

The Unreasonable Effectiveness of Structured Random Orthogonal Embeddings

Krzysztof Choromanski, Mark Rowland, Adrian Weller

Why orthogonal random features ?

➤ In practice they are very effective

- o applied successfully in many ML settings:
LSH, kernel ridge regression, RNNs and more...

➤ Better time and space complexity

- o certain discrete variants provide computational speedups and lower space complexity even though it was **not known** whether they offer also accuracy improvements

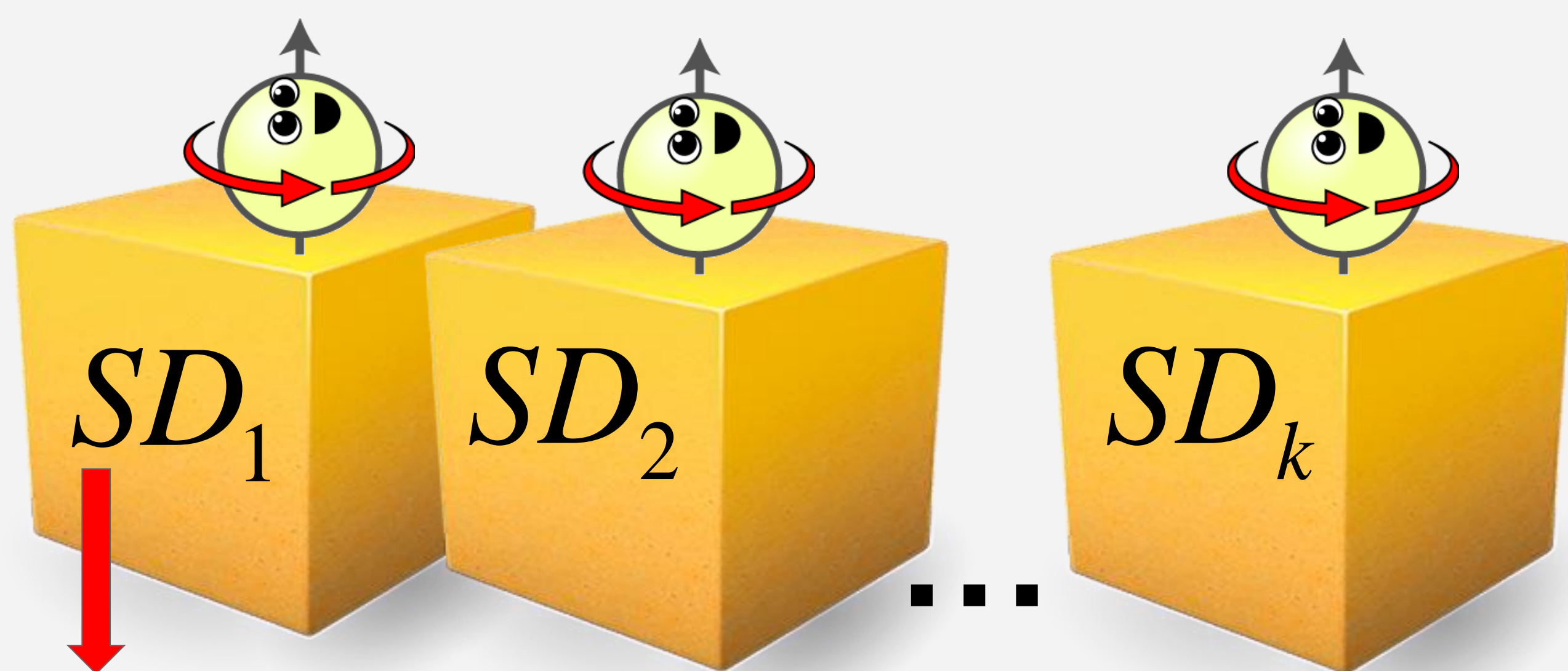
➤ Improving accuracy over state-of-the-art

- o improved accuracy proved previously **only** for the Gaussian kernel and **only** asymptotically for large n



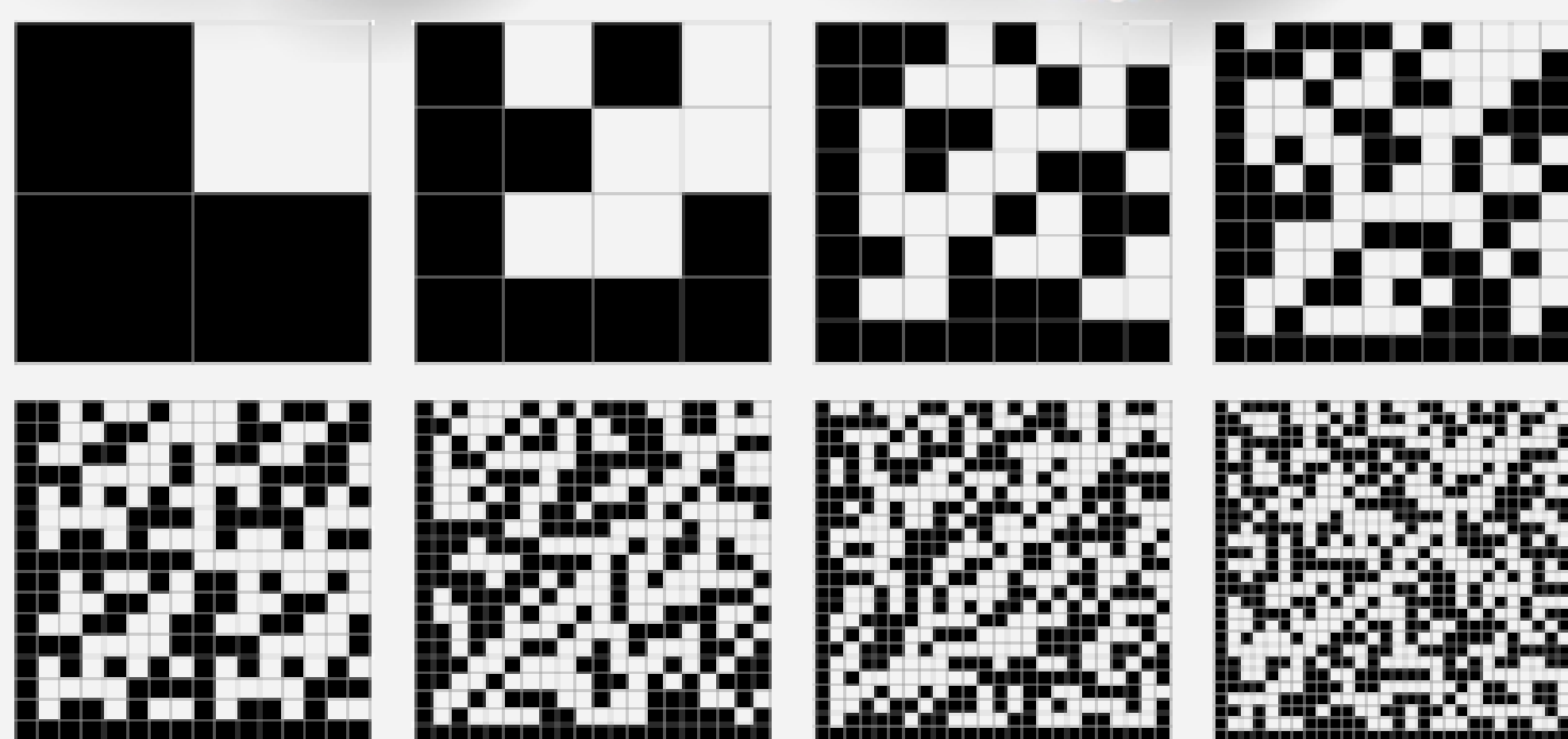
ROM Matrices

Discrete constructions



$$s_{i,j} \in \left\{ -\frac{1}{\sqrt{n}}, +\frac{1}{\sqrt{n}} \right\}$$

- orthogonal rows
- many examples:
Kronecker-product matrices, Walsh-Hadamard matrices, quadratic residue constructions



$$\mathbf{D} = \begin{pmatrix} \lambda_1 & 0 & 0 & \cdots \\ 0 & \lambda_2 & 0 & \cdots \\ 0 & 0 & \lambda_3 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$$

$$\lambda_i \sim \text{Unif}\{-1, +1\}$$

$$\lambda_i \sim \text{Unif}\{-1, +1, -i, +i\}$$

$$\lambda_i \sim \text{Unif}(S^1) \subset \mathbb{C} \rightarrow \text{random complex matrices}$$

Our theoretical results

Free-lunch JLTs

Estimators:

$$\hat{K}_m^{\text{base}}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} (\mathbf{G}\mathbf{x})^\top (\mathbf{G}\mathbf{y}) \rightarrow \text{previous state-of-the-art}$$

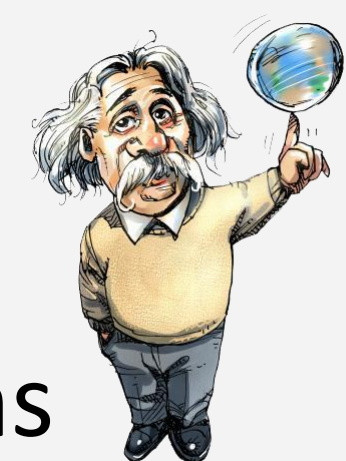
$$\hat{K}_m^{\text{ort}}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} (\mathbf{G}_{\text{ort}}\mathbf{x})^\top (\mathbf{G}_{\text{ort}}\mathbf{y}) \rightarrow \text{our estimators}$$

$$\hat{K}_m^{(k)}(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \left(\mathbf{M}_{\text{SR}}^{(k), \text{sub}} \mathbf{x} \right)^\top \left(\mathbf{M}_{\text{SR}}^{(k), \text{sub}} \mathbf{y} \right)$$

$$\mathbf{M}_{\text{SR}}^{(k)} = \prod_{i=1}^k \text{SD}_i(\mathcal{R}) \rightarrow |\lambda_i| = 1$$

➤ sub: different subsampling strategies for row selection to reduce dimensionality

- o first/last m rows
- o sample uniformly at random with repetitions
- o sample uniformly at random without repetitions



THEOREM

$$\text{MSE}(\hat{K}_m^{(k)}(\mathbf{x}, \mathbf{y})) = \frac{1}{m} \left(\frac{n-m}{n-1} \right) \left(((\mathbf{x}^\top \mathbf{y})^2 + \|\mathbf{x}\|^2 \|\mathbf{y}\|^2) + \sum_{r=1}^{k-1} \frac{(-1)^r 2^r}{n^r} (2(\mathbf{x}^\top \mathbf{y})^2 + \|\mathbf{x}\|^2 \|\mathbf{y}\|^2) + \frac{(-1)^k 2^k}{n^{k-1}} \sum_{i=1}^n x_i^2 y_i^2 \right).$$

➤ Implies that discrete ROMs provide more accurate JLT mechanisms than state-of-the-art, with better time and space complexity

➤ x2 smaller MSEs for complex discrete ROMs

➤ strictly better accuracy also for the angular kernel approximation with Gaussian orthogonal matrices

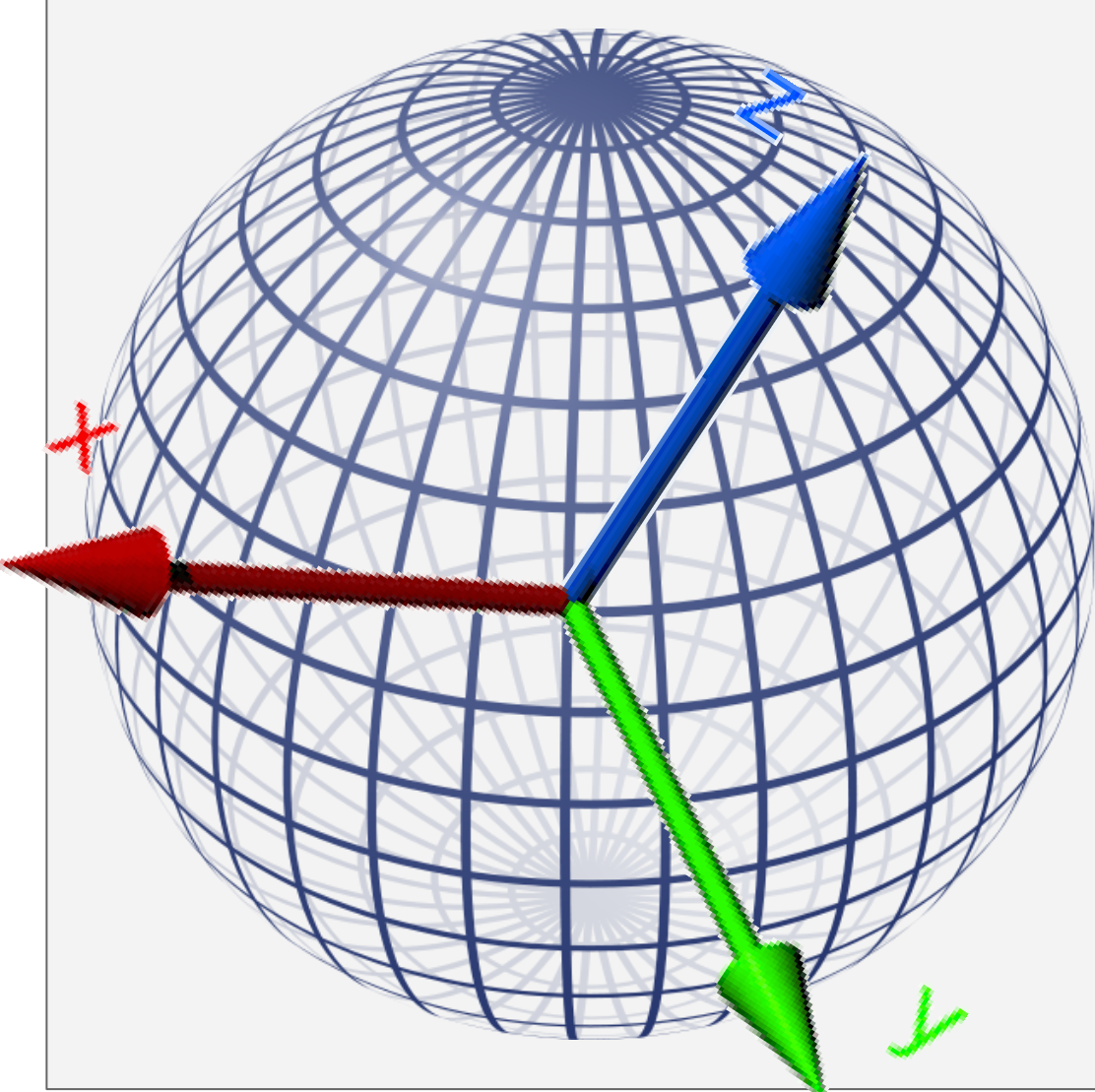
$$K^{\text{ang}}(\mathbf{x}, \mathbf{y}) = 1 - \frac{2\theta_{\mathbf{x}, \mathbf{y}}}{\pi} \rightarrow \text{kernel definition}$$

$$\hat{K}_m(\mathbf{x}, \mathbf{y}) = \frac{1}{m} \text{sgn}(\mathbf{M}\mathbf{x})^\top \text{sgn}(\mathbf{M}\mathbf{y}) \rightarrow \text{general estimator}$$

THEOREM

$$\text{MSE}(\hat{K}_m^{\text{ang,ort}}(\mathbf{x}, \mathbf{y})) < \text{MSE}(\hat{K}_m^{\text{ang,base}}(\mathbf{x}, \mathbf{y}))$$

G_{ort} Continuous constructions



➤ Gaussian orthogonal matrices

- o correspond to truly random rotations in n-dimensional spaces
- o can be easily generated from unstructured Gaussian matrices (with one-time extra cost) via the Gram-Schmidt orthogonalization

Experiments

