# Pre-integration via active subspaces

Sifan Liu and Art B. Owen

Department of Statistics, Stanford University

February 2022

## Abstract

Pre-integration is an extension of conditional Monte Carlo to quasi-Monte Carlo and randomized quasi-Monte Carlo. It can reduce but not increase the variance in Monte Carlo. For quasi-Monte Carlo it can bring about improved regularity of the integrand with potentially greatly improved accuracy. Pre-integration is ordinarily done by integrating out one of $d$ input variables to a function. In the common case of a Gaussian integral one can also pre-integrate over any linear combination of variables. We propose to do that and we choose the first eigenvector in an active subspace decomposition to be the pre-integrated linear combination. We find in numerical examples that this active subspace pre-integration strategy is competitive with pre-integrating the first variable in the principal components construction on the Asian option where principal components are known to be very effective. It outperforms other pre-integration methods on some basket options where there is no well established default. We show theoretically that, just as in Monte Carlo, pre-integration can reduce but not increase the variance when one uses scrambled net integration. We show that the lead eigenvector in an active subspace decomposition is closely related to the vector that maximizes a less computationally tractable criterion using a Sobol' index to find the most important linear combination of Gaussian variables. They optimize similar expectations involving the gradient. We show that the Sobol' index criterion for the leading eigenvector is invariant to the way that one chooses the remaining $d-1$ eigenvectors with which to sample the Gaussian vector.

## 1 Introduction

Pre-integration [12] is a strategy for high dimensional numerical integration in which one variable is integrated out in closed form (or by a very accurate quadrature rule) while the others are handled by quasi-Monte Carlo (QMC) sampling. This strategy has a long history in Monte Carlo (MC) sampling [13, 43], where it is known as conditional Monte Carlo. In the Markov chain Monte Carlo (MCMC) literature, such conditioning is called Rao-Blackwellization, although it does not generally bring the optimal variance reduction that results from the Rao-Blackwell theorem. See [38] for a survey.

The advantage of pre-integration in QMC goes beyond the variance reduction that arises in MC. After pre-integration, a $d$-dimensional integration problem with a discontinuity or a kink (discontinuity in the gradient) can be converted into a much smoother $d-1$-dimensional problem [12]. QMC exploits regularity of the integrand and then smoothness brings a benefit on top of the

1

$L^2$ norm reduction that comes from conditioning. Gilbert, Kuo and Sloan [9] point out that the resulting smoothness depends critically on a monotonicity property of the integrand with respect to the variable being integrated out. Hoyt and Owen [17] give conditions where pre-integration reduces the mean dimension (that we will define below) of the integrand. It can reduce mean dimension from proportional to $\sqrt{d}$ to $O(1)$ as $d \to \infty$ in a sequence of ridge functions with a discontinuity that the pre-integration smooths out. He [14] studied the error rate of pre-integration for scrambled nets applied to functions of the form $f(\boldsymbol{x}) = h(\boldsymbol{x})\mathbf{1}\{\phi(\boldsymbol{x}) \geqslant 0\}$ for a Gaussian variable $\boldsymbol{x}$. That work assumes that $h$ and $\phi$ are smooth functions and $\phi$ is monotone in $x_j$. Then pre-integration has a smoothing effect that when combined with some boundary growth conditions yields an error rate of $O(n^{-1+\varepsilon})$.

Many of the use cases for pre-integration involve integration with respect to the multivariate Gaussian distribution, especially for problems arising in finance. In the Gaussian context, we have more choices for the variable to pre-integrate over. In addition to pre-integration over any one of the $d$ coordinates, pre-integration over any linear combination of the variables remains integration with respect to a univariate Gaussian variable. Our proposal is to pre-integrate over a linear combination of variables chosen to optimize a measure of variable importance derived from active subspaces [3].

When sampling from a multivariate Gaussian distribution by QMC, even without pre-integration, one must choose a square root of the covariance matrix by which to multiply some sampled scalar Gaussian variables. There are numerous choices for that square root. One can sample via the principal component matrix decomposition as [1] and many others do. For integrands defined with respect to Brownian motions, one can use the Brownian bridge construction studied by [23]. These options have some potential disadvantages. It is always possible that the integrand is little affected by the first principal component. In a pessimistic scenario, the integrand could be proportional to a principal component that is orthogonal to the first one. This is a well known pitfall in principal components regression [20]. In a related phenomenon, [34] exhibits an integrand where QMC via the standard construction is more effective than via the Brownian bridge.

Not only might a principal component direction perform poorly, the first principal component is not necessarily well defined. Although the problem may be initially defined in terms of a specific Gaussian distribution, by a change of variable we can rewrite our integral as an expectation with respect to another Gaussian distribution with a different covariance matrix that has a different first principal component. Or, if the problem is posed with a covariance equal to the $d$-dimensional identity matrix, then every unit vector linear combination of variables is a first principal component.

Some proposed methods take account of the specific integrand while formulating a sampling strategy. These include stratifying in a direction chosen from exponential tilting [11], exploiting a linearization of the integrand at $d + 1$ special points starting with the center of the domain [18], and a gradient principal component analysis (GPCA) algorithm [45] that we describe in more detail below.

The problem we consider is to compute an approximation to $\mu = \mathbb{E}(f(\boldsymbol{x}))$ where $\boldsymbol{x} \in \mathbb{R}^d$ has the spherical Gaussian distribution denoted by $\mathcal{N}(0, I)$ and $f$ has a gradient almost everywhere, that is square integrable. Let $C = \mathbb{E}(\nabla f(\boldsymbol{x}) \nabla f(\boldsymbol{x})^{\mathsf{T}}) \in \mathbb{R}^{d \times d}$. The $r$ dimensional active subspace [3] is the space spanned by the $r$ leading eigenvectors of $C$. For other uses of the matrix $C$ in numerical computation, see the references in [4]. For $r = 1$, let $\theta$ be the leading eigenvector of $C$ normalized to be a unit vector. Our use is to pre-integrate $f$ over $\theta^{\mathsf{T}}\boldsymbol{x} \sim \mathcal{N}(0, 1)$.

The eigendecomposition of $C$ is an uncentered principal components decomposition of the gradients, also known as the GPCA. The GPCA method [45] also uses the eigendecomposition of $C$ to define a matrix square root for a QMC sampling strategy to reduce effective dimension, but

it involves no pre-integration. The algorithm in [44] pre-integrates the first variable $x_1$ out of $f$. Then it applies GPCA to the remaining $d-1$ variables in order to find a suitable $(d-1) \times (d-1)$ matrix square root for the remaining Gaussian variables. They pre-integrate over a coordinate variable while we always pre-integrate over the leading eigenvector which is not generally one of the $d$ coordinates. All of the algorithms that involve $C$ take a sample in order to estimate it.

This paper is organized as follows. Section 2 provides some background on RQMC and pre-integration. Section 3 shows that pre-integration never increases the variance of scrambled net integration, extending the well known property of conditional integration to RQMC. This holds even if the points being scrambled are not a digital net. There is presently very little guidance in the literature about which variable to pre-integrate over, beyond the monotonicity condition recently studied in [9] and a remark for ridge functions in [17]. An RQMC variance formula from Dick and Pillichshammer [7] overlaps greatly with a certain Sobol' index and so this section advocates using the variable that maximizes that Sobol' index of variable importance. That index has an effective practical estimator due to Jansen [19]. In Section 4, we describe using the unit vector $\theta$ which maximizes the Sobol' index for $\theta^\mathsf{T} \boldsymbol{x}$. This strategy is to find a matrix square root for which the first column maximizes the criterion from Section 3. We show that this Sobol' index $\underline{\tau}_\theta^2$ is well defined in that it does not depend on how we parameterize the space orthogonal to $\theta$. This index is upper bounded by $\theta^\mathsf{T} \mathbb{E}(\nabla f(\boldsymbol{x}) \nabla f(\boldsymbol{x})^\mathsf{T})\theta$. When $f$ is differentiable, the Sobol' index can be expressed as a weighted expectation of $\theta^\mathsf{T} \nabla f(\boldsymbol{x}) \nabla f(\boldsymbol{x})^\mathsf{T} \theta$. As a result, choosing a projection by active subspaces amounts to optimizing a computationally convenient proxy measure for a Sobol' index of variable importance. We apply active subspace pre-integration to option pricing examples in Section 5. These include an Asian call option and some of its partial derivatives (called the Greeks) as well as a basket option. The following summary is based on the figures in that section. Using the standard construction, the active subspace pre-integration is a clear winner for five of the six Asian option integration problems and is essentially tied with the method from [44] for the one called Gamma. The principal components construction of [1] is more commonly used for this problem and has been extensively studied. In that construction active subspace pre-integration is more accurate than the other methods for Gamma, is worse than some others for Rho and is essentially the same as the best methods in the other four problems. Every one of the pre-integration methods was more accurate than RQMC without pre-integration, consistent with Theorem 3.2. We also looked at a basket option for which there is not a strong default construction as well accepted as the PCA construction is for the Asian option. There in six cases, two basket options and three baseline sampling constructions, active subspace pre-integration was always the most accurate. That section makes some brief comments about the computational cost. In our simulations, the pre-integration methods cost significantly more than the other methods, but the cost factor is not enough to overwhelm their improved accuracy, except that pre-integrating the first variable in the standard construction of Brownian motion usually brought little to no advantage. Section 6 has our conclusions.

## 2   Background

In this section, we introduce some background of RQMC and pre-integration and active subspaces. First we introduce some notations. Additional notations are introduced as needed.

For a positive integer $d$, we let $1{:}d = \{1, 2, \ldots, d\}$. For a subset $u \subseteq 1{:}d$, let $-u = 1{:}d \setminus u$. For an integer $j \in 1{:}d$, we use $j$ to represent $\{j\}$ when the context is clear and $-j = 1{:}d \setminus \{j\}$. Let $|u|$ denote the cardinality of $u$. For $\boldsymbol{x} \in \mathbb{R}^d$, we let $\boldsymbol{x}_u$ be the $|u|$ dimensional vector containing

only the $x_j$ with $j \in u$. For $\boldsymbol{x}, \boldsymbol{z} \in \mathbb{R}^d$ and $u \subseteq 1{:}d$, we let $\boldsymbol{x}_u{:}\boldsymbol{z}_{-u}$ be the $d$ dimensional vector, whose $j$-th entry is $x_j$ if $j \in u$, and $z_j$ if $j \notin u$. We use $\mathbb{N} = \{1, 2, \dots\}$ for the natural numbers and $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$. We denote the density and the cumulative distribution function (CDF) of standard Gaussian distribution $\mathcal{N}(0, 1)$ as $\varphi$ and $\Phi$, respectively. We let $\Phi^{-1}$ denote the inverse CDF of $\mathcal{N}(0, 1)$. We also use $\varphi$ to denote the density of the $d$-dimensional standard Gaussian distribution $\mathcal{N}(0, I_d)$, $\varphi(\boldsymbol{x}) = (2\pi)^{-d/2} \exp(-\|\boldsymbol{x}\|^2/2)$. We use $\mathcal{N}(0, I)$ when the dimension of the random variable is clear from context. For a matrix $\Theta \in \mathbb{R}^{d \times d}$ we often need to select a subset of columns. We do that via $\Theta[u]$, such as $\Theta[1{:}r]$ or $\Theta[2{:}d]$.

## 2.1 QMC and RQMC

For background on QMC see the monograph [7] and the survey article [6]. For a survey of RQMC see [21]. We will emphasize scrambled net integration with a recent description in [32]. Here we give a brief account.

QMC provides a way to estimate $\mu = \int_{[0,1]^d} f(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x}$ with greater accuracy than can be done by MC while sharing with MC the ability to handle larger dimensions $d$ than can be well handled by classical quadrature methods such as those in [5]. The QMC estimate, like the MC one takes the form $\hat{\mu}_n = (1/n) \sum_{i=0}^{n-1} f(\boldsymbol{x}_i)$, except that instead of $\boldsymbol{x}_i \overset{\text{iid}}{\sim} \mathbb{U}[0,1]^d$ the sample points are chosen strategically to get a small value for

$$D_n^* = D_n^*(\boldsymbol{x}_0, \dots, \boldsymbol{x}_{n-1}) = \sup_{\boldsymbol{a} \in [0,1]^d} \left| \frac{1}{n} \sum_{i=0}^{n-1} 1\{\boldsymbol{x}_i \in [\boldsymbol{0}, \boldsymbol{a})\} - \prod_{j=1}^{d} a_j \right|$$

which is known as the star discrepancy. The Koksma-Hlawka inequality (see [15]) is

$$|\hat{\mu}_n - \mu| \leqslant D_n^* \times \|f\|_{\text{HK}} \tag{1}$$

where $\|\cdot\|_{\text{HK}}$ denotes total variation in the sense of Hardy and Krause. It is possible to construct points $\boldsymbol{x}_i$ with $D_n^* = O((\log n)^{d-1}/n)$ or to choose an infinite sequence of them along which $D_n^* = O((\log n)^d/n)$. Both of these are often written as $O(n^{-1+\epsilon})$ for any $\epsilon > 0$ and this rate translates directly to $|\hat{\mu}_n - \mu|$ when $f \in \text{BVHK}[0,1]^d$, the set of functions of bounded variation in the sense of Hardy and Krause. While the logarithmic powers are not negligible they correspond to worst case integrands that are not representative of integrands evaluated in practice, and as we see below, some RQMC methods provide a measure of control against them.

