest eigenvalue $\gamma := \lambda_{\min}(G)$ is positive.*

The error OLSMC estimation error is subject to the following concentration bound.

Theorem 3.1 (Concentration inequality for OLSMC). *Assume Assumptions 1, 2 and 3 hold. Write $\zeta_h = U^2 / \gamma$. Then for all $\delta \in (0, 1)$ and all integer n such that*

$$n \geq 8m \log(8m/\delta) \zeta_h \quad \text{and} \quad n \geq 4 \|h^T G^{-1} h\|_{\infty} (32m + 4 \log(4/\delta))$$

we have, with probability at least $1 - \delta$,

$$|\hat{\alpha}_n^{\text{ols}}(f) - P(f)| \leq \sqrt{2 \log(8/\delta)} \frac{\tau}{\sqrt{n}} + 27m \log(8m/\delta) \zeta_h \frac{\tau}{n}. \quad (3.2)$$

Compared to the bound (3.1) for the oracle estimator, the bound (3.2) for OLSMC has an additional term. This term is due to the additional learning step that is needed to estimate the optimal control variate.

Remark 6 (On the factor ζ_h). The smallest eigenvalue of G being bounded by the mean of the eigenvalues, we have $\gamma \leq m^{-1} \sum_{j=1}^m P(h_j^2) \leq U^2$ and thus $\zeta_h = U^2 / \gamma \geq 1$. Further, the quantity ζ_h does not change if all control variates h_1, \dots, h_m are scaled the same way.

Remark 7 (On the factor $\|h^T G^{-1} h\|_{\infty}$). By the cyclic property of the trace operator, we have $P(h^T G^{-1} h) = P[\text{tr}(G^{-1} h h^T)] = \text{tr}(I_m) = m$, and therefore $\|h^T G^{-1} h\|_{\infty} \geq m$. We thus need $n \geq 128m^2$. The function $h^T G^{-1} h$ remains invariant under invertible linear transformations of the vector h and thus depends only on the control space \mathcal{F}_m . The inequality $n \geq 4 \|h^T G^{-1} h\|_{\infty} (32m + 4 \log(4/\delta))$ is a finite-sample version of the (asymptotic) Newey condition $\|h^T G^{-1} h\|_{\infty} = o(n/m)$ as $n \rightarrow \infty$ in [14].

Remark 8 (Rates). Consider an asymptotic set-up where the number of control variates m tends to infinity with n . The OLSMC method improves upon the basic MC method ($m = 0$), which has rate $1/\sqrt{n}$, as soon as $\tau + \tau\zeta_h m \log(m)/\sqrt{n} \rightarrow 0$. To recover the same order as the one of the oracle estimator $\hat{\alpha}_n^{\text{of}}(f)$, which has rate τ/\sqrt{n} , one must have $m \log(m)\zeta_h = O(\sqrt{n})$ as $n \rightarrow \infty$. This means that m must be not too large compared to n .

The LASSOMC takes advantage of *sparse* regression models. A regression model is sparse whenever many of the coefficients of the parameter vector β are equal to zero, i.e., many of the covariates are useless to predict the output in the presence of the other covariates. The *active set* associated to the coefficient vector $\beta \in \mathbb{R}^m$ is

$$S(\beta) = \{j = 1, \dots, m : \beta_j \neq 0\}.$$

The number of elements in $S^* = S(\beta^*(f))$, denoted by $\ell^* := |S^*|$, quantifies the level of sparsity associated to the regression model. We will see that the LASSOMC improves upon the OLS whenever ℓ^* becomes small compared to m .

We follow the approach of [19, Section 11.4.1], in which the analysis of the LASSO is carried out using a *restricted eigenvalue condition* dealing only with the directions in the active set, discarding the non-active directions. For a vector $\beta \in \mathbb{R}^m$ and an ordered set $S = (k_1, \dots, k_\ell) \subset \{1, \dots, m\}$, let $\beta_S = (\beta_{k_1}, \dots, \beta_{k_\ell})^T$.

Assumption 4 (Linearly independent active control variates). *The active control variates h_k , $k \in S^*$, are linearly independent. As a consequence, the $\ell^* \times \ell^*$ Gram matrix $G_{S^*} = P(h_{S^*} h_{S^*}^T)$ is positive definite and its smallest eigenvalue $\gamma^* := \lambda_{\min}(G_{S^*})$ is strictly positive.*

Needed also will be that the active control functions are orthogonal, in $L_2(P)$, to the inactive ones.

Assumption 5 (Orthogonality between active and inactive controls). *We have $P(h_j h_k) = 0$ for all $j \in S \setminus S^*$ and all $k \in S^*$.*

Recall that the ℓ_1 -penalty of the LASSO is weighted by a regularization parameter $\lambda > 0$.

Theorem 3.2 (Support recovery of LASSOMC). *If Assumptions 1, 2, 4 and 5 hold, then for all $\delta \in (0, 1)$, all integer n such that*

$$n \geq 4 \|h_{S^*}^T G_{S^*}^{-1} h_{S^*}\|_\infty (32\ell^* + 4 \log(5/\delta)) \quad \text{and} \quad n \geq 70(\ell^*)^2 \log(10\ell^* m/\delta) (U^2/\gamma^*)^2,$$

and all λ such that

$$17\sqrt{\log(10m/\delta)} U \frac{\tau}{\sqrt{n}} \leq \lambda < \frac{\gamma^*}{3\sqrt{\ell^*}} \min_{j \in S^*} |\beta_j^*(f)|, \quad (3.3)$$

it holds that, with probability at least $1 - \delta$, the LASSO solution $\hat{\beta}_n^{\text{lasso}}(f)$ in (2.5) is unique and $S(\beta^*(f)) = S(\hat{\beta}_n^{\text{lasso}}(f))$.

Theorem 3.3 (Concentration inequality for LASSOMC). *Under the same conditions as Theorem 3.2, we have, with probability at least $1 - \delta$,*

$$|\hat{\alpha}_n^{\text{lasso}}(f) - P(f)| \leq \sqrt{2 \log(10/\delta)} \frac{\tau}{\sqrt{n}} + 9\lambda\ell^* \sqrt{\log(10m/\delta)} \frac{U/\gamma^*}{\sqrt{n}}.$$

For λ equal to the lower bound in (3.3), we have, on the same event,

$$|\hat{\alpha}_n^{\text{lasso}}(f) - P(f)| \leq \sqrt{2 \log(10/\delta)} \frac{\tau}{\sqrt{n}} + 153\ell^* \log(10m/\delta) (U^2/\gamma^*) \frac{\tau}{n}. \quad (3.4)$$

The benefits of LASSOMC over OLSMC can be observed by comparing the bounds in (3.2) and (3.4). The total number m , of control functions has been replaced by the active number ℓ^* of such functions. Further, the smallest eigenvalue γ^* of G_{S^*} in Assumption 3 is at least as large as the smallest eigenvalue γ of G in Assumption 4.

4 Numerical application

To compare the practical performance of the different MC estimates using control variates, we focus on the standard integration problem over the unit cube $[0, 1]^d$. The goal is to compute $\int_{[0, 1]^d} f(x) dx$. We shall consider various dimensions $d \geq 1$, different integrands $f : [0, 1]^d \rightarrow \mathbb{R}$, and several choices for the computational budget, n , and the number of control variates, m . We shall focus on difficult situations where d is relatively large compared to n . For ease of reproducibility, the code is available upon request.

Control variates. Multivariate control functions with respect to the uniform distribution over $[0, 1]^d$ are easy to construct based on univariate ones. Let (h_1, \dots, h_K) be a vector of one-dimensional control functions, i.e., $\int_0^1 h_k(x) dx = 0$ for each $k = 1, \dots, K$. Without further information on the integrand, the usual way to construct multivariate controls is by forming tensor products of the form $h_\ell(x_1, \dots, x_d) = \prod_{j=1}^d h_{\ell_j}(x_j)$, for a multi-index $\ell = (\ell_1, \dots, \ell_d)$ in $\{0, 1, \dots, K\}^d \setminus \{(0, \dots, 0)\}$, yielding a total number of $m = (K + 1)^d - 1$ control functions. A drawback of such a construction is that the number of control functions grows quickly with K . Alternative approaches yielding smaller control spaces consist of imposing $\ell_j = 0$ for all but a small number (one or two, say) of coordinates $j = 1, \dots, d$ or simply picking at random a desired number, say m , of indices $\ell = (\ell_1, \dots, \ell_d)$.

The set of control variates at our disposal is constructed as follows. Let $K = 12$ and for $k \in \{1, \dots, K\}$, let $h_k(x) = P_k(2x - 1)$ for $x \in [0, 1]$, with P_k the univariate Legendre polynomial (Legendre function of the first kind) of degree k . Because the Legendre polynomials are orthogonal, they provide some numerical stability when inverting the Gram matrix. Let $m_{\max} = 2000$ and let d be such that $Kd \leq m_{\max}$. The first Kd control variates gather all the Legendre polynomials seen as tensor products but with $\ell_j = 0$ for all but a single coordinate. The $m_{\max} - Kd$ others are chosen at random uniformly over the remaining tensor products. The set of control variates constructed in this way is fixed during the whole study.

Integrands. We consider three integrands on $[0, 1]^d$:

$$\begin{aligned} f_1(x_1, \dots, x_d) &= (2/\pi)^{1/2} x_1^{-1} \exp(-(\log x_1)^2/2), \\ f_2(x_1, \dots, x_d) &= (2/\pi)^{d/2} \prod_{j=1}^d \{x_j^{-1} \exp(-(\log x_j)^2/2)\}, \\ f_3(x_1, \dots, x_d) &= 1 + \sin\left(\pi\left(\frac{2}{d} \sum_{j=1}^d x_j - 1\right)\right). \end{aligned}$$

All three functions integrate to 1 on $[0, 1]^d$. The function f_1 depends on the first coordinate only. In contrast, f_2 and f_3 represent more difficult situations. None of the three integrands belongs to the linear span of the control variates constructed in the previous paragraph.

Methods in competition. Besides the methods in Section 2, we also consider methods where the least squares estimator is computed via Principal Component Regression (PCR) or Ridge regression. Together, the methods in competition are thus OLSMC, Principal Component Regression Monte Carlo (PCRMC), Ridge regression Monte Carlo (RidgeMC), LASSOMC, LSLASSOMC with variable selection on the full sample and LSLASSOMCX, which is the same as LSLASSOMC but with a subsampling strategy when computing the active set.

The OLSMC, LASSOMC, RidgeMC and LSLASSOMC/X have been computed using the *sklearn* library [13]. To select the regularization parameter, we use standard cross-validation with the number of folds equal to 3 and the following parameter grids: 20 values from 1e-5 to 1e-1 for the LASSO variations and from 1e-5 to 1e2 for Ridge. The size of the subsample is $N = \lfloor 3\sqrt{n} \rfloor$. This choice accelerates the computation in a substantial manner without deteriorating too much the support recovery property of the LASSO. The implementation of PCRMC is not standard. The number of components selected by the Principal Component Analysis is set equal to the number of active variables given by the LASSO.

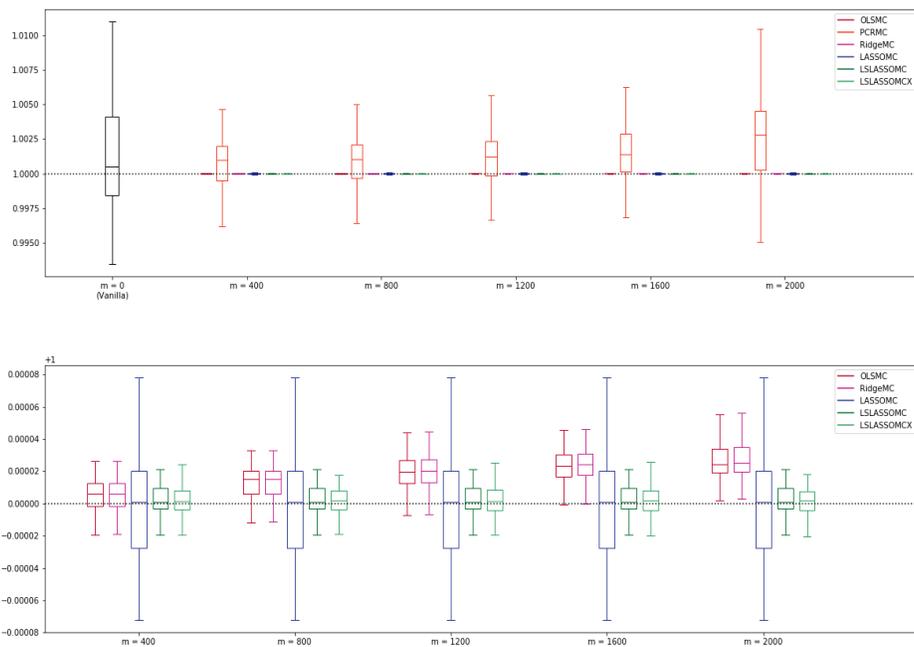


Figure 4.1: Boxplots (based on 50 replications) of the values returned by each of the methods (top) and zooming on the best ones (bottom) for f_1 . The dimension is $d = 5$, the sample size is $n = 5000$ and m (horizontal axis) varies from 400 to 2000.

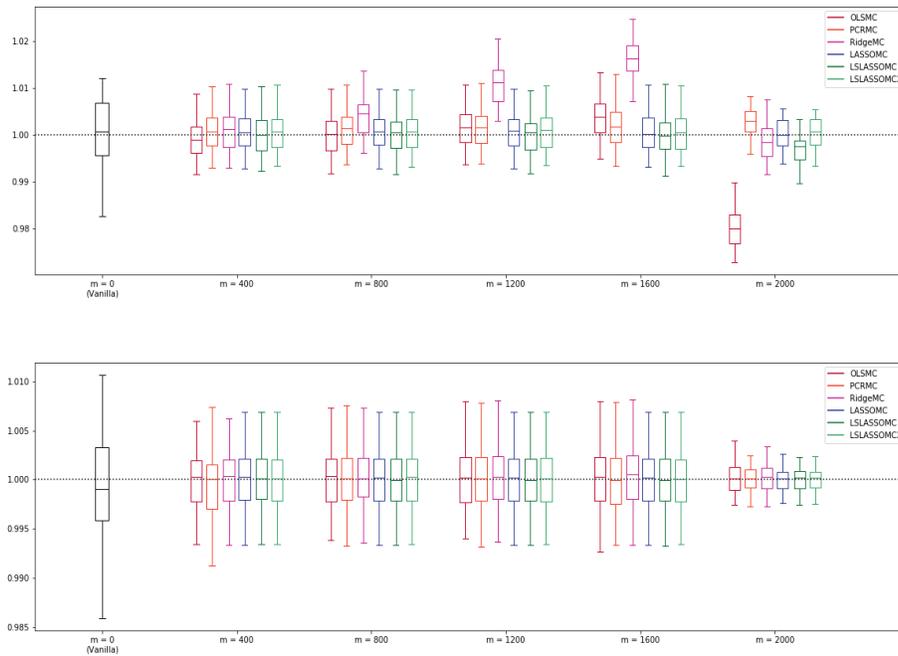


Figure 4.2: Boxplots (based on 50 replications) of the values returned by each of the methods for f_2 (top) and f_3 (bottom). The dimension is $d = 5$, the sample size is $n = 10\,000$ and m (horizontal axis) varies from 400 to 2000.

Parameter configuration. The following configurations of (n, d, m) are considered: $d \in \{5, 10\}$, $n \in \{5\,000, 10\,000\}$, and $m \in \{400, 800, 1\,200, 1\,600, 2\,000\}$. The case $d = 10$ represents a difficult situation as the number of points n is relatively small compared to the dimension. For instance, a grid made of only 3 points in each direction would already comprise 59 049 points.

Results. The figures presented in the paper deal with the case $d = 5$. The corresponding figures for $d = 10$ are given in the Appendix.

In Figure 4.1, boxplots of the values returned by each of the methods are provided for f_1 when $d = 5$ and $n = 5\,000$. The bottom panel zooms in on the most accurate methods. The clear winners are the LSLASSOMC and the LSLASSOMCX. The LASSO variable selection was very stable: almost always the same set of active variables was selected. Whereas the number of sample points used in the selection step of LSLASSOMCX has been reduced drastically compared to the LSLASSOMC (from n to $\lfloor 3\sqrt{n} \rfloor$), the stability of the active set is barely attenuated. Accordingly, the error distributions for LSLASSOMC and LSLASSOMCX are quite similar. In contrast, PCRMC performs quite poorly because the construction of the principal components is done regardless of the integrand and tends to discard information that is carried by relevant control variates.

In Figure 4.2, boxplots of the values returned by each of the methods are provided for f_2 and f_3 when $d = 5$ and $n = 10\,000$. Even if neither function exhibits any sparsity (which, as for f_1 , would favor the LASSO), the three LASSO-based methods are among the most accurate ones. Because f_2 and f_3 are symmetric in their arguments (in contrast to f_1), the PCRMC shows a reasonable performance. The traditional cross-validation approach for the LSLASSOMC tends to select too many control variates, while for the LSLASSOMCX, due to subsampling, it selects a smaller number of variables. This explains the excellent performance of LSLASSOMCX for these examples.

A Auxiliary results

For a function f on S , write $\|f\|_\infty = \sup_{x \in S} |f(x)|$. Let \mathbb{P} and \mathbb{E} denote the probability measure and the corresponding expectation operator on the probability space carrying the random variables X_i .

Lemma 1. *Let X_1, \dots, X_n be independent and identically distributed random variables with distribution P . Let $\varphi_1, \dots, \varphi_m$ be real-valued functions such that $P(\varphi_k) = 0$ and $\varphi_k \in \mathcal{G}(\tau^2)$ for all $k = 1, \dots, m$. Then for all $\delta > 0$, we have with probability at least $1 - \delta$,*

$$\max_{1 \leq k \leq m} \left| \sum_{i=1}^n \varphi_k(X_i) \right| \leq \sqrt{2n\tau^2 \log(2m/\delta)}.$$

Proof. For each $k = 1, \dots, m$, the centered random variable $\sum_{i=1}^n \varphi_k(X_i)$ is sub-Gaussian with variance factor $n\tau^2$. By the union bound and by Chernoff's inequality, we have, for each $t > 0$,

$$\begin{aligned} \mathbb{P} \left(\max_{1 \leq k \leq m} \left| \sum_{i=1}^n \varphi_k(X_i) \right| > t \right) &\leq \sum_{k=1}^m \mathbb{P} \left(\left| \sum_{i=1}^n \varphi_k(X_i) \right| > t \right) \\ &\leq 2m \exp \left(\frac{-t^2}{2n\tau^2} \right). \end{aligned}$$

Set $t = \sqrt{2n\tau^2 \log(2m/\delta)}$ to find the result. \square

Let $\lambda_{\min}(A)$ denote the smallest eigenvalue of the symmetric matrix A .

Lemma 2. *Let X_1, \dots, X_n be independent and identically distributed random variables with distribution P . Let $h = (h_1, \dots, h_m)^T \in L^2(P)^m$ be such that the $m \times m$ Gram matrix $G = P(hh^T)$ satisfies $\lambda_{\min}(G) > 0$. Let $\delta, \eta \in (0, 1)$. If the integers m and n are such that $1 \leq m < n$ and*

$$\|h^T G^{-1} h\|_\infty \leq \frac{n\eta^2}{32m + 4 \log(1/\delta)}, \quad (\text{A.1})$$

then with probability at least $1 - \delta$, the empirical Gram matrix $\hat{G}_n = P_n(hh^T)$ satisfies

$$\lambda_{\min}(\hat{G}_n) > (1 - \eta)\lambda_{\min}(G).$$

Proof. Suppose that the result is true when G is the identity matrix. Then it would be possible to apply the result to the vector of functions $\tilde{h} = G^{-1/2}h$, whose Gram matrix is the identity matrix. We would get that $\lambda_{\min}(P_n(\tilde{h}\tilde{h}^T)) > 1 - \eta$ with probability at least $1 - \delta$. Since $P_n(\tilde{h}\tilde{h}^T) = G^{-1/2}\hat{G}_n G^{-1/2}$ and since $u^T G^{-1} u \leq 1/\lambda_{\min}(G)$ for every unit vector $u \in \mathbb{R}^m$, we have

$$\begin{aligned} \lambda_{\min} \left(P_n(\tilde{h}\tilde{h}^T) \right) &= \min_{u^T u = 1} \left\{ u^T P_n(\tilde{h}\tilde{h}^T) u \right\} = \min_{u^T u = 1} \left\{ \frac{(G^{-1/2}u)^T \hat{G}_n G^{-1/2}u}{(G^{-1/2}u)^T G^{-1/2}u} u^T G^{-1} u \right\} \\ &\leq \lambda_{\min}(\hat{G}_n) / \lambda_{\min}(G). \end{aligned}$$

It would then follow that

$$\lambda_{\min}(\hat{G}_n) \geq \lambda_{\min}(P_n(\tilde{h}\tilde{h}^T)) \lambda_{\min}(G) \geq (1 - \eta)\lambda_{\min}(G).$$

Hence we only need to show the result for $G = I$. Write $\hat{\lambda} = \lambda_{\min}(\hat{G}_n)$. By [21, Theorem 2.2], we have

$$\hat{\lambda} \geq 1 - 4C\sqrt{m/n} + CZ/\sqrt{n},$$

where Z is a centered random variable that satisfies the lower-tail bound

$$\forall t > 0, \quad \mathbb{P}(Z \leq -t) \leq \exp(-t^2/2) \quad (\text{A.2})$$

and where

$$C^2 = \sup_{v \in \mathbb{R}^m: v^T v = 1} P(|v^T h|^4).$$

It follows that

$$\begin{aligned} \mathbb{P}(\hat{\lambda} > 1 - \eta) &\geq \mathbb{P}[1 - 4C\sqrt{m/n} + CZ/\sqrt{n} > 1 - \eta] \\ &= \mathbb{P}[Z > -(\sqrt{m}\eta/C - 4\sqrt{m})]. \end{aligned} \quad (\text{A.3})$$

Write $\kappa = \|h^T h\|_\infty$. For $v \in \mathbb{R}^m$ such that $v^T v = 1$, we have, by the Cauchy–Schwarz inequality,

$$|v^T h|^4 = |v^T h|^2 |v^T h|^2 \leq (v^T v)(h^T h)(v^T h h^T v) \leq \kappa(v^T h h^T v).$$

In view of the isotropy of h , we find, again for all $v \in \mathbb{R}^m$ such that $v^T v = 1$,

$$P(|v^T h|^4) \leq \kappa P(v^T h h^T v) = \kappa v^T P(h h^T) v = \kappa v^T v = \kappa.$$

As a consequence, also $C^2 \leq \kappa$. Condition (A.1) implies that $n\eta^2/C^2 \geq n\eta^2/\kappa > 16m$ and thus $\sqrt{n}\eta/C - 4\sqrt{m} \geq \sqrt{n}\eta/\sqrt{\kappa} - 4\sqrt{m} > 0$. The bounds (A.2) and (A.3) yield

$$\begin{aligned} \mathbb{P}(\hat{\lambda} > 1 - \eta) &\geq 1 - \exp\{-(\sqrt{n}\eta/C - 4\sqrt{m})^2/2\} \\ &\geq 1 - \exp\{-(\sqrt{n}\eta/\sqrt{\kappa} - 4\sqrt{m})^2/2\}. \end{aligned}$$

A sufficient condition for $\mathbb{P}(\hat{\lambda} > 1 - \eta) \geq 1 - \delta$ is therefore that

$$\exp\{-(\sqrt{n}\eta/\sqrt{\kappa} - 4\sqrt{m})^2/2\} \leq \delta,$$

which is in turn equivalent to

$$\sqrt{n}\eta/\sqrt{\kappa} - 4\sqrt{m} \geq \sqrt{2 \log(1/\delta)}$$

and thus to

$$n\eta^2/\kappa \geq [4\sqrt{m} + \sqrt{2 \log(1/\delta)}]^2.$$

This criterion coupled with the elementary inequality $(a + b)^2 \leq 2a^2 + 2b^2$ produces the inequality (A.1) as sufficient condition. \square

Lemma 3. *Let (X, Y) be a pair of centered and uncorrelated random variables. If $X \in \mathcal{G}(\nu)$ and if $|Y| \leq \kappa$, where $\nu > 0$ and $\kappa > 0$, then $XY \in \mathcal{G}(8\kappa^2\nu)$.*

Proof. The proof is based upon a refinement of [2, Theorem 2.1]. Without loss of generality, suppose $\nu = 1 = \kappa$; for the general case, consider the variables $X/\sqrt{\nu}$ and Y/κ .

Let $\lambda \in \mathbb{R}$. Let $(X_1, Y_1), (X_2, Y_2)$ be two independent copies of (X, Y) . Since XY is centered too, we have $\mathbb{E}[e^{-\lambda XY}] \geq 1$ and thus

$$\begin{aligned} \mathbb{E}[e^{\lambda XY}] &\leq \mathbb{E}[e^{\lambda XY}] \mathbb{E}[e^{-\lambda XY}] \\ &= \mathbb{E}[e^{\lambda(X_1 Y_1 - X_2 Y_2)}] = \sum_{q=0}^{\infty} \frac{\lambda^{2q}}{(2q)!} \mathbb{E}[(X_1 Y_1 - X_2 Y_2)^{2q}]. \end{aligned}$$

Note that the odd moments in the series expansion vanish since $X_1 Y_1 - X_2 Y_2$ is symmetric.

To show that $\mathbb{E}[e^{\lambda XY}] \leq e^{4\lambda^2}$ for all λ , it is sufficient to show that, for all $q = 0, 1, 2, \dots$,

$$\frac{\mathbb{E}[(X_1 Y_1 - X_2 Y_2)^{2q}]}{(2q)!} \leq \frac{4^q}{q!}. \quad (\text{A.4})$$

We treat the cases $q = 0, 1, 2$ and $q \geq 3$ separately. A useful inequality will be that, since $X \in \mathcal{G}(1)$, for nonnegative integer q , by [2, Theorem 2.1]

$$\mathbb{E}[X^{2q}] \leq 2^{q+1} q!. \quad (\text{A.5})$$

Moreover, $\mathbb{E}[X^2] \leq \nu = 1$ [2, Exercise 2.16].

- For $q = 0$ there is nothing to show.

- For $q = 1$, we use $\mathbb{E}[X_1 Y_1 X_2 Y_2] = \mathbb{E}[X_1 Y_1] \mathbb{E}[X_2 Y_2] = 0$ and $|Y| \leq \kappa = 1$ to find

$$\mathbb{E}[(X_1 Y_1 - X_2 Y_2)^2] = \mathbb{E}[(X_1 Y_1)^2] + \mathbb{E}[(X_2 Y_2)^2] = 2 \mathbb{E}[(XY)^2] \leq 2 \mathbb{E}[X^2] \leq 2.$$

- For $q = 2$, we use again the fact that the variables $X_1 Y_1$ and $X_2 Y_2$ are independent, identically distributed and centered to get that

$$\mathbb{E}[(X_1 Y_1 - X_2 Y_2)^4] = 2 \mathbb{E}[X^4 Y^4] + 6 \mathbb{E}[X^2 Y^2] \leq 2 \mathbb{E}[X^4] + 6 \leq 2 \cdot 16 + 6 = 38.$$

- For $q \geq 3$, we have, by convexity of the function $x \mapsto x^{2q}$ and by (A.5), that

$$\begin{aligned} \mathbb{E}[(X_1 Y_1 - X_2 Y_2)^{2q}] &\leq 2^{2q-1} (\mathbb{E}[(X_1 Y_1)^{2q}] + \mathbb{E}[(-X_2 Y_2)^{2q}]) \\ &= 2^{2q} \mathbb{E}[(XY)^{2q}] \leq 4^q \mathbb{E}[X^{2q}] \leq 4^q 2^{q+1} q!. \end{aligned}$$

Hence, the inequality (A.4) for integer $q \geq 3$ follows provided that $2^{q+1}(q!)^2 \leq (2q)!$ for all such q . But this is true since, for all integer $q \geq 3$, we have

$$\frac{(2q)!}{(q!)^2} = \prod_{j=1}^q \frac{q+j}{j} = (q+1) \prod_{j=2}^q \frac{q+j}{j} \geq 4 \cdot 2^{q-1} = 2^{q+1}.$$

We have thus verified (A.4) for all integer $q \geq 0$, and thus $\mathbb{E}[e^{\lambda XY}] \leq e^{4\lambda^2} = e^{8\lambda^2/2}$, as required. \square

B Proof of Theorem 3.1

The proof is divided into several steps. In one of them, we use a non-probabilistic property of the OLS estimate

$$\hat{\beta}_n^{\text{ols}}(f) \in \arg \min_{\beta \in \mathbb{R}^m} \|f_c^{(n)} - H_c \beta\|_2^2, \quad (\text{B.1})$$

where $f_c^{(n)} = f^{(n)} - \mathbf{1}_n(\mathbf{1}_n^T f^{(n)})/n$ and $H_c = H - \mathbf{1}_n(\mathbf{1}_n^T H)/n$.

Lemma 4. *If there exists $\nu > 0$ such that $\|H_c u\|_2^2 \geq \nu \|u\|_2^2$ for all $u \in \mathbb{R}^m$, then the minimizer $\hat{\beta}_n^{\text{ols}}(f)$ in (B.1) is unique and*

$$\|\hat{\beta}_n^{\text{ols}}(f) - \beta^*(f)\|_2 \leq \frac{\sqrt{m}}{\nu n} \max_{k=1, \dots, m} |H_{c,k}^T \epsilon_c^{(n)}|, \quad (\text{B.2})$$

where $\epsilon_c^{(n)} = f_c^{(n)} - H_c \beta^*(f)$ and where $H_{c,k}$ is the k -th column of H_c .

Proof. The matrix $H_c^T H_c$ is invertible, since its smallest eigenvalue is bounded from below by ν . The OLS estimate is thus unique and given by

$$\begin{aligned} \hat{\beta}_n^{\text{ols}}(f) &= (H_c^T H_c)^{-1} H_c^T f_c^{(n)} \\ &= (H_c^T H_c)^{-1} H_c^T (H_c \beta^*(f) + \epsilon_c^{(n)}) \\ &= \beta^*(f) + (H_c^T H_c)^{-1} H_c^T \epsilon_c^{(n)}. \end{aligned}$$

The largest eigenvalue of $(H_c^T H_c)^{-1}$ being bounded from above by $(\nu)^{-1}$, we obtain

$$\|\hat{\beta}_n^{\text{ols}}(f) - \beta^*(f)\|_2 = \|(H_c^T H_c)^{-1} H_c^T \epsilon_c^{(n)}\|_2 \leq \frac{1}{\nu} \|H_c^T \epsilon_c^{(n)}\|_2$$

Since $\|x\|_2 \leq \sqrt{m} \max_{k=1, \dots, m} |x_k|$ for $x \in \mathbb{R}^m$, we can conclude. \square

Step 1. — Since $f = P(f) + \beta^*(f)^T h + \epsilon$, the oracle estimate of $P(f)$, which uses the unknown, optimal coefficient vector $\beta^*(f)$, is

$$\hat{\alpha}_n^{\text{or}}(f) = P_n[f - \beta^*(f)^T h] = P(f) + P_n(\epsilon).$$

The difference between the OLS and oracle estimates is

$$\hat{\alpha}_n^{\text{ols}}(f) - \hat{\alpha}_n^{\text{or}}(f) = \left(\hat{\beta}_n^{\text{ols}}(f) - \beta^*(f) \right)^T P_n(h).$$

The estimation error of the OLS estimator can thus be decomposed as

$$\begin{aligned} n(\hat{\alpha}_n^{\text{ols}}(f) - P(f)) &= n(\hat{\alpha}_n^{\text{or}}(f) - P(f)) + \left(\beta^*(f) - \hat{\beta}_n^{\text{ols}}(f) \right)^T n P_n(h) \\ &= \sum_{i=1}^n \epsilon(X_i) + \left(\beta^*(f) - \hat{\beta}_n^{\text{ols}}(f) \right)^T \sum_{i=1}^n h(X_i). \end{aligned}$$

By the triangle and Cauchy–Schwarz inequalities and since $\|x\|_2 \leq \sqrt{m} \max_{k=1, \dots, m} |x_k|$ for $x \in \mathbb{R}^m$, we get

$$n |\hat{\alpha}_n^{\text{ols}}(f) - P(f)| \leq \left| \sum_{i=1}^n \epsilon(X_i) \right| + \|\beta^*(f) - \hat{\beta}_n^{\text{ols}}(f)\|_2 \sqrt{m} \max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \right|. \quad (\text{B.3})$$

To proceed, we will construct an event that has probability at least $1 - \delta$ and on which we can control each of the three terms on the right-hand side simultaneously (Step 2). Most difficult to treat will be the term $\|\beta^*(f) - \hat{\beta}_n^{\text{ols}}(f)\|_2$ (Step 3). Collecting all the inequalities, we will arrive at the stated bound (Step 4).

Step 2. — Let $\delta > 0$ and $n \geq 1$. We construct an event with probability at least $1 - \delta$ on which four inequalities hold simultaneously.

- The empirical Gram matrix of the vector $h = (h_1, \dots, h_m)^T \in L_2(P)$ based on the sample X_1, \dots, X_n is $n^{-1} H^T H$. By Lemma 2 with $\eta = 1/2$, because $4\|h^T G^{-1} h\|_\infty \leq n/(32m + 4 \log(4/\delta))$ by assumption, we have with probability at least $1 - \delta/4$,

$$\|Hu\|_2^2 \geq n\gamma \|u\|_2^2/2. \quad (\text{B.4})$$

- By virtue of Assumptions 1 and 2, we can apply Lemma 3 to find $h_k \epsilon \in \mathcal{G}(C\tau^2 U^2)$ with $C = 8$ ¹. Hence we can apply Lemma 1 to get that, with probability at least $1 - \delta/4$,

$$\max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \epsilon(X_i) \right| \leq \sqrt{2nC\tau^2 U^2 \log(8m/\delta)}. \quad (\text{B.5})$$

- In view of [2, Lemma 2.2] and Assumption 2, we have $h_k \in \mathcal{G}(U^2)$ for all $k = 1, \dots, m$. Hence we can apply Lemma 1 to get that, with probability at least $1 - \delta/4$,

$$\max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \right| \leq \sqrt{2nU^2 \log(8m/\delta)}, \quad (\text{B.6})$$

- Finally, because $\epsilon \in \mathcal{G}(\tau^2)$, with probability at least $1 - \delta/4$,

$$\left| \sum_{i=1}^n \epsilon(X_i) \right| \leq \sqrt{2n\tau^2 \log(8/\delta)}. \quad (\text{B.7})$$

¹We use a generic C so as to allow for easy modifications in case sharper constants can be found.

By the union bound, the event on which the inequalities (B.4), (B.5), (B.6), and (B.7) are all satisfied simultaneously has probability at least $1 - \delta$. For the remainder of the proof, we work on this event, denoted by E .

Step 3. — We will show that, on the event E constructed in Step 2, we have

$$\|\hat{\beta}_n^{\text{ols}}(f) - \beta^*(f)\|_2 \leq \frac{4\sqrt{m}}{\gamma n} \sqrt{2nC\tau^2 U^2 \log(8m/\delta)} \left(1 + \sqrt{(2/C) \log(8/\delta)/n}\right). \quad (\text{B.8})$$

To do so, we will apply Lemma 4 with $\nu = \gamma/4$, but we need to show first that the condition on H_c is satisfied (Step 3.1). Then, we will control the right-hand side in (B.2) (Step 3.2). Inequality (B.10) below together with Lemma 4 with $\nu = \gamma/4$ will then yield (B.8).

Step 3.1. — On the event E , we have

$$\|H_c u\|_2^2 \geq \frac{n\gamma}{4} \|u\|_2^2, \quad \forall u \in \mathbb{R}^m. \quad (\text{B.9})$$

To see why, first note that $n^{-1}H_c^T H_c = n^{-1}H^T H - P_n(h)P_n(h)^T$. The Cauchy–Schwarz inequality gives

$$\|H_c u\|_2^2 = \|Hu\|_2^2 - n(P_n(h)^T u)^2 \geq \|Hu\|_2^2 - n\|P_n(h)\|_2^2 \|u\|_2^2.$$

In view of (B.6), we also have

$$\|P_n(h)\|_2^2 \leq \frac{m}{n^2} \max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \right|^2 \leq \frac{m}{n^2} 2nU^2 \log(8m/\delta) \leq \frac{\gamma}{4},$$

as $n \geq 8m \log(8m/\delta) U^2 / \gamma$ by assumption. In view of (B.4), we get

$$\|H_c u\|_2^2 \geq \frac{n\gamma}{2} \|u\|_2^2 - \frac{n\gamma}{4} \|u\|_2^2 = \frac{n\gamma}{4} \|u\|_2^2.$$

This shows (B.9) with $\nu = \gamma/4$.

Step 3.2. — On the event E , we have

$$\max_{k=1, \dots, m} |H_{c,k}^T \epsilon_c^{(n)}| \leq \sqrt{2nC\tau^2 U^2 \log(8m/\delta)} \left(1 + \sqrt{(2/C) \log(8/\delta)/n}\right). \quad (\text{B.10})$$

Indeed, the left-hand side in (B.10) is

$$\begin{aligned} & \max_{k=1, \dots, m} \left| \sum_{i=1}^n (h_k(X_i) - P_n(h_k)) (\epsilon(X_i) - P_n(\epsilon)) \right| \\ &= \max_{k=1, \dots, m} \left| \left(\sum_{i=1}^n h_k(X_i) \epsilon(X_i) \right) - n P_n(h_k) P_n(\epsilon) \right| \\ &\leq \max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \epsilon(X_i) \right| + n^{-1} \left| \sum_{i=1}^n \epsilon(X_i) \right| \max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \right| \\ &\leq \sqrt{2nC\tau^2 U^2 \log(8m/\delta)} + n^{-1} \sqrt{2n\tau^2 \log(8/\delta)} \sqrt{2nU^2 \log(8m/\delta)}. \end{aligned}$$

The right-hand side can be simplified to the bound in (B.10).

Step 4. — The three terms in the bound (B.3) on the estimation error of the OLS estimate can be controlled by inequalities (B.6), (B.7), and (B.8), all of which hold on the event E . We find

$$\begin{aligned} & n |\hat{\alpha}_n^{\text{ols}}(f) - P(f)| \\ &\leq \sqrt{2n\tau^2 \log(8/\delta)} + \frac{4m}{\gamma n} \sqrt{2nC\tau^2 U^2 \log(8m/\delta)} \left(1 + \sqrt{\frac{2 \log(8/\delta)}{Cn}}\right) \sqrt{2nU^2 \log(8m/\delta)} \\ &= \sqrt{2 \log(8/\delta)} \tau \sqrt{n} + 8\sqrt{C} \gamma^{-1} \tau U^2 m \log(8m/\delta) \left(1 + \sqrt{\frac{2 \log(8/\delta)}{Cn}}\right). \end{aligned}$$

Divide by n and plug in $C = 8$ and $\zeta_h = U^2/\gamma$ to find

$$|\hat{\alpha}_n^{\text{ols}}(f) - P(f)| \leq \sqrt{2 \log(8/\delta)} \tau \sqrt{n} + 16m \sqrt{2 \log(8m/\delta)} \left(1 + \frac{1}{2} \sqrt{\log(8/\delta)/n}\right) \zeta_h \frac{\tau}{n}.$$

The smallest eigenvalue of G is certainly bounded by the mean of its eigenvalues. It follows that $\gamma \leq m^{-1} \sum_{j=1}^m P(h_j^2) \leq U^2$ and thus $\zeta_h \geq 1$. The condition $n \geq 8m \log(8m/\delta) \zeta_h$ thus implies $\log(8/\delta)/n \leq 1/(8m)$ and therefore

$$\frac{1}{2} \sqrt{\log(8/\delta)/n} \leq \frac{1}{2\sqrt{8m}} = \frac{1}{4\sqrt{2m}} \leq \frac{1}{4\sqrt{2}}.$$

Substitute this into the bound on the OLS estimation error to find

$$|\hat{\alpha}_n^{\text{ols}}(f) - P(f)| \leq \sqrt{2 \log(8/\delta)} \tau \sqrt{n} + 16 \left(\sqrt{2} + \frac{1}{4}\right) m \log(8m/\delta) \zeta_h \frac{\tau}{n},$$

and use that $16(\sqrt{2} + 1/4) \leq 27$ to obtain the final inequality. \square

C Proof of Theorem 3.2

Any norm $\|\cdot\|$ on the Euclidean space \mathbb{R}^p can be extended to a matrix norm $\|\cdot\|$ on $\mathbb{R}^{p \times p}$ via $\|A\| := \sup_{\|u\|=1} \|Au\|$. For any vector $\beta \in \mathbb{R}^m$ and any ordered set $S = (k_1, \dots, k_\ell) \subset \{1, \dots, m\}$, let $\beta_S = (\beta_{k_1}, \dots, \beta_{k_\ell})^T$. For any matrix $A \in \mathbb{R}^{n \times m}$ and $k \in \{1, \dots, m\}$, denote by A_k its k -th column. For any matrix $A \in \mathbb{R}^{n \times m}$ and any ordered set $S = (k_1, \dots, k_\ell) \subset \{1, \dots, m\}$, let $A_S = (A_{k_1}, \dots, A_{k_\ell})$.

Recall that $S^* = \{j = 1, \dots, m : \beta_j^*(f) \neq 0\}$ is the set of active control variates, of which there are $\ell^* = |S^*|$. Further, recall that the LASSO estimator of the coefficient vector is

$$\hat{\beta}_n^{\text{lasso}}(f) = \arg \min_{\beta \in \mathbb{R}^m} \frac{1}{2n} \|f_c^{(n)} - H_c \beta\|_2^2 + \lambda \|\beta\|_1, \quad (\text{C.1})$$

where $f_c^{(n)} = f^{(n)} - \mathbb{1}_n (\mathbb{1}_n^T f^{(n)})/n$ and $H_c = H - \mathbb{1}_n (\mathbb{1}_n^T H)/n$ and where $\lambda > 0$ is a regularization parameter. Its support set is $S(\hat{\beta}_n^{\text{lasso}}(f)) = \{k = 1, \dots, m : \hat{\beta}_{n,k}^{\text{lasso}}(f) \neq 0\}$, and these are variables selected by the LASSO. The vector of centered (and observable) errors is

$$\epsilon_c^{(n)} = f_c^{(n)} - H_c \beta^*(f).$$

Our treatment follows from the one exposed in [19, Chapter 11], except that we pay special attention to the randomness of the design matrix H_c , the properties of which follow from the ones of the control variates.

Step 1. — We first establish some (non-probabilistic) properties of $\hat{\beta}_n^{\text{lasso}}(f)$. To this end, we consider the linear regression of the non-active control variates on the active ones: for $k \in S^* = \{j = 1, \dots, m : \beta_j^*(f) = 0\}$, this produces the coefficient vector

$$\hat{\theta}_n^{(k)} \in \arg \min_{\theta \in \mathbb{R}^{\ell^*}} \|H_{c,k} - H_{c,S^*} \theta\|_2.$$

Further, we consider the OLS oracle estimate $\hat{\beta}_n^*$, which is the OLS estimator based upon the active control variables only, i.e.,

$$\hat{\beta}_n^* \in \arg \min_{\beta \in \mathbb{R}^{\ell^*}} \|f_c^{(n)} - H_{c,S^*} \beta\|_2.$$

Our assumptions will imply that, with large probability, H_{c,S^*} has rank ℓ^* , in which case

$$\begin{aligned} \hat{\theta}_n^{(k)} &= (H_{c,S^*}^T H_{c,S^*})^{-1} H_{c,S^*}^T H_{c,k}, \\ \hat{\beta}_n^* &= (H_{c,S^*}^T H_{c,S^*})^{-1} H_{c,S^*}^T f_c^{(n)}. \end{aligned}$$

The following lemma provides a number of (non-probabilistic) properties of $\hat{\beta}_n^{\text{lasso}}(f)$, given certain conditions on H_c and $\epsilon_c^{(n)}$.

Lemma 5. *If H_{c,S^*} has rank ℓ^* and if there exists $\kappa \in (0, 1]$ such that*

$$\max_{k \in \overline{S^*}} \|\hat{\theta}_n^{(k)}\|_1 \leq 1 - \kappa, \quad (\text{C.2})$$

$$\max_{k \in \overline{S^*}} |(H_{c,k} - H_{c,S^*} \hat{\theta}_n^{(k)})^T \epsilon_c^{(n)}| \leq \kappa \lambda n, \quad (\text{C.3})$$

then the minimizer $\hat{\beta}_n^{\text{lasso}}(f)$ in (C.1) is unique, with support $S(\hat{\beta}_n^{\text{lasso}}(f)) \subset S^*$, and it satisfies

$$\max_{k \in \overline{S^*}} |\hat{\beta}_{n,k}^{\text{lasso}}(f) - \beta_k^*(f)| \leq \max_{k \in \overline{S^*}} |\hat{\beta}_{n,k}^* - \beta_k^*(f)| + n\lambda \|(H_{c,S^*}^T H_{c,S^*})^{-1}\|_\infty. \quad (\text{C.4})$$

Proof. The proof of the previous result is actually contained in [19]. The uniqueness of the LASSO solution and the property that it does not select inactive covariates follows directly from the proof of their Theorem 11.3. The only difference is that, in our case, the inequality (C.3) is an assumption whereas in [19] it is a property of the Gaussian fixed design model. The approach in [19] is based upon checking the *strict dual feasibility condition*. The bound (C.4) is Eq. (11.37) in [19]. \square

We slightly modify Lemma 5 to make the conditions (C.2) and (C.3) easier to check and to make the bound (C.4) easier to use.

Lemma 6. *If there exists $\nu > 0$ such that*

$$\|H_{c,S^*} u\|_2^2 \geq n\nu \|u\|_2^2, \quad \forall u \in \mathbb{R}^{\ell^*}, \quad (\text{C.5})$$

and if there exists $\kappa \in (0, 1]$ such that

$$\frac{\ell^*}{\nu n} \max_{k \in \overline{S^*}} \max_{j \in \overline{S^*}} |H_{c,j}^T H_{c,k}| \leq 1 - \kappa, \quad (\text{C.6})$$

$$\max_{k=1,\dots,m} |H_{c,k}^T \epsilon_c^{(n)}| \leq \frac{1}{2} \kappa \lambda n, \quad (\text{C.7})$$

then the minimizer $\hat{\beta}_n^{\text{lasso}}(f)$ in (C.1) is unique, with support $S(\hat{\beta}_n^{\text{lasso}}(f)) \subset S^*$, and it satisfies

$$\max_{k \in \overline{S^*}} |\hat{\beta}_{n,k}^{\text{lasso}}(f) - \beta_k^*(f)| \leq (1 + \kappa/2) \sqrt{\ell^*} \lambda / \nu. \quad (\text{C.8})$$

Proof. By (C.5), the smallest eigenvalue of the $\ell^* \times \ell^*$ matrix $H_{c,S^*}^T H_{c,S^*}$ matrix is positive, so that it is invertible and H_{c,S^*} has rank ℓ^* .

We show that (C.6) implies (C.2). For each $k \in \overline{S^*}$, the vector $\hat{\theta}_n^{(k)}$ has length ℓ^* , so that

$$\|\hat{\theta}_n^{(k)}\|_1 \leq \sqrt{\ell^*} \|\hat{\theta}_n^{(k)}\|_2.$$

Because $\hat{\theta}_n^{(k)}$ is an OLS estimate, we can apply Lemma 4 with $\theta^* = 0$ to get

$$\|\hat{\theta}_n^{(k)}\|_2 \leq \frac{\sqrt{\ell^*}}{\nu n} \max_{j \in \overline{S^*}} |H_{c,j}^T H_{c,k}|.$$

Combining the previous two bounds, we find that (C.6) indeed implies (C.2).

Next we show that (C.7) implies (C.3). For $k \in \overline{S^*}$, we have

$$\begin{aligned} |(H_{c,k} - H_{c,S^*} \hat{\theta}_n^{(k)})^T \epsilon_c^{(c)}| &\leq |H_{c,k}^T \epsilon_n^{(c)}| + |(\hat{\theta}_n^{(k)})^T H_{c,S^*}^T \epsilon_n^{(c)}| \\ &\leq |H_{c,k}^T \epsilon_n^{(c)}| + \|\hat{\theta}_n^{(k)}\|_1 \max_{j \in \overline{S^*}} |H_{c,j}^T \epsilon_n^{(c)}|. \end{aligned}$$

Using (C.2) and (C.7) we deduce (C.3).

The conditions of Lemma 5 have been verified, and so its conclusion holds. We simplify the two terms in the upper bound (C.4). First, we apply Lemma 4 to the oracle OLS estimator $\hat{\beta}_n^*$. Second, for any matrix $A \in \mathbb{R}^{p \times p}$, we have $\|A\|_\infty \leq \sqrt{p}\|A\|_2$ (e.g., [7, page 365]), and this we apply to $(H_{c,S^*}^T H_{c,S^*})^{-1}$. In this way, the upper bound in (C.4) is dominated by

$$\|\hat{\beta}_n^* - \beta^*(f)\|_2 + n\lambda \cdot \sqrt{\ell^*} \|(H_{c,S^*}^T H_{c,S^*})^{-1}\|_2 \leq \frac{\sqrt{\ell^*}}{n\nu} \max_{k \in S^*} |H_{c,k}^T \epsilon_c^{(n)}| + n\lambda \cdot \sqrt{\ell^*} \cdot \frac{1}{n\nu},$$

since the largest eigenvalue of $(H_{c,S^*}^T H_{c,S^*})^{-1}$ is at most $(n\nu)^{-1}$. Use (C.7) to further simplify the right-hand side, yielding (C.8). \square

Step 2. — Let $\delta \in (0, 1)$ and $n = 1, 2, \dots$. In a similar way as in the proof of Theorem 3.1, we construct an event of probability at least $1 - \delta$. This time, we need five inequalities to hold simultaneously.

- Because $n \geq 4\|h_{S^*}^T G_{S^*}^{-1} h_{S^*}\|_\infty (32\ell^* + 4\log(5/\delta))$ by assumption, we have, by Lemma 2 with $\eta = 1/2$, with probability at least $1 - \delta/5$,

$$\|H_{S^*} u\|_2^2 \geq n\gamma^* \|u\|_2^2 / 2, \quad \forall u \in \mathbb{R}^{\ell^*}. \quad (\text{C.9})$$

- By virtue of Assumptions 1 and 2, we have $h_k \epsilon \in \mathcal{G}(CU^2\tau^2)$, where $C = 8$. Hence we can apply Lemma 1 to get that, with probability at least $1 - \delta/5$,

$$\max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \epsilon(X_i) \right| \leq \sqrt{2Cn\tau^2 U^2 \log(10m/\delta)}. \quad (\text{C.10})$$

- Similarly, because $h_k \in \mathcal{G}(U^2)$, with probability at least $1 - \delta/5$,

$$\max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \right| \leq \sqrt{2nU^2 \log(10m/\delta)}, \quad (\text{C.11})$$

- Because $\epsilon \in \mathcal{G}(\tau^2)$, with probability at least $1 - \delta/5$,

$$\left| \sum_{i=1}^n \epsilon(X_i) \right| \leq \sqrt{2n\tau^2 \log(10/\delta)}. \quad (\text{C.12})$$

- Finally, because $P(h_k h_j) = 0$ and $\|h_k h_j\|_\infty \leq U^2$ for all $k \in \overline{S^*}$ and $j \in S^*$, we have $h_k h_j \in \mathcal{G}(U^4)$ for such k and j , and thus, with probability at least $1 - \delta/5$,

$$\max_{k \in \overline{S^*}} \max_{j \in S^*} \left| \sum_{i=1}^n h_k(X_i) h_j(X_i) \right| \leq \sqrt{2nU^4 \log(10\ell^* m/\delta)}. \quad (\text{C.13})$$

By the union bound, the event, say E , on which (C.9), (C.10), (C.11), (C.12) and C.13 are satisfied simultaneously has probability at least $1 - \delta$. We work on the event E for the rest of the proof.

Step 3. — On the event E , we have

$$\|H_{c,S^*} u\|_2^2 \geq n\alpha\gamma^* \|u\|_2^2, \quad \forall u \in \mathbb{R}^{\ell^*} \quad (\text{C.14})$$

where $\alpha \in (0, 1/2)$ is an absolute constant whose value will be fixed later in the proof.

The proof of (C.14) is similar to the one of (B.9). We have

$$H_{c,S^*}^T H_{c,S^*} = H_{S^*}^T H_{S^*} - n P_n(h_{S^*}) P_n(h_{S^*})^T$$

and thus, by the Cauchy–Schwarz inequality and by (C.9),

$$\begin{aligned}\|H_{c,S^*}u\|_2^2 &\geq \|H_{S^*}u\|_2^2 - n\|P_n(h_{S^*})\|_2^2\|u\|_2^2 \\ &\geq n(\gamma^*/2 - \|P_n(h_{S^*})\|_2^2)\|u\|_2^2.\end{aligned}$$

In view of (C.11), we have

$$\|P_n(h_{S^*})\|_2^2 \leq \frac{\ell^*}{n^2} \cdot 2nU^2 \log(10m/\delta) = 2\ell^* \log(10m/\delta)U^2/n.$$

We thus get

$$\|H_{c,S^*}u\|_2^2 \geq n\gamma^* \left[\frac{1}{2} - \frac{2\ell^* \log(10m/\delta)U^2/\gamma^*}{n} \right] \|u\|_2^2$$

A sufficient condition for (C.14) is thus that the term in square brackets is at least α , i.e.,

$$n \geq \frac{2}{1/2 - \alpha} \ell^* \log(10m/\delta)U^2/\gamma^*$$

Since $\ell^* \geq 1$ and $U^2 \geq \gamma^*$, a condition of the form

$$n \geq \rho(\ell^*)^2 \log(10\ell^*m/\delta)(U^2/\gamma^*)^2 \quad (\text{C.15})$$

is thus sufficient, provided $\rho > 2/(1/2 - \alpha)$. In Step 6(ii), we will choose α in such a way that the constant $\rho = 70$ appearing in the statement of the theorem is sufficient.

Step 4. — On the event E , we have

$$\max_{k \in \overline{S^*}} \max_{j \in S^*} |H_{c,j}^T H_{c,k}| \leq \sqrt{2nU^4 \log(10\ell^*m/\delta)} + 2U^2 \log(10m/\delta). \quad (\text{C.16})$$

Indeed, the left-hand side is bounded by

$$\begin{aligned}&\max_{k \in \overline{S^*}} \max_{j \in S^*} \left| \left(\sum_{i=1}^n h_k(X_i) h_j(X_i) \right) - nP_n(h_k)P_n(h_j) \right| \\ &\leq \max_{k \in \overline{S^*}} \max_{j \in S^*} \left| \sum_{i=1}^n h_k(X_i) h_j(X_i) \right| + \frac{1}{n} \max_{k \in \overline{S^*}} \left| \sum_{i=1}^n h_k(X_i) \right| \max_{j \in S^*} \left| \sum_{i=1}^n h_j(X_i) \right| \\ &\leq \max_{k \in \overline{S^*}} \max_{j \in S^*} \left| \sum_{i=1}^n h_k(X_i) h_j(X_i) \right| + \frac{1}{n} \max_{k=1, \dots, m} \left| \sum_{i=1}^n h_k(X_i) \right|^2 \\ &\leq \sqrt{2nU^4 \log(10\ell^*m/\delta)} + \frac{1}{n} \cdot 2nU^2 \log(10m/\delta),\end{aligned}$$

which is (C.16).

Step 5. — On the event E , we have

$$\max_{k=1, \dots, m} |H_{c,k}^T \epsilon_n^{(c)}| \leq \sqrt{2nC\tau^2 U^2 \log(10m/\delta)} \left(1 + \sqrt{(2/C) \log(10/\delta)/n} \right). \quad (\text{C.17})$$

The proof is the same as the one of (B.10).

Step 6. — We will verify that on the event E , the three assumptions of Lemma 6 are satisfied with $\kappa = 1/2$ and $\nu = \alpha\gamma^*$, with α as in Step 3. We will make use of the inequality²

$$\forall (a, b, c) \in (0, \infty)^3, \forall x \geq \sqrt{b^2 + 4ac}/a, \quad ax^2 \geq bx + c. \quad (\text{C.18})$$

²The convex parabola $x \mapsto ax^2 - bx - c$ has zeroes at $x_{\pm} = (b \pm \sqrt{b^2 + 4ac})/(2a)$, and $x_- < 0 < x_+ < \sqrt{b^2 + 4ac}/a$.

(i) Eq. (C.5) with $\nu = \alpha\gamma^*$ is just (C.14).

(ii) Eq. (C.6) with $\nu = \alpha\gamma^*$ and $\kappa = 1/2$ follows from (C.16) provided we have

$$\frac{\ell^*}{\alpha\gamma^*n} \cdot \left(\sqrt{2nU^4 \log(10\ell^*m/\delta)} + 2U^2 \log(10m/\delta) \right) \leq 1 - \frac{1}{2}.$$

To check whether the latter inequality is satisfied, we apply (C.18) with $x = \sqrt{n}$ and

$$\begin{aligned} a &= \alpha\gamma^*/(2\ell^*), \\ b &= \sqrt{2U^4 \log(10\ell^*m/\delta)}, \\ c &= 2U^2 \log(10m/\delta). \end{aligned}$$

Sufficient is that $n = x^2$ is bounded from below by $(b^2 + 4ac)/a^2 = (b/a)^2 + 4c/a$, which is

$$\begin{aligned} \frac{2U^4 \log(10\ell^*m/\delta)}{(\alpha\gamma^*/(2\ell^*))^2} + 4 \frac{2U^2 \log(10m/\delta)}{\alpha\gamma^*/(2\ell^*)} \\ = \frac{8}{\alpha^2} (\ell^*)^2 \log(10\ell^*m/\delta) \cdot (U^2/\gamma^*)^2 + \frac{16}{\alpha} \ell^* \log(10m/\delta) \cdot U^2/\gamma^*. \end{aligned}$$

But $\ell^* \geq 1$ and $\gamma^* \leq (1/\ell^*) \sum_{j \in S^*} P(h_j^2) \leq U^2$, so that a sufficient condition is that

$$n \geq \left(\frac{8}{\alpha^2} + \frac{16}{\alpha} \right) (\ell^*)^2 \log(10\ell^*m/\delta) \cdot (U^2/\gamma^*)^2.$$

The constant ρ in (C.15) must thus be such that

$$\rho \geq \max \left(\frac{2}{1/2 - \alpha}, \frac{8}{\alpha^2} + \frac{16}{\alpha} \right).$$

The minimum of the right-hand side as a function of $\alpha \in (0, 1/2)$ occurs at $\alpha = \sqrt{2}/3$ and is equal to $2/(1/2 - \sqrt{2}/3) \approx 69.94113$. Taking $\rho = 70$ as in the assumption on n is thus sufficient.

(iii) Eq. (C.7) with $\kappa = 1/2$ follows from (C.17), since (recall $C = 8$)

$$\sqrt{16n\tau^2 U^2 \log(10m/\delta)} \left(1 + \sqrt{\log(10/\delta)/(4n)} \right) \leq \lambda n/4$$

by the assumed lower bound on λ . Indeed, since $\ell^* \geq 1$ and $U^2 \geq \gamma^*$, the assumed lower bound for n implies that $n \geq 70 \log(10m/\delta)$ and thus

$$\frac{\log(10/\delta)}{4n} \leq \frac{\log(10/\delta)}{280 \log(10m/\delta)} \leq \frac{1}{280}.$$

Since $16 \cdot (1 + 1/\sqrt{280}) \approx 16.95618$, the assumed lower bound for λ suffices.

Step 7. — By the previous step, the conclusions of Lemma 6 with $\kappa = 1/2$ and $\nu = \alpha\gamma^*$ hold on the event E , where $\alpha = \sqrt{2}/3$ was specified in Step 6(ii). The minimizer $\hat{\beta}_n^{\text{lasso}}$ in (C.1) is thus unique and we have $S(\hat{\beta}_n^{\text{lasso}}) \subset S^*$.

To show the reverse inclusion, we need to verify that $|\hat{\beta}_{n,k}^{\text{lasso}}(f)| > 0$ for all $k \in S^*$. To this end, we apply (C.8) with $\kappa = 1/2$ and $\nu = \alpha\gamma^*$, which becomes

$$\max_{k \in S^*} |\hat{\beta}_{n,k}^{\text{lasso}}(f) - \beta_k^*(f)| \leq (5/4) \sqrt{\ell^*} \lambda / (\alpha\gamma^*).$$

For any $k \in S^*$, we thus have

$$|\hat{\beta}_{n,k}^{\text{lasso}}(f)| \geq \min_{j \in S^*} |\beta_j^*(f)| - (5/(4\alpha)) \sqrt{\ell^*} \lambda / \gamma^*.$$

But for $\alpha = \sqrt{2}/3$, we have $5/(4\alpha) \approx 2.6516$. As $\min_{j \in S^*} |\beta_j^*(f)| > 3\sqrt{\ell^*} \lambda / \gamma^*$ by the assumed upper bound for λ , we find $|\hat{\beta}_{n,k}^{\text{lasso}}(f)| > 0$, as required. \square

D Proof of Theorem 3.3

Recall $\hat{\beta}_n^{\text{lasso}}$ in (C.1). We start with a deterministic property of the LASSO.

Lemma 7. *If*

$$n\lambda \geq 2 \max_{k=1,\dots,m} \left| \sum_{i=1}^n (h_k(X_i) - P_n(h_k))(\epsilon(X_i) - P_n(\epsilon)) \right|, \quad (\text{D.1})$$

then, writing $\hat{u} = \hat{\beta}_n^{\text{lasso}}(f) - \beta^*(f)$, we have

$$\|H_c \hat{u}\|_2^2 \leq 3n\lambda\sqrt{\ell^*} \|\hat{u}_{S^*}\|_2. \quad (\text{D.2})$$

Proof. This is just a reformulation of the reasoning on p. 298 in [19]. The vector \hat{v} at the right-hand side of their Eq. (11.23) is in fact \hat{v}_S . \square

We work on the same event E as in Step 2 of the proof of Theorem 3.2, i.e., (C.9), (C.10), (C.11), (C.12) and C.13 are all satisfied. This event has probability $1 - \delta$. On this event, the LASSO solution $\hat{\beta}_n^{\text{lasso}}(f)$ is unique and $S(\hat{\beta}_n^{\text{lasso}}(f)) = S^*$, as shown in the proof of Theorem 3.2.

We show that on the event E and given the assumed bounds on n and λ , the inequality (D.1) is satisfied. As in Step 3.2 of the proof of Theorem 3.3, we have, on the event E , with $C = 8$ and with $\delta/4$ there replaced by $\delta/5$ here, the bound

$$\max_{k=1,\dots,m} \left| \sum_{i=1}^n (h_k(X_i) - P_n(h_k))(\epsilon(X_i) - P_n(\epsilon)) \right| \leq 4\sqrt{\log(10m/\delta)}\tau U\sqrt{n} \left[1 + \frac{1}{4}\sqrt{\frac{\log(10/\delta)}{n}} \right].$$

Since $\ell^* \geq 1$ and $U^2 \geq (1/\ell^*) \sum_{k \in S^*} P(h_k^2) \geq \gamma^*$, the assumed lower bounds on n imply that $n \geq 70 \log(10/\delta)$, so that the factor in square brackets is bounded by $1 + 1/(4 \cdot 8) = 33/32$. But as $n\lambda/2 \geq (17/2)\sqrt{\log(10m/\delta)}\tau U\sqrt{n}$ by assumption, the inequality (D.1) is clearly satisfied.

Set $\hat{u} = \hat{\beta}_n^{\text{lasso}}(f) - \beta^*(f)$. On the event E , the inequality (D.1) holds, and by Lemma 7, then also (D.2). Because $\hat{u}_k = 0$ whenever $k \notin S^*$, we have, as in Step 1 of the proof of Theorem 3.1, with $\hat{\beta}_n^{\text{ols}}(f)$ replaced by $\hat{\beta}_n^{\text{lasso}}(f)$, the inequality

$$\begin{aligned} n |\hat{\alpha}_n^{\text{lasso}}(f) - P(f)| &\leq \left| \sum_{i=1}^n \epsilon(X_i) \right| + \|\hat{u}_{S^*}\|_2 \sqrt{\ell^*} \max_{k \in S^*} \left| \sum_{i=1}^n h_k(X_i) \right| \\ &\leq \sqrt{2n\tau^2 \log(10/\delta)} + \|\hat{u}_{S^*}\|_2 \sqrt{2n\ell^*U^2 \log(10m/\delta)}. \end{aligned} \quad (\text{D.3})$$

Using (C.14) with $\alpha = \sqrt{2}/3$ [as in Step 6(ii) of the proof of Theorem 3.2] and then (D.2), we find

$$\begin{aligned} \|\hat{u}_{S^*}\|_2^2 &\leq \frac{1}{\alpha\gamma^*n} \|H_{c,S^*}\hat{u}_{S^*}\|_2^2 \\ &= \frac{3}{\sqrt{2}\gamma^*n} \|H_c \hat{u}\|_2^2 \leq \frac{3}{\sqrt{2}\gamma^*} \cdot 3\lambda\sqrt{\ell^*} \|\hat{u}_{S^*}\|_2. \end{aligned}$$

It follows that

$$\|\hat{u}_{S^*}\|_2 \leq \frac{9\lambda}{\sqrt{2}\gamma^*} \sqrt{\ell^*}.$$

Injecting this bound into (D.3), we obtain

$$|\hat{\alpha}_n^{\text{lasso}}(f) - P(f)| \leq \sqrt{2 \log(10/\delta)} \frac{\tau}{\sqrt{n}} + 9\lambda\ell^* \sqrt{\log(10m/\delta)} \frac{U/\gamma^*}{\sqrt{n}}.$$

Setting λ equal to its minimal value, the lower bound simplifies as stated. \square

E Additional graphs

We provide two additional illustrations. Figures E.1 and E.2 correspond to the same experiments as the ones of Figure 4.1 and 4.2 in the paper, respectively, except that the underlying dimension is now $d = 10$. The remarks given in Section 4, paragraph “Results”, extend to the present situation.

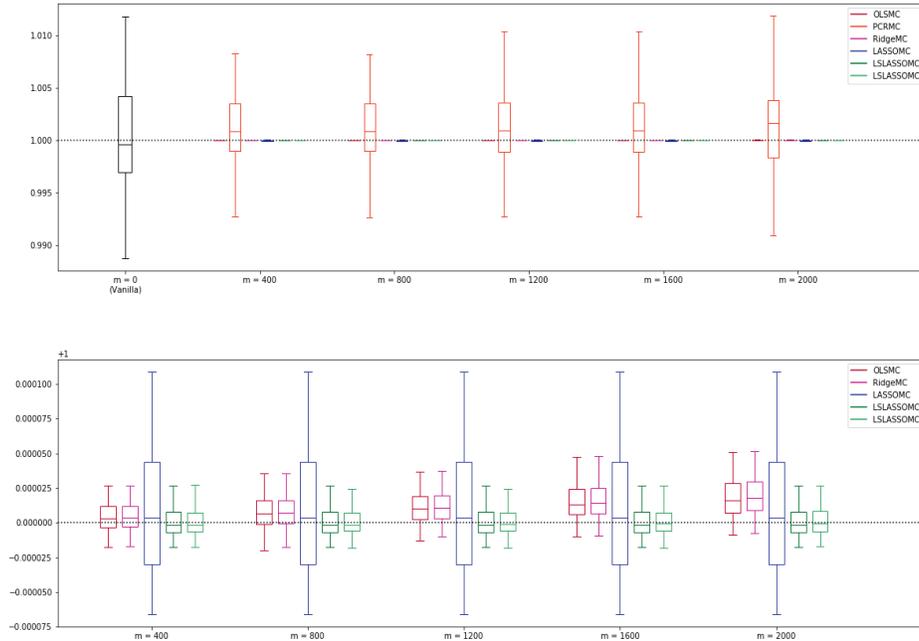


Figure E.1: Boxplots (based on 50 replications) of the values returned by each of the methods (top) and zooming on the best ones (bottom) for f_1 . The dimension is $d = 10$, the sample size is $n = 5000$ and m (horizontal axis) varies from 400 to 2000.

References

- [1] Alexandre Belloni, Victor Chernozhukov, et al. Least squares after model selection in high-dimensional sparse models. *Bernoulli*, 19(2):521–547, 2013.
- [2] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration Inequalities*. Oxford University Press, 2013.
- [3] LLdiko E Frank and Jerome H Friedman. A statistical view of some chemometrics regression tools. *Technometrics*, 35(2):109–135, 1993.
- [4] Paul Glasserman. *Monte Carlo Methods in Financial Engineering*, volume 53. Springer Science & Business Media, 2013.
- [5] Peter W Glynn and Roberto Szechtman. Some new perspectives on the method of control variates. In *Monte Carlo and Quasi-Monte Carlo Methods 2000*, pages 27–49. Springer, 2002.
- [6] Robert M Gower, Nicolas Le Roux, and Francis Bach. Tracking the gradients using the Hessian: A new look at variance reducing stochastic methods. *arXiv preprint arXiv:1710.07462*, 2017.

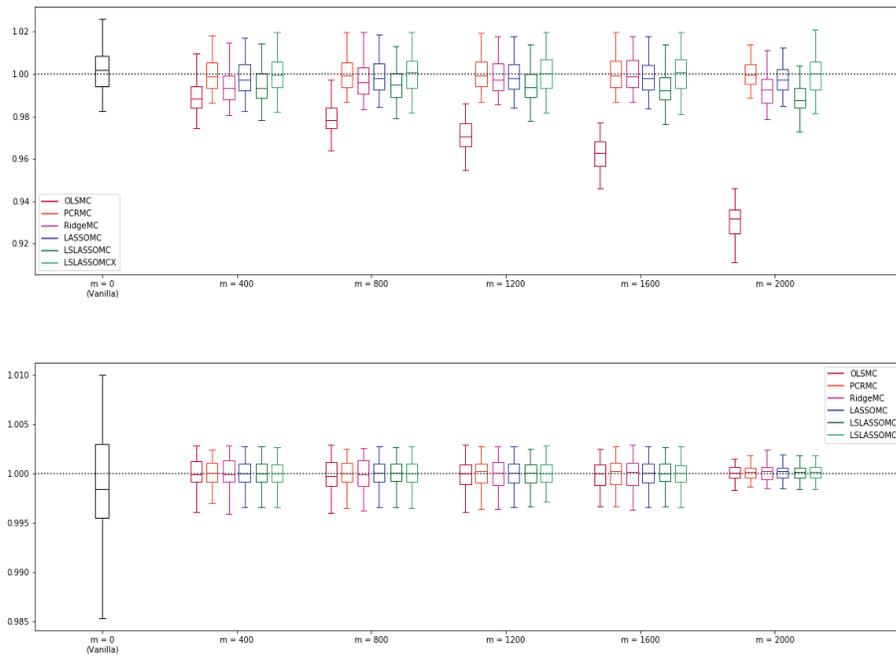


Figure E.2: Boxplots (based on 50 replications) of the values returned by each of the methods for f_2 (top) and f_3 (bottom). The dimension is $d = 10$, the sample size is $n = 10\,000$ and m (horizontal axis) varies from 400 to 2000.

- [7] Roger A Horn and Charles R Johnson. *Matrix Analysis*. Cambridge University Press, 2012.
- [8] Tang Jie and Pieter Abbeel. On a connection between importance sampling and the likelihood ratio policy gradient. In *Advances in Neural Information Processing Systems*, pages 1000–1008, 2010.
- [9] Hao Liu, Yihao Feng, Yi Mao, Dengyong Zhou, Jian Peng, and Qiang Liu. Action-depedent control variates for policy optimization via Stein’s identity. *arXiv preprint arXiv:1710.11198*, 2017.
- [10] Chris J Oates, Mark Girolami, and Nicolas Chopin. Control functionals for Monte Carlo integration. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):695–718, 2017.
- [11] Art Owen and Yi Zhou. Safe and effective importance sampling. *Journal of the American Statistical Association*, 95(449):135–143, 2000.
- [12] Art B. Owen. *Monte Carlo Theory, Methods and Examples*. 2013.
- [13] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [14] François Portier and Johan Segers. Monte Carlo integration with a growing number of control variates. *arXiv preprint arXiv:1801.01797*, 2018.
- [15] Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial Intelligence and Statistics*, pages 814–822, 2014.
- [16] Walter Rudin. *Real and Complex Analysis*. Tata McGraw-hill education, 2006.

- [17] Leah F South, Chris J Oates, Antonietta Mira, and Christopher Drovandi. Regularised zero-variance control variates for high-dimensional variance reduction. *arXiv preprint arXiv:1811.05073*, 2018.
- [18] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58(1):267–288, 1996.
- [19] Robert Tibshirani, Martin Wainwright, and Trevor Hastie. *Statistical Learning with Sparsity: The Lasso and Generalizations*. Chapman and Hall/CRC, 2015.
- [20] Chong Wang, Xi Chen, Alexander J Smola, and Eric P Xing. Variance reduction for stochastic gradient optimization. In *Advances in Neural Information Processing Systems*, pages 181–189, 2013.
- [21] Pavel Yaskov. Lower bounds on the smallest eigenvalue of a sample covariance matrix. *Electronic Communications in Probability*, 19(83):1–10, 2014.