





## MONTE CARLO (MC) METHODS

Consider an expectation of the form:  $\mathbb{E}_{X \sim \mu} \left[ f(X) \right] \,,$ 

where:  $\mu \in \mathscr{P}(\mathbb{R}^d)$  is an isotropic probability distribution,  $f : \mathbb{R}^d \to \mathbb{R}$  is measurable, and  $\mu$ integrable. A **standard MC estimator** is given by:

 $\frac{1}{N} \sum_{i=1}^{N} f(X_i^{\text{iid}}), \text{ where } (X_i^{\text{iid}})_{i=1}^{N} \overset{\text{i.i.d.}}{\sim} \mu.$ 

Applications in Machine Learning:

- 1. dimensionality reduction techniques (JLTs)
- 2. scaling kernel methods via random features
- 3. locality sensitive hashing algorithms (LSH)
- 4. ES methods for Reinforcement Learning
- 5. and many more...

How to improve **computational complexity** and **accuracy** of the above baseline ?

## APPROXIMATE MC METHODS

To improve computational complexity of orthogonal MC estimators, several constructions related to sampling from distributions "close" to Haar distribution on  $\mathcal{O}(d)$  were proposed:

Hadamard-Rademacher Chains:

$$\mathbf{X}_T = \prod_{i=1}^T \mathbf{H} \mathbf{D}_i,$$

with i.i.d. random diagonal matrices  $(\mathbf{D}_i)_{i=1}^T$  with i.i.d. diagonal entries  $Unif(\{\pm 1\})$ , and where  $\mathbf{H} \in \mathbb{R}^{2^L \times 2^L}$  stands for the normalized *Kroneckerproduct* Hadamard matrix defined as:



L times

for  $\otimes$  denoting the Kronecker-product operator.

**Computational Complexity:** Due to Fast Walsh-Hadamard Transform,  $\mathbf{X}_T \mathbf{v}$  for any vector  $\mathbf{v} \in \mathbb{R}^{2^L}$  can be computed in time:  $\mathcal{O}(T \times L \times 2^L)$ .

Many related constructions: e.g. Butterfly matrices generalizing Hadamard-Rademacher chains.

## UNIFYING ORTHOGONAL MONTE CARLO METHODS KRZYSZTOF CHOROMANSKI<sup>1\*</sup>, MARK ROWLAND<sup>2\*</sup>, WENYU CHEN<sup>3</sup>, ADRIAN WELLER<sup>2,4</sup> <sup>1</sup>Google Brain, <sup>2</sup>University of Cambridge, <sup>3</sup>Massachusetts Institute of Technology, <sup>4</sup>The Alan Turing Institute, \*EQUAL CONTRIBUTION

## **ORTHOGONAL MC ESTIMATORS**

Consider the following orthogonal MC estimator:

$$\frac{1}{N}\sum_{i=1}^{N} f(X_i^{\text{ort}}), \text{ where } (X_i^{\text{ort}}) \sim \mu \text{ and } X_i^{\text{ort}} \perp X_j^{\text{ort}}.$$

Note that for  $N \leq d$ :

- if  $\mu$  is isotropic, it can be constructed via Gram-Schmidt orthogonalization,
- the constructions are expensive:  $\mathcal{O}(d^3)$  time
- condition:  $X_i^{\text{ort}} \sim \mu$  implies **unbiasedness**.
- $(X_i^{\text{ort}})_{i=1}^N$  are renormalized rows of a matrix sampled from Haar measure on  $\mathcal{O}(d)$ .

Statistical improvements: It often holds

$$\operatorname{MSE}\left(\frac{1}{N}\sum_{i=1}^{N}f(X_{i}^{\operatorname{ort}})\right) < \operatorname{MSE}\left(\frac{1}{N}\sum_{i=1}^{N}f(X_{i}^{\operatorname{iid}})\right).$$

## **GIVENS ROTATIONS MATRICES**

A *d*-dimensional **Givens rotation** is an orthogonal matrix specified by two distinct indices  $i, j \in$ [d], and an angle  $\theta \in [0, 2\pi)$ . The Givens rotation is then given by the matrix  $\mathbf{G}[i, j, \theta]$  satisfying

if  $k = l \in \{i, j\}$  $\cos(\theta)$ if k = i, l = j $-\sin(\theta)$  $\mathbf{G}[i,j,\theta]_{k,l} = \mathbf{\zeta}$  $\sin( heta)$ if k = j, l = iif  $k = l \notin \{i, j\}$ otherwise.