The QMC methods we study are digital nets and sequences. To define them, for an integer base $b \geqslant$ let

$$E(\boldsymbol{k}, \boldsymbol{c}) = \prod_{j=1}^{d} \left[ \frac{c_j}{b^{k_j}}, \frac{c_j + 1}{b^{k_j}} \right) \tag{2}$$

for $\boldsymbol{k} = (k_1, \dots, k_d)$ and $\boldsymbol{c} = (c_1, \dots, c_d)$ where $k_j \in \mathbb{N}_0$ and $0 \leqslant a_j < b^{k_j}$. The sets $E(\boldsymbol{k}, \boldsymbol{c})$ are called elementary intervals in base $b$. They have volume $b^{-|\boldsymbol{k}|}$ where $|\boldsymbol{k}| = \sum_{j=1}^{d} k_j$. For integers $m \geqslant t \geqslant 0$, the points $\boldsymbol{x}_0, \dots, \boldsymbol{x}_{n-1}$ with $n = b^m$ are a $(t, m, d)$-net in base $b$ if every elementary interval $E(\boldsymbol{k}, \boldsymbol{c})$ with $|\boldsymbol{k}| \leqslant m - t$ contains $b^{m-|\boldsymbol{k}|}$ of those points, which is exactly $n$ times the volume of $E(\boldsymbol{k}, \boldsymbol{c})$.

For each $\boldsymbol{k}$ the sets $E(\boldsymbol{k}, \boldsymbol{c})$ partition $[0,1)^d$ into $b^{|\boldsymbol{k}|}$ congruent half open sets. If $|\boldsymbol{k}| \leqslant m - t$, then the $n$ points $\boldsymbol{x}_i$ are perfectly stratified over that partition. The power of digital nets is that the points $\boldsymbol{x}_i$ satisfy $\binom{m-t+d-1}{d-1}$ such stratifications simultaneously. They attain $D_n^* = O((\log n)^{d-1}/n)$ after approximating each $[\boldsymbol{0}, \boldsymbol{a})$ by sets $E(\boldsymbol{k}, \boldsymbol{c})$ [24]. Given $b$ and $m$ and $d$, a smaller $t$ is the better. It is not always possible to get $t = 0$.

For an integer $t \geqslant 0$, the infinite sequence $(\boldsymbol{x}_i)_{i \geqslant 0}$ is a $(t, s)$-sequence in base $b$ if for all integers $m \geqslant t$ and $r \geqslant 0$, the points $\boldsymbol{x}_{rb^m}, \ldots, \boldsymbol{x}_{rb^m + b^m - 1}$ are a $(t, m, d)$-net in base $b$. This means that the first $b^m$ points are a $(t, m, d)$-net as are the second $b^m$ points and if we take the first $b$ such nets together we get a $(t, m+1, d)$-net. Similarly the first $b$ such $(t, m+1, d)$-nets form a $(t, m+2, d)$-net and so on ad infinitum.

The nets and sequences in common use are called digital nets and sequences owing to a digital strategy used in their construction, where $x_{ij}$ is expanded in base $b$ and there are rules for constructing those base $b$ digits from the base $b$ representation of $i$. See [24]. The most commonly used nets and sequences are those of Sobol' [39]. Sobol' sequences are $(t, d)$-sequences in base 2 and sometimes using base 2 brings computational advantages for computers that work in base 2. The value of $t$ is nondecreasing in $d$. The first $2^m$ points of a $(t, d)$-sequence can be a $(t', m, d)$-net for some $t' < t$.

While the Koksma-Hlawka inequality (1) shows that QMC is asymptotically better than MC for $f \in \mathrm{BVHK}[0,1]^d$ it is not usable for error estimation, and furthermore, many integrands of interest such as unbounded ones have $V_{\mathrm{HK}}(f) = \infty$. RQMC methods can address both problems. In RQMC the points $\boldsymbol{x}_i \sim \mathbb{U}[0,1]^d$ individually while collectively they have a small $D_n^*$. The uniformity property makes

$$\hat{\mu}_n = \frac{1}{n} \sum_{i=0}^{n-1} f(\boldsymbol{x}_i) \tag{3}$$

an unbiased estimate of $\mu$ when $f \in L^1[0,1]^d$. If $f \in L^2[0,1]^d$ then $\mathrm{Var}(\hat{\mu}_n) < \infty$ and it can be estimated by using independent repeated RQMC evaluations.

Scrambled $(t, m, d)$-nets have $\boldsymbol{x}_i \sim \mathbb{U}[0,1]^d$ individually and the collective condition is that they form a $(t, m, d)$-net with probability one. For an infinite sequence of $(t, m, d)$-nets one can use scrambled $(t, d)$-sequences. See [26] for both of these. The estimate $\hat{\mu}_n$ taken over a scrambled $(t, d)$-sequence satisfies a strong law of large numbers if $f \in L^2[0,1]^{1+\epsilon}$ [32]. If $f \in L^2[0,1]^d$, then $\mathrm{Var}(\hat{\mu}_n) = o(1/n)$ giving the method asymptotically unbounded efficiency versus MC which has variance $\sigma^2/n$ for $\sigma^2 = \mathrm{Var}(f(\boldsymbol{x}))$. For smooth enough $f$, $\mathrm{Var}(\hat{\mu}_n) = O(n^{-3}(\log n)^{d-1})$ [28, 30] with the sharpest sufficient condition in [46]. Also there exists $\Gamma < \infty$ with $\mathrm{Var}(\hat{\mu}_n) \leqslant \Gamma \sigma^2/n$ for all $f \in L^2[0,1]^d$ [29]. This bound involves no powers of $\log(n)$.

## 2.2 Pre-integration with respect to $\mathbb{U}[0,1]^d$ and $\mathcal{N}(0, I)$

For $j \in 1{:}d$, $\int_{[0,1]^d} f(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = \int_{[0,1]^{d-1}} \int_0^1 f(\boldsymbol{x}) \, \mathrm{d}x_j \, \mathrm{d}\boldsymbol{x}_{-j}$ which we can also write as $\mathbb{E}(f(\boldsymbol{x})) = \mathbb{E}(\mathbb{E}(f(\boldsymbol{x}) \,|\, \boldsymbol{x}_{-j}))$ for $\boldsymbol{x} \sim \mathbb{U}[0,1]^d$. For $\boldsymbol{x} \in [0,1]^d$, define

$$g(\boldsymbol{x}) = g_j(\boldsymbol{x}) = \int_0^1 f(\boldsymbol{x}) \, \mathrm{d}x_j.$$

It simplifies the presentation of some of our results, especially Theorem 3.2 on variance reduction, to keep $g$ defined as above on $[0,1]^d$ even though it does not depend at all on $x_j$ and could be

written as a function of $\boldsymbol{x}_{-j} \in [0,1]^{d-1}$ instead. In pre-integration

$$\hat{\mu}_n = \hat{\mu}_{n,j} = \frac{1}{n} \sum_{i=0}^{n-1} g_j(\boldsymbol{x}_i)$$

which as we noted in the introduction is conditional MC except that we now use RQMC inputs.

Pre-integration can bring some advantages for RQMC. The plain MC variance of $g(\boldsymbol{x})$ is no larger than that of $f(\boldsymbol{x})$ and is generally smaller unless $f$ does not depend at all on $x_j$. Thus the bound $\Gamma \sigma^2/n$ is reduced. Next, the pre-integrated integrand $g$ can be much smoother than $f$ and (R)QMC improves on MC by exploiting smoothness. For example [12, 44] observe that for some option pricing integrands, pre-integrating certain variable can remove the discontinuities in the integrand or its gradient.

The integrands we consider here are defined with respect to a Gaussian random variable. We are interested in $\mu = \mathbb{E}(f(\boldsymbol{z}))$ for $\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0}, \Sigma)$ and a nonsingular covariance $\Sigma \in \mathbb{R}^{d \times d}$. Letting $R_0 \in \mathbb{R}^{d \times d}$ with $R_0 R_0^{\mathsf{T}} = \Sigma$ we can write $\mu = \mathbb{E}(f(R_0 \boldsymbol{z}))$ for $\boldsymbol{z} \sim \mathcal{N}(0, I)$. For an orthogonal matrix $Q \in \mathbb{R}^{d \times d}$ we also have $Q\boldsymbol{z} \sim \mathcal{N}(0, I)$. Then taking $\boldsymbol{z} = \Phi^{-1}(\boldsymbol{x})$ componentwise leads us to the estimate

$$\hat{\mu} = \frac{1}{n} \sum_{i=0}^{n-1} f(R\Phi^{-1}(\boldsymbol{x}_i)) \quad \text{with } R = R_0 Q$$

for RQMC points $\boldsymbol{x}_i$. The choice of $Q$ or equivalently $R$ does not affect the MC variance of $\hat{\mu}$ but it can change the RQMC variance. We will consider some examples later. The mapping $\Phi^{-1}$ from $\mathbb{U}[0,1]^d$ to $\mathcal{N}(0, I)$ can be replaced by another one such as the Box-Muller transformation. The choice of transformation does not affect the MC variance but does affect the RQMC variance. Most researchers use $\Phi^{-1}$ but [25] advocates for Box-Muller.

When we are using pre-integration for a problem defined with respect to a $\mathcal{N}(0, \Sigma)$ random variable we must choose $R$ and then the coordinate $j$ over which to pre-integrate. Our approach is to choose $R$ to make coordinate $j = 1$ as important as we can, using active subspaces.

## 2.3 The ANOVA decomposition

For $f \in L^2[0,1]^d$ we can define an analysis of variance (ANOVA) decomposition from [16, 40, 8]. For details see [31, Appendix A.6]. This decomposition writes

$$f(\boldsymbol{x}) = \sum_{u \subseteq 1:d} f_u(\boldsymbol{x})$$

where $f_u$ depends on $\boldsymbol{x}$ only through $x_j$ with $j \in u$ and also $\int_0^1 f_u(\boldsymbol{x}) \, \mathrm{d}x_j = 0$ whenever $j \in u$. The decomposition is orthogonal in that $\int_{[0,1]^d} f_u(\boldsymbol{x}) f_v(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x} = 0$ if $u \neq v$. The term $f_\varnothing$ is the constant function everywhere equal to $\mu = \int_{[0,1]^d} f(\boldsymbol{x}) \, \mathrm{d}\boldsymbol{x}$. To each effect $f_u$ there corresponds a variance component

$$\sigma_u^2 = \mathrm{Var}(f_u(\boldsymbol{x})) = \begin{cases} \int_{[0,1]^d} f_u(\boldsymbol{x})^2 \, \mathrm{d}\boldsymbol{x}, & |u| > 0 \\ 0, & \text{else.} \end{cases}$$

For $|u| \geqslant 2$, the effect $f_u$ is called a $|u|$-fold interaction. The variance components sum to $\sigma^2 = \mathrm{Var}(f(\boldsymbol{x}))$. We will use the ANOVA decomposition below when describing how to choose a pre-integration variable.

Sobol' indices [41] are derived from the ANOVA decomposition. For $u \subseteq 1{:}d$ these are

$$\underline{\tau}_u^2 = \sum_{v \subseteq u} \sigma_v^2 \quad \text{and} \quad \overline{\tau}_u^2 = \sum_{v \subseteq 1:d} \mathbf{1}_{\{u \cap v \neq \varnothing\}} \sigma_v^2.$$

They provide two ways to judge the importance of the set of variables $x_j$ for $j \in u$. They're usually normalized by $\sigma^2$ to get an interpretation as a proportion of variance explained.

The mean dimension of $f$ is $\nu(f) = \sum_{u \subseteq 1:d} |u| \sigma_u^2 / \sigma^2$. It satisfies $\nu(f) = \sum_{j=1}^d \overline{\tau}_j^2 / \sigma^2$. A ridge function takes the form $f(\boldsymbol{x}) = h(\Theta^\mathsf{T} \boldsymbol{x})$ for $\Theta \in \mathbb{R}^{d \times r}$ with $\Theta^\mathsf{T} \Theta = I$. For $\boldsymbol{x} \sim \mathcal{N}(0, I)$, the variance of $f$ does not depend on $d$ and the mean dimension is $O(1)$ as $d \to \infty$ [17] if $h$ is Lipschitz. If $h$ has a step discontinuity then it is possible to have $\nu(f) = \Omega(\sqrt{d})$ reduced to $O(1)$ by pre-integration over a component variable $x_j$ with $\theta_j$ bounded away from zero as $d \to \infty$ [17].

# 3 Pre-integration and scrambled net variance

Conditional MC can reduce but not increase the variance of plain MC integration. Here we show that the same thing holds for scrambled nets using the nested uniform scrambling of [26]. The affine linear scrambling of [22] has the same variance and hence the same result. We assume that $f \in L^2[0,1)^d$. The half-open interval is just a notational convenience. For any $f \in L^2[0,1]^d$ we could set $f(\boldsymbol{x}) = 0$ for any $\boldsymbol{x} \in D = [0,1]^d \setminus [0,1)^d$ and get an equivalent function with the same integral and, almost surely, the same RQMC estimate because all $\boldsymbol{x}_i \sim \mathbb{U}[0,1]^d$ avoid $D$ with probability one.

We will pre-integrate over one of the $d$ components of $\boldsymbol{x} \in [0,1)^d$. It is also possible to pre-integrate over multiple components and reduce the RQMC variance each time, though the utility of that strategy is limited by the availability of suitable closed forms or effective quadratures.

## 3.1 Walsh function expansions

To get variance formulas for scrambled nets we follow Dick and Pillichshammer [7] who work with a Walsh function expansion of $L^2[0,1)^d$ for which they credit Pirsic [36]. Let $\omega_b = e^{2\pi i / b}$ with $i$ being the imaginary unit. For $k \in \mathbb{N}_0$ write $k = \sum_{j \geqslant 0} \kappa_j b^j$ for base $b$ digits $\kappa_j \in \{0, 1, \dots, b-1\}$. For $x \in [0,1)$ write $x = \sum_{k \geqslant 1} \xi_j b^{-j}$ for base $b$ digits $\xi_j \in \{0, 1, \dots, b-1\}$. Countably many $x$ have an expansion terminating in either infinitely many 0s or in infinitely many $b-1$s. For those we always choose the expansion terminating in 0s.

Using the above notation we can define the $k$-th $b$-adic Walsh function ${}_b\mathrm{wal}_k : [0,1) \to \mathbb{C}$ as

$$_b\mathrm{wal}_k(x) = \omega_b^{\sum_{j \geqslant 1} \xi_j \kappa_{j-1}}.$$

The summation in the exponent is finite because $k < \infty$. Note that ${}_b\mathrm{wal}_0(x) = 1$ for all $x \in [0,1)$. For $\boldsymbol{x} = (x_1, \dots, x_d) \in [0,1)^d$ and $\boldsymbol{k} = (k_1, \dots, k_d) \in \mathbb{N}_0^d$, the $d$-dimensional Walsh functions are defined as

$$_b\mathrm{wal}_{\boldsymbol{k}}(\boldsymbol{x}) = \prod_{j=1}^d {}_b\mathrm{wal}_{k_j}(x_j).$$

The Walsh series expansion of $f(x)$ is

$$f(\boldsymbol{x}) \sim \sum_{\boldsymbol{k} \in \mathbb{N}_0^d} \hat{f}(\boldsymbol{k})_b \mathrm{wal}_{\boldsymbol{k}}(\boldsymbol{x}), \text{ where } \hat{f}(\boldsymbol{k}) = \int_{[0,1)^d} f(\boldsymbol{x}) \overline{_b \mathrm{wal}_{\boldsymbol{k}}(\boldsymbol{x})} \, \mathrm{d}\boldsymbol{x}.$$

The $d$-dimensional $b$-adic Walsh function system is a complete orthonormal basis in $L_2([0,1)^d)$ [7, Theorem A.11] and the series expansion converges to $f$ in $L^2$.

While our integrand is real valued, it will also satisfy an expansion written in terms of complex numbers. For real valued $f$,

$$\mathrm{Var}(f) = \sum_{\boldsymbol{k} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} |\hat{f}(\boldsymbol{k})|^2.$$

The variance under scrambled nets is different. To study it we group the Walsh coefficients. For $\boldsymbol{\ell} \in \mathbb{N}_0^d$ let

$$C_{\boldsymbol{\ell}} = \left\{ \boldsymbol{k} \in \mathbb{N}_0^d \mid \lfloor b^{k_j - 1} \rfloor \leqslant k_j < b^{\ell_j}, 1 \leqslant j \leqslant d \right\}.$$

Then define

$$\beta_{\boldsymbol{\ell}}(\boldsymbol{x}) = \sum_{\boldsymbol{k} \in C_{\boldsymbol{\ell}}} \hat{f}(\boldsymbol{k})_b \mathrm{wal}_{\boldsymbol{k}}(\boldsymbol{x}).$$

The functions $\beta_{\boldsymbol{\ell}}$ are orthogonal in that $\int_{[0,1)^d} \beta_{\boldsymbol{\ell}}(\boldsymbol{x}) \overline{\beta_{\boldsymbol{\ell}'}(\boldsymbol{x})} \, \mathrm{d}\boldsymbol{x} = 0$ when $\boldsymbol{\ell}' \neq \boldsymbol{\ell}$. For $\boldsymbol{\ell} \neq \boldsymbol{0}$, $\beta_{\boldsymbol{\ell}}(\boldsymbol{x})$ has variance

$$\sigma_{\boldsymbol{\ell}}^2 = \int_{[0,1)^d} |\beta_{\boldsymbol{\ell}}(\boldsymbol{x})|^2 \, \mathrm{d}\boldsymbol{x} = \sum_{\boldsymbol{k} \in C_{\boldsymbol{\ell}}} |\hat{f}(\boldsymbol{k})|^2.$$

If $\boldsymbol{x}_i$ are a scrambled version of original points $\boldsymbol{a}_i \in [0,1)^d$ then under this scrambling

$$\mathrm{Var}(\hat{\mu}_n) = \frac{1}{n} \sum_{\boldsymbol{\ell} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} \Gamma_{\boldsymbol{\ell}} \sigma_{\boldsymbol{\ell}}^2 \tag{4}$$

for a collection of gain coefficients $\Gamma_{\boldsymbol{\ell}} \geqslant 0$ that depend on the $\boldsymbol{a}_i$. This expression can also be obtained through a base $b$ Haar wavelet decomposition [27]. Our $\Gamma_{\boldsymbol{\ell}}$ equals $nG_{\boldsymbol{\ell}}$ from [7]. The variance of $\hat{\mu}$ under IID MC sampling is $(1/n) \sum_{\boldsymbol{\ell} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} \sigma_{\boldsymbol{\ell}}^2$ so $\Gamma_{\boldsymbol{\ell}} < 1$ corresponds to integrating the term $\beta_{\boldsymbol{\ell}}(\boldsymbol{x})$ with less variance than MC does.

If scrambling of [26] or [22] is applied to $\boldsymbol{a}_i$ then

$$\Gamma_{\boldsymbol{\ell}} = \frac{1}{n} \sum_{i,i'=0}^{n-1} \prod_{j=1}^d \frac{b \mathbf{1} \left\{ \lfloor b^{\ell_j} a_{i,j} \rfloor = \lfloor b^{\ell_j} a_{i',j} \rfloor \right\} - \mathbf{1} \left\{ \lfloor b^{\ell_i - 1} a_{i,j} \rfloor = \lfloor b^{\ell_j - 1} a_{i',j} \rfloor \right\}}{b - 1}. \tag{5}$$

This holds for any $\boldsymbol{a}_i$ not just digital nets. When $\boldsymbol{a}_i$ are the first $b^m$ points of a $(t,d)$-sequence in base $b$ then $\Gamma = \sup_{\boldsymbol{\ell}} \Gamma_{\boldsymbol{\ell}} < \infty$ (uniformly in $m$) so that $\mathrm{Var}(\hat{\mu}) \leqslant \Gamma \sigma^2 / n$. Similarly for any $\boldsymbol{\ell} \in \mathbb{N}_0^d$ we have $\Gamma_{\boldsymbol{\ell}} \to 0$ as $n = b^m \to \infty$ in a $(t,d)$-sequence in base $b$ from which $\mathrm{Var}(\hat{\mu}_n) = o(1/n)$. For a $(t,m,s)$-net in base $b$, one can show that the gain coefficients $\Gamma_{\boldsymbol{\ell}} = 0$ for all $\boldsymbol{\ell}$ with $|\boldsymbol{\ell}| \leqslant m - t$.

## 3.2 Walsh decomposition after pre-integration

**Proposition 3.1.** *For $f \in L^2[0,1)^d$ and $j \in 1{:}d$, let $g$ be $f$ pre-integrated over $x_j$. Then for $\boldsymbol{k} \in \mathbb{N}_0^d$,*

$$\hat{g}(\boldsymbol{k}) = \begin{cases} \hat{f}(\boldsymbol{k}), & k_j = 0 \\ 0, & k_j > 0. \end{cases} \tag{6}$$

*Proof.* We write

$$\hat{g}(\boldsymbol{k}) = \int_{[0,1)^{d-1}} \int_0^1 g(\boldsymbol{x}) \overline{{}_b\mathrm{wal}_{\boldsymbol{k}}(\boldsymbol{x})} \, \mathrm{d}x_j \, \mathrm{d}\boldsymbol{x}_{-j} = \int_{[0,1)^{d-1}} g(\boldsymbol{x}) \int_0^1 \overline{{}_b\mathrm{wal}_{\boldsymbol{k}}(\boldsymbol{x})} \, \mathrm{d}x_j \, \mathrm{d}\boldsymbol{x}_{-j}$$

because $g(\boldsymbol{x})$ does not depend on $x_j$. If $k_j > 0$, then the inner integral vanishes establishing the second clause in (6). If $k_j = 0$ then ${}_b\mathrm{wal}_{k_j}(x_j) = 1$ for all $x_j$ and the inner integral equals $\prod_{\ell \neq j} {}_b\mathrm{wal}_{k_j}(x_j) = {}_b\mathrm{wal}_{\boldsymbol{k}}(\boldsymbol{x})$ establishing the second clause. $\qquad\square$

**Theorem 3.2.** *For $\boldsymbol{a}_0, \dots, \boldsymbol{a}_{n-1} \in [0,1)^d$ let $\boldsymbol{x}_0, \dots, \boldsymbol{x}_{n-1}$ be a scrambled version of them using the algorithm from [26] or [22]. Let $f \in L^2[0,1)^d$ and for $j \in 1{:}d$, let $g$ be $f$ pre-integrated over $x_j$. Then*

$$\mathrm{Var}\Big(\frac{1}{n}\sum_{i=0}^{n-1} g(\boldsymbol{x}_i)\Big) \leqslant \mathrm{Var}\Big(\frac{1}{n}\sum_{i=0}^{n-1} f(\boldsymbol{x}_i)\Big).$$

*Proof.* With either $f$ or $g$ we have the same gain coefficients $\Gamma_{\boldsymbol{\ell}}$ for $\boldsymbol{\ell} \in \mathbb{N}_0^d$. However

$$\sigma_{\boldsymbol{\ell}}^2(g) = \sum_{\boldsymbol{k} \in C_{\boldsymbol{\ell}}} |\hat{g}(\boldsymbol{k})|^2 = \sum_{\boldsymbol{k} \in C_{\boldsymbol{\ell}}, k_j=0} |\hat{f}(\boldsymbol{k})|^2 \leqslant \sum_{\boldsymbol{k} \in C_{\boldsymbol{\ell}}} |\hat{f}(\boldsymbol{k})|^2 = \sigma_{\boldsymbol{\ell}}^2(f).$$

The result now follows from (4). $\qquad\square$

Theorem 3.2 shows that pre-integration does not increase the variance under scrambling. This holds whether or not the underlying points are a digital net, though of course the main case of interest is for scrambling of digital nets and sequences.

Pre-integration has another benefit that is not captured by Theorem 3.2. By reducing the input dimension from $d$ to $d-1$ we might be able to find better points in $[0,1]^{d-1}$ than the $(t,m,d)$-net we would otherwise use in $[0,1]^d$. Those improved points might have some smaller gain coefficients or they might be a $(t', m, s-1)$ net in base $b$ with $t' < t$.

For scrambled net sampling, reducing the dimension reduces an upper bound on the variance. For any function $f \in L^2[0,1)^d$, the variance using a scrambled $(t,m,d)$-net in base $b$ is at most $b^{t+d}$ times the MC variance. Reducing the dimension reduces the bound to $b^{t+d-1}$ times the MC variance. For digital constructions in base 2, there are sharper bounds on these ratios, $2^{t+d-1}$ and $2^{t+d-2}$, respectively [33] and for some nets described there, even lower bounds apply.

As remarked above, pre-integration over a variable $x_j$ that $f(\boldsymbol{x})$ uses will reduce the variance under scrambled net sampling. This reduction does not require $f$ to be monotone in $x_j$, though such cases have the potential to bring a greater improvement. Pre-integration can either increase or decrease the mean dimension because

$$\nu(f) = \frac{\sum_{\boldsymbol{k} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} |\hat{f}(\boldsymbol{k})|^2 \times \|\boldsymbol{k}\|_0}{\sum_{\boldsymbol{k} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} |\hat{f}(\boldsymbol{k})|^2} \quad \text{and} \quad \nu(g_j) = \frac{\sum_{\boldsymbol{k} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} \boldsymbol{1}_{\{k_j>0\}} |\hat{f}(\boldsymbol{k})|^2 \times \|\boldsymbol{k}\|_0}{\sum_{\boldsymbol{k} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} \boldsymbol{1}_{\{k_j>0\}} |\hat{f}(\boldsymbol{k})|^2}$$

9

and pre-integration could possibly reduce the denominator by a greater proportion than it reduces the numerator.

## 3.3 Choice of $x_j$

In order to choose $x_j$ to pre-integrate over, we can look at the variance reduction we get. Pre-integrating over $x_j$ reduces the scrambled net variance by

$$\frac{1}{n} \sum_{\boldsymbol{\ell} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} \Gamma_{\boldsymbol{\ell}} \sigma_{\boldsymbol{\ell}}^2 \mathbf{1}_{\{\ell_j > 0\}} = \frac{1}{n} \sum_{\boldsymbol{k} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} \Gamma_{\boldsymbol{\ell}} |\hat{f}(\boldsymbol{k})|^2 \mathbf{1}_{\{k_j > 0\}}. \tag{7}$$

Evaluating this quantity for each $j \in 1{:}d$ might be more expensive than getting a good estimate of $\mu$. However we don't need to find the best $j$. Any $j$ where $f$ depends on $x_j$ will bring some improvement. Below we develop a principled and computationally convenient choice by choosing the $j$ which is most important as measured by a Sobol' index [41] from global sensitivity analysis [37].

A convenient proxy replacement for (7) is

$$\frac{1}{n} \sum_{\boldsymbol{k} \in \mathbb{N}_0^d \setminus \{\boldsymbol{0}\}} |\hat{f}(\boldsymbol{k})|^2 \mathbf{1}_{\{k_j > 0\}} = \frac{1}{n} \sum_{u \subseteq 1{:}d} \sigma_u^2 \mathbf{1}_{\{j \in u\}} \tag{8}$$

where $\sigma_u^2$ is the ANOVA variance component for the set $u$. The equality above follows because the ANOVA can be defined, as Sobol' [40] did, by collecting up the terms involving $x_j$ for $j \in u$ from the orthogonal decomposition. Sobol' used Haar functions. The right hand side of (8) equals $\overline{\tau}_j^2/n$. It counts all the variance components in which variable $j$ participates. From the orthogonality properties of ANOVA effects it follows that

$$\overline{\tau}_j^2 = \frac{1}{2} \int_{[0,1]^{d+1}} \big( f(\boldsymbol{z}_j{:}\boldsymbol{x}_{-j}) - f(\boldsymbol{x}) \big)^2 \, \mathrm{d}\boldsymbol{z}_j \, \mathrm{d}\boldsymbol{x}.$$

The Jansen estimator [19] is an estimate of the above integral that can be done by a $d+1$ dimensional MC or QMC or RQMC sampling algorithm. Our main interest in this Sobol' index estimator is that we use it as a point of comparison to the use of active subspaces in choosing a projection of a Gaussian vector along which to pre-integrate.

## 4 Active subspace method

Without loss of generality an expectation defined with respect to $\boldsymbol{x} \sim \mathcal{N}(0, \Sigma)$ for nonsingular $\Sigma \in \mathbb{R}^{d \times d}$ can be written as an expectation with respect to $\boldsymbol{x} \sim \mathcal{N}(0, I)$. For a unit vector $\theta \in \mathbb{R}^d$, we will pre-integrate over $\boldsymbol{x}^{\mathsf{T}} \theta \sim \mathcal{N}(0, 1)$ and then the problem is to make a principled choice of $\theta$. It would not be practical to seek an optimal choice.

Our proposal is to use active subspaces [3]. As mentioned in the introduction we let

$$C = \mathbb{E}(\nabla f(\boldsymbol{x}) \nabla f(\boldsymbol{x})^{\mathsf{T}})$$

and then let $\Theta[1{:}r]$ comprise the $r$ leading eigenvectors of $C$. The original use for active subspaces is to approximate $f(\boldsymbol{x}) \approx \tilde{f}(\Theta[1{:}r]^{\mathsf{T}} \boldsymbol{x})$ for some function $\tilde{f}$ on $\mathbb{R}^r$. It is well known that one can

construct functions where the active subspace will be a bad choice over which to approximate. For instance, with $f(\boldsymbol{x}) = \sin(10^6 x_1) + 100 x_2$ the $r = 1$ active subspace provides a function of $x_1$ alone while a function of $x_2$ alone can provide a better approximation than a function of $x_1$ alone can. Active subspaces remain useful for approximation because the motivating problems are not so pathological and there is a human in the loop to catch such things. They also have an enormous practical advantage that one set of evaluations of $\nabla f$ can be used in the search for $\Theta$ instead of having every candidate $\Theta$ require its own evaluations of $\nabla f$. Using active subspaces for integration retains that advantage.

In our setting, we take $r = 1$ and pre-integrate over $\theta^\mathsf{T}\boldsymbol{x}$ where $\theta$ is the leading eigenvector of $C$. That is $\theta$ maximizes $\theta^\mathsf{T}\mathbb{E}(\nabla f(\boldsymbol{x})\nabla f(\boldsymbol{x})^\mathsf{T})\theta$ over $d$ dimensional unit vectors. Now suppose that instead of using $f(\boldsymbol{x})$ we use $f_Q(\boldsymbol{x}) = f(Q\boldsymbol{x})$ for an orthogonal matrix $Q \in \mathbb{R}^{d\times d}$. Then $\mathbb{E}(\nabla f_Q(\boldsymbol{x})\nabla f_Q(\boldsymbol{x})^\mathsf{T}) = Q^\mathsf{T}CQ$ which is similar to $C$. It has the same eigenvalues and the leading eigenvector is $\tilde{\theta} = Q^\mathsf{T}\theta$. Furthermore, Theorem 3.1 of [47] shows that the invariance extends to the whole eigendecomposition.

## 4.1  Connection to a Sobol' index

The discussion in Section 3 motivates pre-integration of $f(\boldsymbol{x})$ for $\boldsymbol{x} \sim \mathcal{N}(0, I)$ over a linear combination $\theta^\mathsf{T}\boldsymbol{x}$ having the largest Sobol' index over unit vectors $\theta$. For $\theta_1 = \theta$, let $\Theta = (\theta_1, \theta_2, \ldots, \theta_d) \in \mathbb{R}^{d\times d}$ be an orthogonal matrix and write

$$f_\Theta(\boldsymbol{x}) = f(\Theta\boldsymbol{x}) = f(x_1\theta_1 + x_2\theta_2 + \cdots + x_d\theta_d).$$

Then we define $\bar{\tau}^2_\theta(f)$ to be $\bar{\tau}^2_1$ in the ANOVA of $f_\Theta(\boldsymbol{x})$. First we show that $\bar{\tau}^2_\theta$ does not depend on the last $d - 1$ columns of $\Theta$.

Let $z$, $\tilde{z}$ and $\boldsymbol{y}$ be independent with distributions $\mathcal{N}(0, 1)$, $\mathcal{N}(0, 1)$ and $\mathcal{N}(0, I_{d-1})$ respectively. Let $\boldsymbol{x} = z\theta_1 + \Theta_{-1}\boldsymbol{y}$ and $\tilde{\boldsymbol{x}} = \tilde{z}\theta_1 + \Theta_{-1}\boldsymbol{y}$. Using the Jansen formula for $\theta = \theta_1$,

$$\bar{\tau}^2_\theta = \frac{1}{2}\int_\mathbb{R}\int_\mathbb{R}\int_{\mathbb{R}^{d-1}} \big(f(z\theta_1 + \Theta_{-1}\boldsymbol{y}) - f(\tilde{z}\theta_1 + \Theta_{-1}\boldsymbol{y})\big)^2 \varphi(\boldsymbol{y})\,\mathrm{d}\boldsymbol{y}\varphi(\tilde{z})\,\mathrm{d}\tilde{z}\varphi(z)\,\mathrm{d}z. \tag{9}$$

Now for an orthogonal matrix $Q \in \mathbb{R}^{(d-1)\times(d-1)}$, let

$$\widetilde{\Theta} = \Theta\begin{pmatrix} 1 & \boldsymbol{0}^\mathsf{T}_{d-1} \\ \boldsymbol{0}_{d-1} & Q \end{pmatrix} = \begin{pmatrix} \tilde{\theta}_1 & \tilde{\theta}_2 & \cdots & \tilde{\theta}_d \end{pmatrix}$$

where $\tilde{\theta}_1 = \theta_1$. In this parameterization we get

$$\bar{\tau}^2_\theta = \frac{1}{2}\int_\mathbb{R}\int_\mathbb{R}\int_{\mathbb{R}^{d-1}} \big(f(z\theta_1 + \widetilde{\Theta}_{-1}\boldsymbol{y}) - f(\tilde{z}\theta_1 + \widetilde{\Theta}_{-1}\boldsymbol{y})\big)^2 \varphi(\boldsymbol{y})\,\mathrm{d}\boldsymbol{y}\varphi(\tilde{z})\,\mathrm{d}\tilde{z}\varphi(z)\,\mathrm{d}z$$

$$= \frac{1}{2}\int_\mathbb{R}\int_\mathbb{R}\int_{\mathbb{R}^{d-1}} \big(f(z\theta_1 + \Theta_{-1}Q\boldsymbol{y}) - f(\tilde{z}\theta_1 + \Theta_{-1}Q\boldsymbol{y})\big)^2 \varphi(\boldsymbol{y})\,\mathrm{d}\boldsymbol{y}\varphi(\tilde{z})\,\mathrm{d}\tilde{z}\varphi(z)\,\mathrm{d}z$$

which matches (9) after a change of variable. There is an even stronger invariance property in this setup. The random variable $\mathbb{E}(f_\Theta(\boldsymbol{x})\,|\,\boldsymbol{x}_{-1})$ (random because it depends on $\boldsymbol{x}_{-1}$) has a distribution that does not depend on $\theta_2, \ldots, \theta_d$.

**Theorem 4.1.** *Let $\boldsymbol{x} \sim \mathcal{N}(0, I_d)$ for $f$ with $\mathbb{E}(f(\boldsymbol{x})^2) < \infty$ and let $\Theta \in \mathbb{R}^{d\times d}$ be an orthogonal matrix with columns $\theta_j$ for $j = 1, \ldots, d$. Then the distribution of $\mathbb{E}(f_\Theta(\boldsymbol{x})\,|\,\boldsymbol{x}_{-1})$ does not depend on the last $d - 1$ columns of $\Theta$.*

11

*Proof.* Let $\Theta = (\theta_1 \ \Theta_{-1})$ and $\widetilde\Theta = (\theta_1 \ \widetilde\Theta_{-1})$ be orthogonal $d \times d$ matrices. Define

$$f_\Theta(\boldsymbol{x}) = f(\Theta\boldsymbol{x}) = f(x_1\theta_1 + \Theta_{-1}\boldsymbol{x}_{-1}) \quad \text{and} \quad f_{\widetilde\Theta}(\boldsymbol{x}) = f(\widetilde\Theta\boldsymbol{x}) = f(x_1\theta_1 + \widetilde\Theta_{-1}\boldsymbol{x}_{-1}).$$

Now $\Theta_{-1}\boldsymbol{x}_{-1}$ and $\widetilde\Theta_{-1}\boldsymbol{x}_{-1}$ both have the same $\mathcal{N}(0, I - \theta_1\theta_1^\mathsf{T})$ distribution independently of $x_1\theta_1$. Choose $A_1 \subset \mathbb{R}$ and $A_0$ in the support of $\mathcal{N}(0, I - \theta_1\theta_1^\mathsf{T})$ with positive probability under that (singular) distribution. Then $\Pr(f_\Theta(\boldsymbol{x}) \in A_1 \,|\, \Theta_{-1}\boldsymbol{x}_{-1} \in A_0) = \Pr(f_{\widetilde\Theta}(\boldsymbol{x}) \in A_1 \,|\, \widetilde\Theta_{-1}\boldsymbol{x}_{-1} \in A_0)$.

For $A \subset \mathbb{R}$ and $B \subset \mathbb{R}^{d-1}$,

$$
\begin{aligned}
\Pr(f_\Theta(\boldsymbol{x}) \in A \,|\, \boldsymbol{x}_{-1} \in B) &= \Pr(f_\Theta(\boldsymbol{x}) \in A \,|\, \Theta_{-1}\boldsymbol{x}_{-1} \in \Theta_{-1}B) \\
&= \Pr(f_{\widetilde\Theta}(\boldsymbol{x}) \in A \,|\, \widetilde\Theta_{-1}\boldsymbol{x}_{-1} \in \Theta_{-1}B) \\
&= \Pr(f_{\widetilde\Theta}(\boldsymbol{x}) \in A \,|\, \Theta_{-1}^\mathsf{T}\widetilde\Theta_{-1}\boldsymbol{x}_{-1} \in B)
\end{aligned}
$$

where $\Theta_{-1}^\mathsf{T}\widetilde\Theta_{-1}\boldsymbol{x}_{-1}$ has the same $\mathcal{N}(0, I_{d-1})$ distribution that $\boldsymbol{x}_{-1}$ has. Then for $C \subset [0, 1]$,

$$
\begin{aligned}
\Pr\big(\Pr(f_\Theta(\boldsymbol{x}) \in A \,|\, \boldsymbol{x}_{-1} \in B) \in C\big) &= \Pr\big(\Pr(f_{\widetilde\Theta}(\boldsymbol{x}) \in A \,|\, \Theta_{-1}^\mathsf{T}\widetilde\Theta_{-1}\boldsymbol{x}_{-1} \in B) \in C\big) \\
&= \Pr\big(\Pr(f_{\widetilde\Theta}(\boldsymbol{x}) \in A \,|\, \boldsymbol{x}_{-1} \in B) \in C\big).
\end{aligned}
$$

It follows that the distribution of $\Pr(f_{\widetilde\Theta}(\boldsymbol{x}) \in A \,|\, \boldsymbol{x}_{-1})$ is the same as the distribution of $\Pr(f_\Theta(\boldsymbol{x}) \in A \,|\, \boldsymbol{x}_{-1})$. Integrating over $\mathbb{R}$ we get that $\mathbb{E}(f_{\widetilde\Theta}(\boldsymbol{x}) \in A \,|\, \boldsymbol{x}_{-1})$ has the same distribution as $\mathbb{E}(f_\Theta(\boldsymbol{x}) \in A \,|\, \boldsymbol{x}_{-1})$ does. $\qquad\square$

Another consequence of Theorem 4.1 is that $\underline\tau_\theta^2(f) = \underline\tau_1^2(f_\Theta)$ is unaffected by $\theta_2, \ldots, \theta_d$. Because the variance of $f$ is unchanged by making an orthogonal matrix transformation of its inputs, the normalized Sobol' indices $\underline\tau_\theta^2/\sigma^2$ and $\overline\tau_\theta^2/\sigma^2$ are also invariant.

Finding the optimal $\theta$ would ordinarily require an expensive search because every estimate of $\overline\tau_\theta^2$, for a given $\theta$ would require its own collection of evaluations of $f$. Using a Poincaré inequality in [42] we can bound that Sobol' index by

$$\overline\tau_\theta^2(f) \leqslant \mathbb{E}((\theta^\mathsf{T}\nabla f(\boldsymbol{x}))^2) = \theta^\mathsf{T} C\theta.$$

The active subspace direction thus maximizes an upper bound on the Sobol' index for a projection. Next we develop a deeper correspondence between these two measures.

For a unit vector $\theta \in \mathbb{R}^d$, we can write $f(\boldsymbol{x}) = f(\theta\theta^\mathsf{T}\boldsymbol{x} + (I - \theta\theta^\mathsf{T})\boldsymbol{x})$. If $\boldsymbol{x}, \boldsymbol{z}$ are independent $\mathcal{N}(0, I)$ vectors then we can change the component of $\boldsymbol{x}$ parallel to $\theta$ by changing the argument of $f$ to be $\theta\theta^\mathsf{T}\boldsymbol{z} + (I - \theta\theta^\mathsf{T})\boldsymbol{x}$. This leaves the resulting point unchanged in the $d-1$ dimensional space orthogonal to $\theta$. Let $\tilde{x} = \theta^\mathsf{T}\boldsymbol{x}$ and $\tilde{z} = \theta^\mathsf{T}\boldsymbol{z}$. Then $\tilde{x}, \tilde{z} \sim \mathcal{N}(0, 1)$ and $(I - \theta\theta^\mathsf{T})\boldsymbol{x} \sim \mathcal{N}(0, I - \theta\theta^\mathsf{T})$ are all independent. If $f$ is differentiable, then by the mean value theorem

$$f(\theta\theta^\mathsf{T}\boldsymbol{z} + (I - \theta\theta^\mathsf{T})\boldsymbol{x}) - f(\theta\theta^\mathsf{T}\boldsymbol{x} + (I - \theta\theta^\mathsf{T})\boldsymbol{x}) = \theta^\mathsf{T}\nabla f(\theta\tilde{y} + (I - \theta\theta^\mathsf{T})\boldsymbol{x})(\tilde{z} - \tilde{x})$$

for a real number $\tilde{y}$ between $\tilde{x}$ and $\tilde{z}$. Using the Jansen formula, the Sobol' index for this projection is

$$\frac{1}{2}\theta^\mathsf{T}\mathbb{E}\Big((\tilde{z} - \tilde{x})^2 \nabla f(\theta\tilde{y} + (I - \theta\theta^\mathsf{T})\boldsymbol{x})\nabla f(\theta\tilde{y} + (I - \theta\theta^\mathsf{T})\boldsymbol{x})^\mathsf{T}\Big)\theta \tag{10}$$

which matches $\theta^\mathsf{T}\mathbb{E}(\nabla f(\boldsymbol{x})\nabla f(\boldsymbol{x})^\mathsf{T})\theta$ over a $d-1$ dimensional subspace but differs from it as follows. First, it includes a weight factor $(\tilde{z} - \tilde{x})^2$ that puts more emphasis on pairs of inputs where $\theta^\mathsf{T}\boldsymbol{x}$

and $\theta^\mathsf{T} z$ are far from each other. Second, the evaluation point projected onto $\theta$ equals $\tilde{y}$ which lies between two independent $\mathcal{N}(0,1)$ variables instead of having the $\mathcal{N}(0,1)$ distribution, and just where it lies between them depends on details of $f$ and there could be more than one such $\tilde{y}$ for some $f$. The formula simplifies in an illustrative way for quadratic functions $f$.

**Proposition 4.2.** *If $f : \mathbb{R}^d \to \mathbb{R}$ is a quadratic function and $\theta \in \mathbb{R}^d$ is a unit vector, then the Sobol' index $\overline{\tau}_\theta^2$ is*

$$\theta^\mathsf{T} \mathbb{E}\left( \nabla f\left(\frac{\theta\theta^\mathsf{T} x}{\sqrt{2}} + (I - \theta\theta^\mathsf{T})x\right)\nabla f\left(\frac{\theta\theta^\mathsf{T} x}{\sqrt{2}} + (I - \theta\theta^\mathsf{T})x\right)^\mathsf{T}\right)\theta. \tag{11}$$

*Proof.* If $f$ is quadratic, then $\tilde{y} = (\tilde{z} + \tilde{x})/2 \sim \mathcal{N}(0, 1/2)$ and $\tilde{z} - \tilde{x} \sim \mathcal{N}(0, 2)$ and $(I - \theta\theta^\mathsf{T})x$ are all independent. Then $\mathbb{E}((\tilde{z} - \tilde{x})^2) = 2$ and $\theta\tilde{y}$ has the same distribution as $\theta\theta^\mathsf{T} x/\sqrt{2}$ which is also independent of $(\tilde{z} - \tilde{x})$ and $(I - \theta\theta^\mathsf{T})x$. Making those substitutions in (10) yields (11). $\qquad\square$

The Sobol' index in equation (11) matches the quantity optimized by the first active subspace apart from the divisor $\sqrt{2}$ affecting one of the $d$ dimensions. We can also show directly that for $x \sim \mathcal{N}(0, I)$ and $f(x) = (1/2)x^\mathsf{T} A x + b^\mathsf{T} x$ for a symmetric matrix $A$, that the Sobol' criterion reduces to $\theta^\mathsf{T} A^2\theta + (\theta^\mathsf{T} b)^2 - (1/2)(\theta^\mathsf{T} A\theta)^2$ compared to an active subspace criterion of $\theta^\mathsf{T} A^2\theta + (\theta^\mathsf{T} b)^2$.

## 4.2 Active subspace

Because $C = \mathbb{E}(\nabla f(x)\nabla f(x)^\mathsf{T})$ is positive semi-definite (PSD), it has the eigen-decomposition $C = \Theta D \Theta^\mathsf{T}$, where $\Theta = (\theta_1, \ldots, \theta_d) \in \mathbb{R}^{d \times d}$ is an orthogonal matrix consisted of eigenvectors of $C$, and $D = \mathrm{diag}(\lambda_1, \ldots, \lambda_d)$ with $\lambda_1 \geqslant \ldots \geqslant \lambda_d \geqslant 0$ being the eigenvalues. Constantine et al. [4] prove that there exists a constant $c$ such that

$$\mathbb{E}\left(\left(f(x) - \mathbb{E}(f(x)\,|\,\Theta[1{:}r]^\mathsf{T} x)\right)^2\right) \leqslant c(\lambda_{r+1} + \cdots + \lambda_d) \tag{12}$$

for all $f$ with a square integrable gradient. In general, the Poincaré constant $c$ depends on the support of the function and the probability measure. But for multivariate standard Gaussian distribution, the Poincaré constant is always 1 [2, 35]. This is because $\Theta[1{:}r]^\mathsf{T} x$ and $\Theta[-(1{:}r)]^\mathsf{T} x$ are independent standard Gaussian variables thus

$$\mathbb{E}\left(\left(f(x) - \mathbb{E}(f(x)\,|\,\Theta[1{:}r]^\mathsf{T} x)\right)^2 \,|\, \Theta[1{:}r]^\mathsf{T} x\right) \leqslant \lambda_{r+1} + \cdots + \lambda_d$$

for all $\Theta[1{:}r]^\mathsf{T} x$.

In our problem, we take $r = 1$ and use $\mathbb{E}(f(x)\,|\,\theta^\mathsf{T} x)$ where $\theta$ is the first column of $\Theta$. Because we will end up with a $d - 1$ dimensional integration problem it is convenient in an implementation to make $\theta^\mathsf{T} x$ the last variable not the first. For instance, one would use the first $d - 1$ components in a Sobol' sequence not components 2 through $d$. Taking $\theta$ to be the first column of $\Theta$, we compute with

$$g(x_{-d}) = \int_{-\infty}^{\infty} f(\theta x_d + \Psi x_{-d})\, \mathrm{d}x_d$$

using a quadrature rule of negligible error or a closed form expression if a suitable one is available for an orthonormal matrix $\Psi \in \mathbb{R}^{d \times (d-1)}$ that is orthogonal to $\theta$. We then integrate this $g$ over $d - 1$ variables by RQMC. We can use $\Psi = \Theta[2{:}d]$. Or if we want to avoid the cost of computing

---
**Algorithm 1:** pre-integration with active subspace

---
**Input** : Integrand $f$, number of samples $M$ to compute $\hat{C}$, number of samples $n$ to compute $\hat{\mu}$

**Output:** An estimate $\hat{\mu}$ of $\int_{\mathbb{R}^d} f(\boldsymbol{x})\varphi(\boldsymbol{x})d\boldsymbol{x}$

/* Find active subspaces                                                        */

Take $\boldsymbol{x}_0, \ldots, \boldsymbol{x}_{M-1} \sim \mathcal{N}(0, I_d)$ by RQMC.

Compute $\widehat{C} = \frac{1}{M} \sum_{i=0}^{M-1} \nabla f(\boldsymbol{x}_i)\nabla f(\boldsymbol{x}_i)^{\mathsf{T}}$.

Compute the eigen-decomposition $\widehat{C} = \widehat{\Theta}\widehat{D}\widehat{\Theta}^{\mathsf{T}}$.

/* Pre-integration                                                              */

Let $\theta$ be the first column of $\widehat{\Theta}$

Compute the pre-integrated function $\mathbb{E}(f(\boldsymbol{x})|\theta^{\mathsf{T}}\boldsymbol{x})$ by a closed form or quadrature rule

/* RQMC integration                                                             */

Take $\boldsymbol{x}_0, \ldots, \boldsymbol{x}_{n-1} \sim \mathcal{N}(0, I_{d-1})$ by RQMC.

Let $\hat{\mu} = \frac{1}{n} \sum_{i=0}^{n-1} \tilde{f}(\boldsymbol{x}_i)$.

---

the full eigendecomposition of $C$ we can find $\theta_1$ by a power iteration and then use a Householder transformation

$$\Theta = I - 2\boldsymbol{w}\boldsymbol{w}^{\mathsf{T}}, \quad \text{where} \quad \boldsymbol{w} = \frac{\theta - e_1}{\|\theta - e_1\|}$$

and $e_1 = (1, 0, 0, \ldots, 0)^{\mathsf{T}}$. This $\Theta$ is an orthogonal matrix whose first column is $\theta$ and again we can choose $\Psi = \Theta[2{:}d]$. In our numerical work, we have used $\Theta[2{:}d]$ instead of the Householder transformation because of the effective dimension motivation for those eigenvectors given by [45].

In practice, we must estimate $C$. In the above description, we replace $C$ by

$$\widehat{C} = \frac{1}{M} \sum_{i=0}^{M-1} \nabla f(\boldsymbol{x}_i)\nabla f(\boldsymbol{x}_i)^{\mathsf{T}}, \tag{13}$$

for an RQMC generated sample with $\boldsymbol{x}_i \sim \mathcal{N}(0, I_d)$ and then define $\theta$ and $\Theta_{-1}$ using $\widehat{C}$ in place of $C$. We summarize the procedure in Algorithm 1.

Using our prior notation we can now describe the approach of [44] more precisely. They first pre-integrate one variable in closed form producing a $d-1$ dimensional integrand. They then apply gradient GPCA to the pre-integrated function to find a good $d-1$ dimensional rotation matrix. [4]. That is, they first find $h(\boldsymbol{x}_{2:d}) := \mathbb{E}(f(\boldsymbol{x})|\boldsymbol{x}_{2:d})$, then compute

$$\widehat{\widetilde{C}} = \frac{1}{M} \sum_{i=0}^{M-1} \nabla h(\boldsymbol{x}_i)\nabla h(\boldsymbol{x}_i)^{\mathsf{T}} \in \mathbb{R}^{(d-1)\times(d-1)}, \quad \boldsymbol{x}_i \sim \mathcal{N}(0, I_{d-1}), \tag{14}$$

using RQMC points $\boldsymbol{x}_i$. Then they find the eigen-decomposition $\widehat{\widetilde{C}} = \widehat{V}\widehat{\Lambda}\widehat{V}^{\mathsf{T}}$. Finally, they use RQMC to integrate the function $h(\widehat{V}\boldsymbol{x})$ where $\boldsymbol{x} \sim \mathcal{N}(0, I_{d-1})$. The main difference is that they apply pre-integration to the original integrand $f(\boldsymbol{x})$ while we apply pre-integration to the rotated integrand $f_\Theta(\boldsymbol{x}) = f(\Theta\boldsymbol{x})$. They conduct GPCA in the end as an approach to reduce effective dimension, while we conduct a similar GPCA (active subspace method) at the beginning to find the important subspace.

# 5 Application to option pricing

Here we study some Gaussian integrals arising from financial valuation. We assume that an asset price $S_t$, such as a stock, follows a geometric Brownian motion satisfying the stochastic differential equation (SDE)

$$\mathrm{d}S_t = rS_t\,\mathrm{d}t + \sigma S_t\,\mathrm{d}B_t,$$

where $B_t$ is a Brownian motion. Here, $r$ is the interest rate and $\sigma > 0$ is the constant volatility for the asset. For an initial price $S_0$, the SDE above has a unique solution

$$S_t = S_0\exp\Big(\Big(r - \frac{\sigma^2}{2}\Big)t + \sigma B_t\Big).$$

Suppose the maturity time of the option is $T$. In practice, we simulate discrete Brownian motion. We call $B$ a $d$-dimensional discrete Brownian motion if $B$ follows a multivariate Gaussian distribution with mean zero and covariance $\Sigma$ with $\Sigma_{ij} = \Delta t\min(i,j)$, where $\Delta t = T/d$ is the length of each time interval and $1 \leqslant i, j \leqslant d$. To sample a discrete Brownian motion, we can first find a $d \times d$ matrix $R$ such that $RR^{\mathsf{T}} = \Sigma$, then generate a standard Gaussian variable $\boldsymbol{z} \sim \mathcal{N}(0, I_d)$, and let $B = R\boldsymbol{z}$. Taking $R$ to be the lower triangular matrix in the Cholesky decomposition of $\Sigma$ yields the *standard construction*. Using the usual eigen-decomposition $\Sigma = U\Lambda U^{\mathsf{T}}$, we can take $R = U\Lambda^{1/2}$. This is called the *principal component analysis (PCA)* construction. For explicit forms of both these choices of $R$, see [10].

## 5.1 Option with one asset

When we use the matrix $R$, we can approximate $S_{j\Delta t}$ by

$$S_j = S_0\exp\Big(\Big(r - \frac{\sigma^2}{2}\Big)j\Delta t + \sigma B_j\Big), \quad 1 \leqslant j \leqslant d$$

where $B = R\boldsymbol{z}$ is the discrete Brownian motion. The arithmetic average of the stock price is given by

$$\bar{S}(R, \boldsymbol{z}) = \frac{S_0}{s}\sum_{j=1}^{d}\exp\Big(\Big(r - \frac{\sigma^2}{2}\Big)j\Delta t + \sigma\sum_{k=1}^{d}R_{jk}z_k\Big).$$

Then the expected payoff of the arithmetic average Asian call option with strike price $K$ is $\mathbb{E}\big((\bar{S}(A, \boldsymbol{z}) - K)_+\big)$, where the expectation is taken over $\boldsymbol{z} \sim \mathcal{N}(0, I_d)$.

Suppose that we want to marginalize over $z_1$ before computing the expectation $\mathbb{E}((\bar{S}(A, \boldsymbol{z}) - K)_+)$. If $R_{j1} > 0$ for all $1 \leqslant j \leqslant d$, then $\bar{S}(R, \boldsymbol{z})$ is increasing in $z_1$ for any value of $\boldsymbol{z}_{2:s}$. If we can find $\gamma = \gamma(\boldsymbol{z}_{2:s})$ such that

$$\bar{S}(R, (\gamma, \boldsymbol{z}_{2:s})) = K, \tag{15}$$

then the pre-integration step becomes

$$
\begin{aligned}
&\mathbb{E}((\bar{S}(A, \boldsymbol{z}) - K)_+ \,|\, \boldsymbol{z}_{2:s}) \\
&= \int_{z_1 \geqslant \gamma(\boldsymbol{z}_{2:s})} (\bar{S}(R, (z_1, \boldsymbol{z}_{2:s})) - K)\varphi(z_1)\,\mathrm{d}z_1 \\
&= \frac{S_0}{d} \sum_{j=1}^{d} \exp\Big(\Big(r - \frac{\sigma^2}{2}\Big)j\Delta t + \sigma \sum_{k=2}^{d} R_{jk} z_k + \frac{\sigma^2 R_{j1}^2}{2}\Big) \bar{\Phi}(\gamma - \sigma R_{j1}) - K\bar{\Phi}(\gamma), \quad (16)
\end{aligned}
$$

where $\bar{\Phi}(x) = 1 - \Phi(x)$. In practice, Equation (15) can be solved by a root finding algorithm. For example, Newton iteration usually converges in only a few steps.

The condition that $R_{j1} > 0$ for $1 \leqslant j \leqslant d$ is satisfied when we use the standard construction or PCA construction of Brownian motion. Using the active subspace method, we are using $\tilde{R} = R\Theta$ in the place of $R$ where $\Theta$ consists of the eigenvectors of $C = \mathbb{E}(\nabla f(\boldsymbol{z})\nabla f(\boldsymbol{z})^{\mathsf{T}})$, and here $f(\boldsymbol{z}) = (\bar{S}(A, \boldsymbol{z}) - K)_+$. If every $\tilde{R}_{j1}$ is negative then we replace $\tilde{R}$ by $-\tilde{R}$. We have not proved that the components of the first column of $\tilde{R}$ must all have the same sign, but that has always held for the integrands in our simulations. As a result we have not had to use a numerical quadrature.

We compare Algorithm 1 with other methods in the option pricing example considered in [44] and [14]. Apart from the payoff function of call option, we also consider the Greeks: Delta, Gamma, Rho, Theta, and Vega. These are defined in [44]. We take the parameters $d = 50$, $T = 1$, $\sigma = 0.4$, $r = 0.1$, $S_0 = K = 100$ the same as in [14]. We consider 4 methods:

- `AS+pre-int`: our proposed active subspace pre-integration method (Algorithm 1), which applies active subspace method to find the direction to pre-integrate,

- `pre-int+DimRed`: the method proposed in [44], which first pre-integrates $z_1$ and applies GPCA to conduct dimension reduction for the other $d - 1$ variables,

- `pre-int`: pre-integrating $z_1$ with no dimension reduction,

- `RQMC`: usual RQMC, and

- `MC`: plain Monte Carlo.

We vary $n$ from $2^3$ to $2^{17}$. For each $n$, we repeat the simulation 50 times and compute the root mean squared error (RMSE) defined by

$$
\sqrt{\frac{1}{50} \sum_{k=1}^{50} (\hat{\mu}^{(k)} - \mu)^2},
$$

where $\hat{\mu}^{(k)}$ is the estimate in the $k$-th replicate. The true value $\mu$ is approximated by applying `pre-int+DimRed` using PCA construction with $n = 2^{17}$ RQMC samples and averaging over 30 independent replicates. We chose this one because it works well and we wanted to avoid using `AS-pre-int` in case the model used for ground truth had some advantage in reported accuracy.

The root mean square error (RMSE) is plotted versus the sample size on the log-log scale. For the methods `AS+pre-int` and `pre-int+DimRed`, we use $M = 128$ samples to estimate $C$ as in (13).

16

We approximate the gradients of the original integrand and the pre-integrated integrand by the finite difference

$$\nabla f(x) \approx \left( \frac{f(x + \varepsilon \mathbf{e}_1) - f(x)}{\varepsilon}, \ldots, \frac{f(x + \varepsilon \mathbf{e}_{d-1}) - f(x)}{\varepsilon} \right)^{\mathsf{T}}, \quad \varepsilon = 10^{-6},$$

matching the choice in [44]. We chose a small value of $M$ to keep the costs comparable to plain RQMC. Also because $\theta$ is a local optimum of $\theta^{\mathsf{T}} C \theta$, we have $\hat{\theta}^{\mathsf{T}} C \hat{\theta} = \theta^{\mathsf{T}} C \theta + O(\|\hat{\theta} - \theta\|^2)$ so there are diminishing returns to accurate estimation of $\theta$. Finally, any $\theta$ where $f$ varies along $\theta^{\mathsf{T}} \boldsymbol{x}$ brings a variance reduction.

We consider both the standard construction (Figure 1) and the PCA construction (Figure 2) of Brownian motion. Several observations are in order:

(a) With the standard construction, `AS+pre-int` dominates all the other methods for five of the six test functions and is tied for best with `pre-int+DimRed` for the other one (Gamma).

(b) With the PCA construction, `AS+pre-int`, `pre-int+DimRed` and `pre-int` are the best methods for the payoff, Delta, Theta, and Vega and are nearly equally good.

(c) For Rho, `pre-int+DimRed` and `pre-int` are best, while for Gamma, `AS+pre-int` is best.

The performance of active subspace pre-integration is the same under either the standard or the PCA construction by invariance. For these Asian options it is already well known that the PCA is especially effective. Active subspace pre-integration finds something almost as good without special knowledge, coming out better in one example, worse in another and essentially the same in the other four.

## 5.2    Basket option

A basket option depends on the weighted average of several assets. Suppose that the $L$ assets $S^{(\ell)}, \ldots, S^{(L)}$ follow the SDE

$$\mathrm{d}S_t^{(\ell)} = r_\ell S_t^{(\ell)} \, \mathrm{d}t + \sigma_\ell S_t^{(\ell)} \, \mathrm{d}B_t^{(\ell)},$$

where $\{B^{(\ell)}\}_{1 \leqslant \ell \leqslant L}$ are standard Brownian motions with correlation

$$\mathrm{Corr}(B_t^{(\ell)}, B_t^{(k)}) = \rho_{\ell k}$$

for all $t > 0$. For some nonnegative weights $w_1 + \ldots + w_L = 1$, the payoff function of the Asian basket call option is given by

$$\left( \sum_{\ell=1}^{L} w_\ell \bar{S}^{(\ell)} - K \right)_+$$

where $\bar{S}^{(\ell)}$ is the arithmetic average of $S_t^{(\ell)}$ in the time interval $[0, T]$. Here, we only consider $L = 2$ assets. To generate $B^{(1)}, B^{(2)}$ with correlation $\rho$, we can generate two independent standard Brownian motions $W^{(1)}, W^{(2)}$ letting

$$B^{(1)} = W^{(1)} \quad \text{and} \quad B^{(2)} = \rho W^{(1)} + \sqrt{1 - \rho^2} W^{(2)}.$$
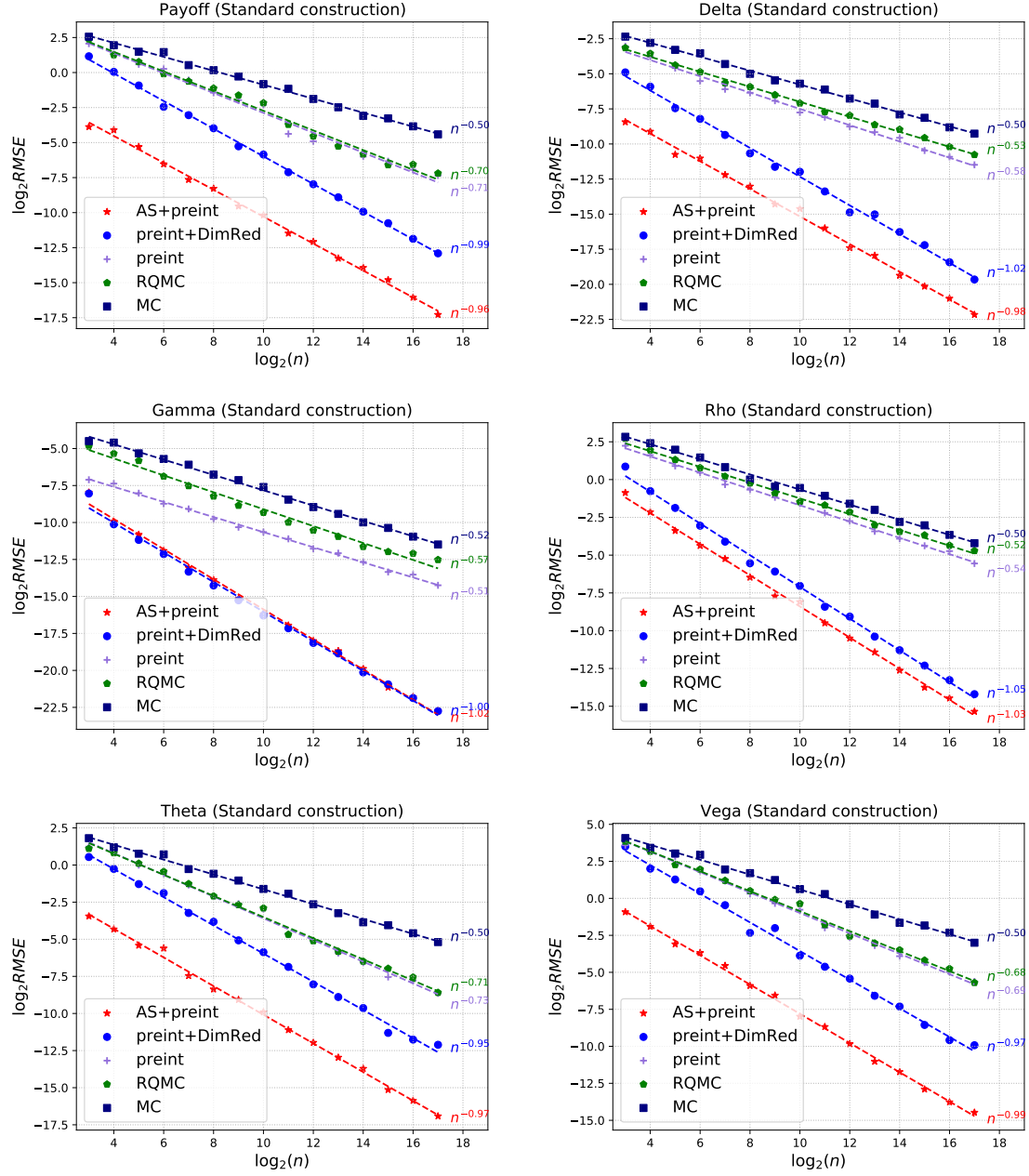
17

Figure 1: Single asset option. Standard construction of Brownian motion with $d = 50$. This and subsequent figures include least squares estimated slopes on the log–log scale.
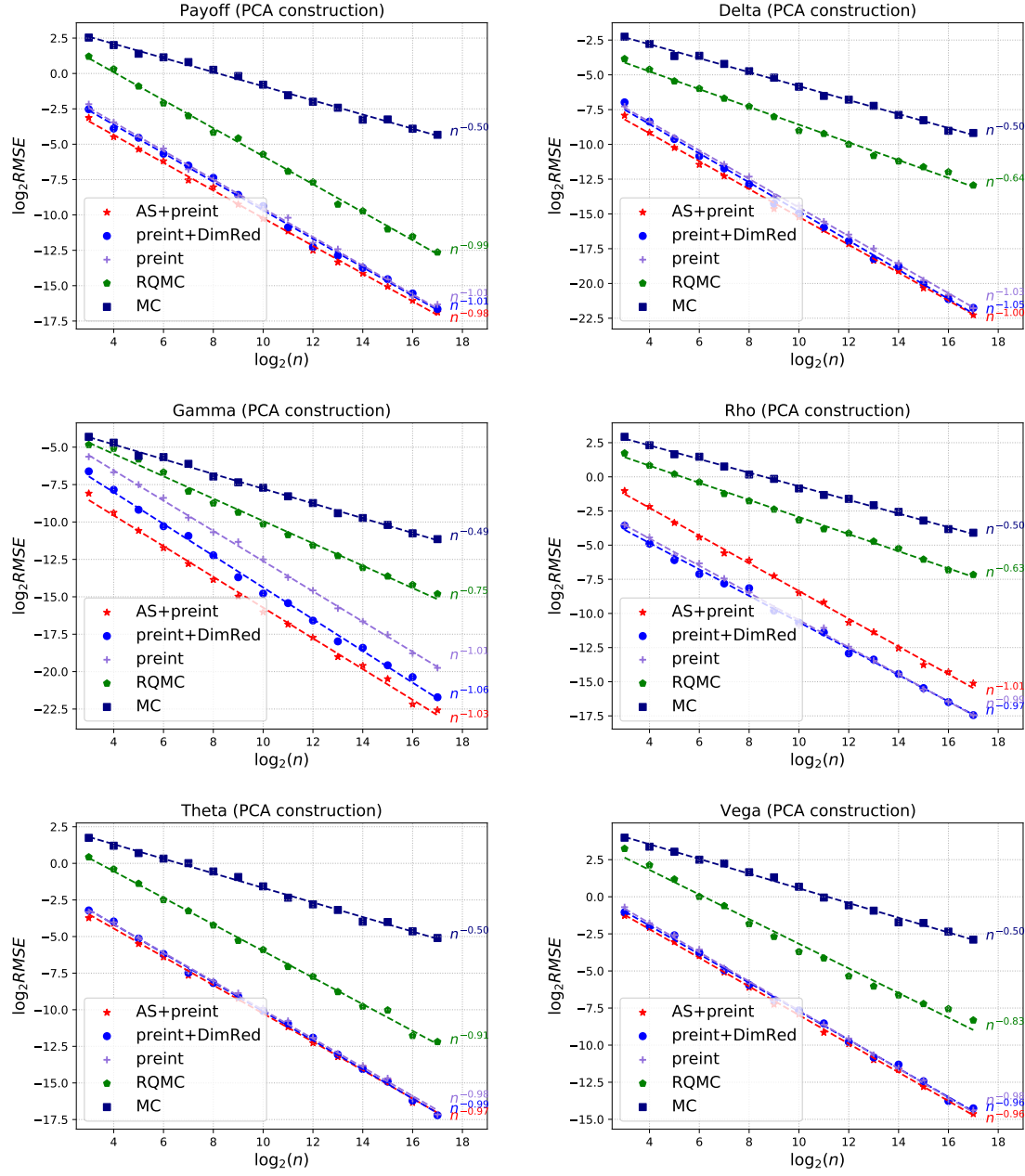
Figure 2: Single asset option. PCA construction of Brownian motion with $d = 50$.

Following the same discretization as before, we can generate $(\boldsymbol{z}^{\mathsf{T}}, \tilde{\boldsymbol{z}}^{\mathsf{T}}) \sim \mathcal{N}(0, I_{2d})$. Then for time steps $j = 1, \ldots, d$, let

$$S_j^{(1)} = S_0^{(1)} \exp\left(\left(r_1 - \frac{\sigma_1^2}{2}\right) j \Delta t + \sigma_1 \sum_{k=1}^{d} R_{jk} z_k\right), \quad \text{and}$$

$$S_j^{(2)} = S_0^{(2)} \exp\left(\left(r_2 - \frac{\sigma_2^2}{2}\right) j \Delta t + \sigma_2 \left(\rho \sum_{k=1}^{d} R_{jk} z_k + \sqrt{1 - \rho^2} \sum_{k=1}^{d} R_{jk} \tilde{z}_k\right)\right).$$

Again, the $d \times d$ matrix $R$ can be constructed by the standard construction or the PCA construction. Thus, the expected payoff can be written as

$$\mathbb{E}\left(\left(\frac{w_1}{d} \sum_{j=1}^{d} S_j^{(1)} + \frac{w_2}{d} \sum_{j=1}^{d} S_j^{(2)} - K\right)_+\right),$$

where the expectation is taken over $(\boldsymbol{z}^{\mathsf{T}}, \tilde{\boldsymbol{z}}^{\mathsf{T}})^{\mathsf{T}} \sim \mathcal{N}(0, I_{2d})$.

In the pre-integration step, we choose to integrate out $z_1$. This can be easily carried out as in equations (15) and (16) provided that the first column of $\tilde{R}$ is nonnegative. We take $d = 50$, $T = 1$, $\rho = 0.5$, $S_0^{(1)} = S_0^{(2)} = 100$ and $K = 95$. The RMSE values are plotted in Figure 3. In the left panel, we take $r_1 = 0.1$, $r_2 = 0.2$, $\sigma_1 = 0.2$, $\sigma_2 = 0.4$, $w_1 = 0.8$ and $w_2 = 0.2$. In the right panel, we take $r_1 = 0.2$, $r_2 = 0.1$, $\sigma_1 = 0.4$, $\sigma_2 = 0.2$, $w_1 = 0.2$ and $w_8 = 0.8$. This reverses the roles of the two assets which will make a difference when one pre-integrates over the first of the $2d$ inputs. We use $M = 128$ as before. In the top row of Figure 3, the matrix $R$ is obtained by the standard construction, while in the bottom row, $R$ is obtained by the PCA construction.

A few observations are in order:

(a) For the standard construction, pre-integrating over $z_1$ brings little improvement over plain RQMC. But for PCA construction, pre-integrating over $z_1$ brings a big variance reduction.

(b) The dimension reduction technique from [44] largely improves the RMSE from pre-integration without dimension reduction. This improvement is particularly significant for the standard construction.

(c) Active subspace pre-integration, `AS+preint`, has the best performance for both the standard and PCA constructions. It is even better than pre-integrating out the first principal component of Brownian motion with dimension reduction. In this example, the active subspace method is able to find a better direction than the first principal component over which to pre-integrate.

This problem required $L = 2$ Brownian motions and the above examples used the same decomposition to sample them both. A sharper principal components analysis would merge the two Brownian motions into a single $2d$-dimensional process and use the principal components from their joint covariance matrix. We call this the full PCA construction as described next.

The processes $B^{(1)} = \sigma_1 R \boldsymbol{z}$, and $B(2) = \sigma_2 R(\rho \boldsymbol{z} + \sqrt{1 - \rho^2} \tilde{\boldsymbol{z}})$, have joint distribution

$$\begin{pmatrix} B^{(1)} \\ B^{(2)} \end{pmatrix} \sim \mathcal{N}\left(0, \begin{pmatrix} \sigma_1^2 \Sigma & \rho \sigma_1 \sigma_2 \Sigma \\ \rho \sigma_1 \sigma_2 \Sigma & \sigma_2^2 \Sigma \end{pmatrix}\right).$$
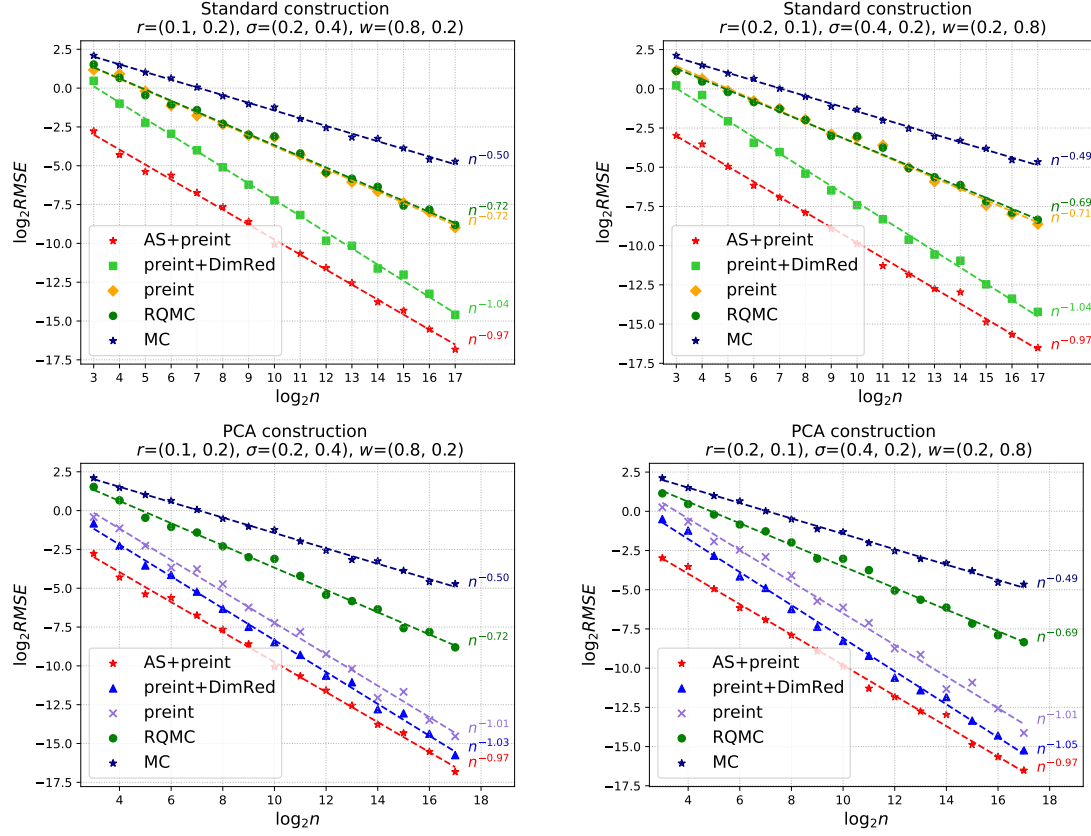
Figure 3: Basket option pricing.

Let $\widetilde{\Sigma}$ be the joint covariance matrix above. An alternative generation method is to pick $\tilde{R}$ with $\tilde{R}\tilde{R}^{\mathsf{T}} = \tilde{\Sigma}$, and let

$$\begin{pmatrix} B^{(1)} \\ B^{(2)} \end{pmatrix} = \tilde{R} \begin{pmatrix} \boldsymbol{z} \\ \tilde{\boldsymbol{z}} \end{pmatrix}, \quad \text{where} \quad \begin{pmatrix} \boldsymbol{z} \\ \tilde{\boldsymbol{z}} \end{pmatrix} \sim \mathcal{N}(0, I_{2d}).$$

The matrix $\tilde{R}$ can be found by either a Cholesky decomposition or eigendecomposition of $\tilde{R}$. We call this method full standard construction or full PCA construction. Taking $B^{(1)} = \sigma_1 R\boldsymbol{z}$, and $B^{(2)} = \sigma_2 R(\rho \boldsymbol{z} + \sqrt{1 - \rho^2}\tilde{\boldsymbol{z}})$, it is equivalent to taking

$$\tilde{R} = \begin{pmatrix} \sigma_1 I_d & \mathbf{0} \\ \sigma_2 \rho I_d & \sqrt{1 - \rho^2}\sigma_2 I_d \end{pmatrix} \begin{pmatrix} R & \mathbf{0} \\ \mathbf{0} & R \end{pmatrix}.$$

We call this the ordinary PCA or standard construction depending on whether $R$ is from the PCA or standard construction.

Figure 4 compares results using this full PCA construction. All methods apply pre-integration before using RQMC. For active subspace pre-integration, we pre-integrate along the direction in
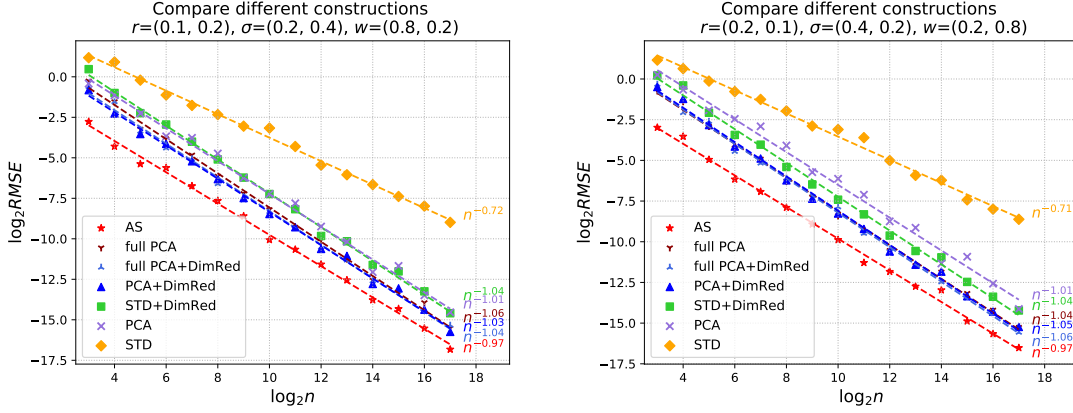
21

Figure 4: These plots compare pre-integration strategies for two basket options. We compare active subspace pre-integration to strategies with full and partial PCA and standard contructions of Brownian motion.

$\mathbb{R}^{2d}$ found by the active subspace method. For "full PCA", we pre-integrate along the first principal component of $\widehat{\Sigma}$. We can see that active subspace pre-integration has a better RMSE than even the full PCA construction with dimension reduction, which we consider to be the natural generalization of the PCA construction to this setting.

## 5.3 Timing

Pre-integration changes the computational cost of RQMC and this has not been much studied in the literature. Our figures compare the RMSE of different samplers as a function of the number $n$ of evaluations. Efficiency depends also on running times. Here we make study of running times for the pre-integration methods we studied. It is brief in order to conserve space. Running times depend on implementation details that could make our efficiency calculations differ from others'. All of timings we report were conducted on a MacBook Pro with 8 cores and 16GB memory. We simulated them all with 10 replicates and $n = 2^{15}$. The standard errors of the average evaluation times were negligible, about 0.3%–0.8% of the corresponding means. We also computed times to find $\widehat{C}$ and $\widehat{\widetilde{C}}$ (defined in equations (13) and (14)). Those had small standard errors too, and their costs are a small fraction of the total. The costs of finding the eigendecompositons are negligible in our examples.

Table 1 shows the results of our timings. For the 6 integrands, we find that pre-integration raises the cost of computation by roughly 12 to 16-fold. That extra cost is not enough to offset the gains from pre-integration with large $n$ except for pre-integration of the first component in the standard construction. That variable is not very important and we might well have expected pre-integration to be unhelpful for it.

Most of the pre-integrated computations took about the same amount of time but plain pre-integration and pre-integration with dimension reduction take slighly more time for the standard construction. Upon investigation we found that those methods took slightly more Newton iterations on average. In hindsight this makes sense. The Newton iterations started at 0. In the standard

|  | MC | RQMC | pre-int | AS+pre-int | pre-int+DimRed | $\widehat{C}$ | $\widehat{\widetilde{C}}$ |
|---|---|---|---|---|---|---|---|
| Payoff | 0.6 | 0.6 | 7.5 | 6.9 | 7.4 | 0.3 | 2.9 |
|  | 0.6 | 0.6 | 6.9 | 6.9 | 6.9 | 0.3 | 2.7 |
| Delta | 0.5 | 0.4 | 5.7 | 5.2 | 5.7 | 0.2 | 2.2 |
|  | 0.5 | 0.4 | 5.1 | 5.2 | 5.2 | 0.2 | 2.0 |
| Gamma | 0.6 | 0.6 | 12.2 | 12.2 | 12.3 | 0.2 | 4.7 |
|  | 0.6 | 0.6 | 11.7 | 12.3 | 11.7 | 0.2 | 4.5 |
| Rho | 0.9 | 0.9 | 10.6 | 10.1 | 10.6 | 0.3 | 4.1 |
|  | 0.9 | 0.9 | 10.1 | 10.2 | 10.1 | 0.4 | 3.9 |
| Theta | 1.0 | 1.0 | 13.9 | 13.3 | 14.0 | 0.4 | 5.4 |
|  | 1.0 | 1.0 | 13.4 | 13.4 | 13.4 | 0.4 | 5.2 |
| Vega | 0.6 | 0.6 | 9.7 | 9.1 | 9.7 | 0.3 | 3.7 |
|  | 0.6 | 0.6 | 9.1 | 9.1 | 9.1 | 0.3 | 3.5 |

Table 1: Time in seconds to compute all Asian option integrands and methods. All simulations take $2^{15}$ samples and the times are averaged over 10 replicates. Each integrand has two rows, where the top row uses the standard construction and the bottom row uses the PCA construction. The right two columns are the times used to estimate $\widehat{C}$ and $\widehat{\widetilde{C}}$ by $M = 128$ samples.

construction, the first variable does not have a very strong effect on the option value, requiring it to take larger values to reach the positivity threshold $x_1 = \gamma(\boldsymbol{x}_{2:d})$ and hence (a few) more iterations.

Another component of cost for active subspace methods is in computing an approximation to $C$ and $\widetilde{C}$. We used $M = 128$ function evaluations. Those are about $2d$ times as expensive as evaluations because they use divided differences. One advantage of active subspace pre-integration is that it uses divided differences of the original integrand. Pre-integration plus dimension reduction requires divided differences of the pre-integrated function with associated Newton searches, and that costs more.

## 6 Discussion

In this paper we have studied a kind of conditional RQMC known as pre-integration. We found that, just like conditional MC, the procedure can reduce variance but cannot increase it. We proposed to pre-integrate over the first component in an active subspace approximation, which is also known as the gradient PCA approximation. We showed a close relationship between this choice of pre-integration variable and what one would get using a computationally infeasible but well motivated choice by maximizing the Sobol' index of a linear combination of variables.

In the numerical examples of option pricing, we see that active subspace pre-integration achieves a better RMSE than previous methods when using the standard construction of Brownian motion. For the PCA construction, the proposed method has comparable accuracy in four of six cases, is better once and is worse once. For those six integrands, the PCA construction is already very good. Active subspace pre-integration still provides an automatic way to choose the pre-integration

direction. Even using the standard construction it is almost as good as pre-integration with the PCA construction. It can be used in settings where there is no strong incumbent decomposition analogous to the PCA for the Asian option. We saw it perform well for basket options. We even saw an improvement for Gamma in the very well studied case of the Asian option with the PCA construction.

We note that active subspaces use an uncentered PCA analysis of the matrix of sample gradients. One could use instead a centered analysis of $\mathbb{E}((\nabla f(\boldsymbol{x}) - \eta)(\nabla f(\boldsymbol{x}) - \eta)^\mathsf{T})$ where $\eta = \mathbb{E}(\nabla f(\boldsymbol{x}))$. The potential advantage of this is that $\nabla f(\boldsymbol{x}) - \eta$ is the gradient of $f(\boldsymbol{x}) - \eta^\mathsf{T} \boldsymbol{x}$ which subtracts a linear approximation from $f$ before searching for $\theta$. The rationale for this alternative is that RQMC might already do well integrating a linear function and we would then want to choose a direction $\theta$ that performs well for the nonlinear part of $f$. In our examples, we found very little difference between the two methods and so we proposed the simpler uncentered active subspace pre-integration.

# Acknowledgments

# References

[1] P. Acworth, M. Broadie, and P. Glasserman. A comparison of some Monte Carlo techniques for option pricing. In H. Niederreiter, P. Hellekalek, G. Larcher, and P. Zinterhof, editors, *Monte Carlo and quasi-Monte Carlo methods '96*, pages 1–18. Springer, 1997.

[2] L. H. Y. Chen. An inequality for the multivariate normal distribution. *Journal of Multivariate Analysis*, 12(2):306–315, 1982.

[3] P. G. Constantine. *Active subspaces: Emerging ideas for dimension reduction in parameter studies*. SIAM, Philadelphia, 2015.

[4] P. G. Constantine, E. Dow, and Q. Wang. Active subspace methods in theory and practice: applications to kriging surfaces. *SIAM Journal on Scientific Computing*, 36(4):A1500–A1524, 2014.

[5] P. J. Davis and P. Rabinowitz. *Methods of Numerical Integration*. Academic Press, San Diego, 2nd edition, 1984.

[6] J. Dick, F. Y. Kuo, and I. H. Sloan. High-dimensional integration: the quasi-Monte Carlo way. *Acta Numerica*, 22:133–288, 2013.

[7] J. Dick and F. Pillichshammer. *Digital sequences, discrepancy and quasi-Monte Carlo integration*. Cambridge University Press, Cambridge, 2010.

[8] B. Efron and C. Stein. The jackknife estimate of variance. *Annals of Statistics*, 9(3):586–596, 1981.

[9] A. D. Gilbert, F. Y. Kuo, and I. H. Sloan. Preintegration is not smoothing when monotonicity fails. Technical report, arXiv:2112.11621, 2021.

[10] P. Glasserman. *Monte Carlo Methods in Financial Engineering*, volume 53. Springer, 2004.

[11] P. Glasserman, P. Heidelberger, and P. Shahabuddin. Asymptotically optimal importance sampling and stratification for pricing path-dependent options. *Mathematical finance*, 9(2):117–152, 1999.

[12] A. Griewank, F. Y. Kuo, H. Leövey, and I. H. Sloan. High dimensional integration of kinks and jumps—smoothing by preintegration. *Journal of Computational and Applied Mathematics*, 344:259–274, 2018.

[13] J. M. Hammersley. Conditional Monte Carlo. *Journal of the ACM (JACM)*, 3(2):73–76, 1956.

[14] Z. He. On the error rate of conditional quasi–Monte Carlo for discontinuous functions. *SIAM Journal on Numerical Analysis*, 57(2):854–874, 2019.

[15] F. J. Hickernell. Koksma-Hlawka inequality. *Wiley StatsRef: Statistics Reference Online*, 2014.

[16] W. Hoeffding. A class of statistics with asymptotically normal distribution. *Annals of Mathematical Statistics*, 19(3):293–325, 1948.

[17] C. R. Hoyt and A. B. Owen. Mean dimension of ridge functions. *SIAM Journal on Numerical Analysis*, 58(2):1195–1216, 2020.

[18] J. Imai and K. S. Tan. Minimizing effective dimension using linear transformation. In *Monte Carlo and Quasi-Monte Carlo Methods 2002*, pages 275–292. Springer, 2004.

[19] M. J. W. Jansen. Analysis of variance designs for model output. *Computer Physics Communications*, 117(1–2):35–43, 1999.

[20] I. T. Jolliffe. A note on the use of principal components in regression. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 31(3):300–303, 1982.

[21] P. L'Ecuyer and C. Lemieux. A survey of randomized quasi-Monte Carlo methods. In *Modeling Uncertainty: An Examination of Stochastic Theory, Methods, and Applications*, pages 419–474. Kluwer Academic Publishers, 2002.

[22] J. Matoušek. On the $L^2$–discrepancy for anchored boxes. *Journal of Complexity*, 14(4):527–556, 1998.

[23] B. Moskowitz and R. E. Caflisch. Smoothness and dimension reduction in quasi-Monte Carlo methods. *Mathematical and Computer Modelling*, 23(8-9):37–54, 1996.

[24] H. Niederreiter. *Random Number Generation and Quasi-Monte Carlo Methods*. SIAM, Philadelphia, PA, 1992.

[25] G. Ökten and A. Göncü. Generating low-discrepancy sequences from the normal distribution: Box–Muller or inverse transform? *Mathematical and Computer Modelling*, 53(5-6):1268–1281, 2011.

[26] A. B. Owen. Randomly permuted $(t, m, s)$-nets and $(t, s)$-sequences. In *Monte Carlo and Quasi-Monte Carlo Methods in Scientific Computing*, pages 299–317, New York, 1995. Springer-Verlag.

[27] A. B. Owen. Monte Carlo variance of scrambled net quadrature. *SIAM Journal of Numerical Analysis*, 34(5):1884–1910, 1997.

[28] A. B. Owen. Scrambled net variance for integrals of smooth functions. *Annals of Statistics*, 25(4):1541–1562, 1997.

[29] A. B. Owen. Scrambling Sobol' and Niederreiter-Xing points. *Journal of Complexity*, 14(4):466–489, December 1998.

[30] A. B. Owen. Local antithetic sampling with scrambled nets. *Annals of Statistics*, 36(5):2319–2343, 2008.

[31] A. B. Owen. Monte Carlo Theory, Methods and Examples. statweb.stanford.edu/~owen/mc, 2013.

[32] A. B. Owen and D. Rudolf. A strong law of large numbers for scrambled net integration. *SIAM Review*, 63(2):360–372, 2021.

[33] Z. Pan and A. B. Owen. The nonzero gain coefficients of Sobol's sequences are always powers of two. Technical Report arXiv:2106.10534, Stanford University, 2021.

[34] A. Papageorgiou. The Brownian bridge does not offer a consistent advantage in quasi-Monte Carlo integration. *Journal of Complexity*, 18(1):171–186, 2002.

[35] M. T. Parente, J. Wallin, and B. Wohlmuth. Generalized bounds for active subspaces. *Electronic Journal of Statistics*, 14(1):917–943, 2020.

[36] G. Pirsic. Schnell konvergierende Walshreihen über gruppen. Master's thesis, University of Salzburg, 1995. Institute for Mathematics.

[37] S. Razavi, A. Jakeman, A. Saltelli, C. Prieur, B. Iooss, E. Borgonovo, E. Plischke, S. L. Piano, T. Iwanaga, W. Becker, S. Tarantola, J. H. A. Guillaume, J. Jakeman, H. Gupta, N. Melillo, G. Rabitti, V. Chabirdon, Q. Duan, X. Sun, S. Smith, R. Sheikholeslami, N. Hosseini, M. Asadzade, A. Puy, S. Kucherenko, and H. R. Maier. The future of sensitivity analysis: an essential discipline for systems modeling and policy support. *Environmental Modelling & Software*, 137:104954, 2021.

[38] C. P. Robert and G. O. Roberts. Rao-Blackwellization in the MCMC era. Technical Report arXiv:2101.01011, University of Warwick, 2021.

[39] I. M. Sobol'. The distribution of points in a cube and the accurate evaluation of integrals. *USSR Computational Mathematics and Mathematical Physics*, 7(4):86–112, 1967.

[40] I. M. Sobol'. *Multidimensional Quadrature Formulas and Haar Functions*. Nauka, Moscow, 1969. (In Russian).

[41] I. M. Sobol'. Sensitivity estimates for nonlinear mathematical models. *Mathematical Modeling and Computational Experiment*, 1:407–414, 1993.

[42] I. M. Sobol' and S. Kucherenko. Derivative based global sensitivity measures and their link with global sensitivity indices. *Mathematics and Computers in Simulation (MATCOM)*, 79(10):3009–3017, 2009.

[43] H. F. Trotter and J. W. Tukey. Conditional Monte Carlo for normal samples. In *Symposium on Monte Carlo Methods*, pages 64–79, New York, 1956. Wiley.

[44] Y. Xiao and X. Wang. Conditional quasi-Monte Carlo methods and dimension reduction for option pricing and hedging with discontinuous functions. *Journal of Computational and Applied Mathematics*, 343:289–308, 2018.

[45] Y. Xiao and X. Wang. Enhancing quasi-Monte Carlo simulation by minimizing effective dimension for derivative pricing. *Computational Economics*, 54(1):343–366, 2019.

[46] R.-X. Yue and S.-S. Mao. On the variance of quadrature over scrambled nets and sequences. *Statistics & Probability Letters*, 44(3):267–280, 1999.

[47] C. Zhang, X. Wang, and Z. He. Efficient importance sampling in quasi-Monte Carlo methods for computational finance. *SIAM Journal on Scientific Computing*, 43(1):B1–B29, 2021.