Givens rotation  $\mathbf{G}[i, j, \theta]$  composed on the right with a reflection in the j coordinate will be termed a **Givens reflection** and written  $G[i, j, \theta]$ . Givens rotations and reflections will be generically referred to as **Givens transformations**.

Lemma 1 (Pillai, Smith 2016) There exists C > 0such that the total variation distance TV between distribution  $\mathcal{D}_{Giv}^d$  on the d-sphere induced by the product of  $Cd \log(d)$  independent Givens random rotations acting on the  $L_2$ -normalized input vector  $\mathbf{x}$  and the distribution  $\mathcal{D}^d_{\mathrm{Haar}}$  related to Haar measure on that *sphere satisfies:* 

 $\lim_{d \to \infty} \mathrm{TV}(\mathcal{D}_{\mathrm{Giv}}^d, \mathcal{D}_{\mathrm{Haar}}^d) = 0.$ 



**Figure 1:** Row 1: the matrix  $\widetilde{\mathbf{F}}^{1,3}$  expressed as a commuting product of Givens reflections. Row 2: As above, but for matrix  $\widetilde{\mathbf{F}}^{2,3}$ . Row 3: the matrix  $\widetilde{\mathbf{F}}^{3,3}$  expressed as a product of commuting Givens rotations. Row 4: the normalised Hadamard matrix  $H_3$  written as a product of  $\tilde{\mathbf{F}}^{1,3}$ ,  $\tilde{\mathbf{F}}^{2,3}$  and  $\tilde{\mathbf{F}}^{3,3}$ . White/black represent 0/1 elements and grey/blue - elements in (0, 1) and (-1, 0).



Figure 2: Comparison of different Monte Carlo methods on the task of Gaussian kernel approximation.

## ON THE HUNT FOR UNIFYING APPROXIMATE ORTHOGONAL MCS



**Kac's random walk matrices** (random  $I_t, J_t, \theta_t$ ):  $\mathbf{K}_T = \prod \mathbf{G}[I_t, J_t, \theta_t],$ 

The Hadamard-Rademacher Chains: Each block

where:  $\widetilde{\mathbf{F}}^{j,L}$ 

### Hadamard-MultiRademacher matrices:

# **EXPERIMENTS: FROM KERNEL APPROXIMATION TO RL**





Figure 3: Comparison of different Monte Carlo methods for gradient estimation in ES algorithms for RL.



# The Alan Turing Institute

 $HD_i$  of the chain can be rewritten as:

$$\mathbf{H}\mathbf{D}_t = \left(\prod_{i=1}^{L-1} \widetilde{\mathbf{F}}^{i,L}\right) \left(\widetilde{\mathbf{F}}^{L,L}\mathbf{D}_t\right).$$

$$= \prod_{\substack{\lambda \in \mathbb{F}_2^L \\ \lambda_j = 0}} \widetilde{\mathbf{G}}[\boldsymbol{\lambda}, \boldsymbol{\lambda} + \mathbf{e}_j, \pi/4] \in \mathscr{O}(2^L) \,.$$

for the canonical basis  $\mathbf{e}_1, \ldots, \mathbf{e}_L$  of  $\mathbb{F}_2^L$ .

In this expression, we may interpret  $\widetilde{\mathbf{F}}^{L,L}\mathbf{D}_t$  as a product of random Givens transformations with a deterministic, structured choice of rotation axes.

$$\prod_{i=1}^{L} \left( \widetilde{\mathbf{F}}^{i,L} \mathbf{D}_i \right)$$

Denote the kernel estimator using Kac's random walk matrices with k Givens rotations as:  $\widehat{K}_{kac}^k$ and the unstructured baseline as  $\widehat{K}_{\text{base}}$ . We have:

### Theorem 1 (Kac's random walk for RBF kernels) Let $K_d : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ be a Gaussian kernel. Then there exists a constant C > 0 such that for $\mathbf{x}, \mathbf{y}$ , $k = C \cdot d \log d$ and d large enough we have:

