

Random Features for Kernel Approximation: A Survey

Blanca Cano Camarero

Universidad Autónoma de Madrid

April 25, 2023



Overview

- 1 Last week queries
- 2 Article information
- 3 Theoretical Foundation of Random Features
- 4 Taxonomy of random features based algorithms
- 5 Theoretical Analysis
- 6 Data dependent algorithm
- 7 Computational cost

On the Error of Random Fourier Features

- ▶ Authors: Danica J. Sutherland and Jeff Schneider
- ▶ Year 2015.
- ▶ Cited by 185.
- ▶ See Sutherland and Schneider (2015)

Why is interesting

- ▶ Which of two embedding is interesting.
- ▶ Error considerations.

Two embedding based on Fourier transform

One of the form:

$$\tilde{z}(x) = \sqrt{\frac{2}{D}} [\sin(w_1^T x) \cos(w_1^T x) \dots \sin(w_{D/2}^T x) \cos(w_{D/2}^T x)]^T, w_i \sim P \quad (1)$$

and another of the form

$$\check{z}(x) = \sqrt{\frac{2}{D}} [\cos(w_1^T x + b_1) \dots \cos(w_D^T x + b_D)]^T, w_i \sim P, b_i \sim \text{Unif}([0, 2\pi]) \quad (2)$$

Kernel approximation

$$\tilde{s}(x, y) = \frac{1}{D/2} \sum_{i=1}^{D/2} \cos(w_i^T (x - y)) \quad (3)$$

$$\check{s}(x, y) = \frac{1}{D} \sum_{i=1}^D \cos(w_i^T (x - y)) + \cos(w_i^T (x + b) + 2b_i) \quad (4)$$

Models comparatives

- ▶ of the first 100 papers cit- ing Rahimi and Recht (2007) in a Google Scholar search, used either \tilde{z} or the equivalent formulation.
- ▶ The article show that \tilde{z} is superior for Gaussian kernel. Providing a complementary view of the quality of these embeddings:
 - ▶ Variance of each embedding \tilde{s} for RBF uniformly lower variance.
 - ▶ Uniform convergence bounds, tightening constants in the original \tilde{z} (one worse).
 - ▶ L_2 convergence of each approximation \tilde{z} superior again.
 - ▶ Finally empirically.

Random Features for Kernel Approximation: A Survey on Algorithms, Theory, and Beyond

- ▶ Fanghui Liu and Xiaolin Huang and Yudong Chen and Johan A. K. Suykens
- ▶ Year 2021
- ▶ Cited by 94!!!

Liu et al. (2021)

Theoretical Foundation of Random Features

Theorem

Bochner's Theorem A continuous and shift-invariant function $k : \mathbb{R}^d \times \mathbb{R}^d \longrightarrow \mathbb{R}$ is positive definite if and only if it can be represented as

$$k(x - x') = \int_{\mathbb{R}^d} \exp iw^T(x - x') \mu_k(dw) \quad (5)$$

where μ_k is a positive finite measure on the frequencies w .

According to Bochner's theorem, the spectral distribution μ_k of a stationary kernel k is the finite measure induced by a Fourier transform.

By setting $k(0) = 1$, we may normalize μ_k to a probability density p (the Fourier transform associated with k), hence

$$k(x - x') = \int_{\mathbb{R}^d} \exp iw^T(x - x')\mu(dw) \quad (6)$$

$$= E_{w \sim p} \left[\exp(iw^T x) \exp(iw^T x')^*, \right] \quad (7)$$

The kernels used in practice are typically real-valued and thus the imaginary part can be discarded.

$$\varphi_p(x) = \left[\exp(-i w_1^T x), \dots, \exp(-i w_s^T x) \right]^T. \quad (8)$$

Two-layer network correspondence to random features approximation

Given a two-layer network

$$f(x; \theta) = \sqrt{\frac{2}{s}} \sum_{j=1}^s a_j \sigma(w_j^T x)$$

for some activation function σ and $x \in \mathbb{R}^d$ when $w \sim \mathcal{N}(0, I_d)$ and only the second layer are optimized this correspond to random features approximation

$$k(x, x') = E_{w \sim \mathcal{N}(0, I_d)} \left[\sigma(w^T x) \sigma(w^T x'), \right] \quad (9)$$

depends on the kernel type such that

$$\varphi(x) := \sigma(Wx). \quad (10)$$

Some examples:

- ▶ Correspond to a standard RFF model for a Gaussian kernel:

$$\sigma(x) = [\cos(x) \sin(x)] \quad (11)$$

- ▶ First order arc-cosine kernel termed as

$$k(x, x') = k(u) = \frac{1}{\pi}(u(\pi - \arccos(u)) + \sqrt{1 - u^2}) \quad (12)$$

where

$$u = \frac{\langle x, x' \rangle}{\|x\| \|x'\|} \quad (13)$$

if $\sigma(x) = \max\{0, x\}$.

For fully connected deep neural network

Pregunta: Una red neuronal de dos capas representa un kernel. Ese kernel se puede aproximar a una red neuronal de una capa oculta. LUEGO SE ESTARÍA TRANSFORMANDO LAS REDES NEURALES DE UNA CAPA A OTRA.

Referencias que ver Daniely et al. (2017) and Lee et al. (2018) and from the origin article references (12, 63,64,65,66,67 and 68).

Commonly used kernels in Random Features

Actually not interesting

Taxonomy of random features based algorithms

$$\varphi(x) = \left[a_1 \exp(-i w_1^T x), \dots, a_s \exp(-i w_s^T x), \right]^T \quad (14)$$

They can be grouped into two categories, data-independent algorithms and data-dependent algorithms, based on whether or not the selection of w_i and a_i is independent of the training data

See figure 2 of the article.

And Rahimi and Recht (2008) is improved by Le et al. (2014)

Error approaches

- ▶ **Approximation:** how many random features are needed to ensure a high quality estimator in kernel approximation?
- ▶ **Generalization** how many random features are needed to incur no loss of empirical risk and expected risk in a learning estimator?
 - ▶ See convergences on table 3 of Liu et al. (2021).
 - ▶ Rates and required random features in table 4.
 - ▶ Comparison of learning rates and required random features for expected risk with a Lipschitz continuous loss function. Table 5.

Other comparison and interesting fields

- ▶ Space and time complexities see table 2 page 8 of Liu et al. (2021).
- ▶ Trends: High-dimensional Random Features in Over parameterized settings (page 21)
 - ▶ Double descent in the test error curve (as we saw in DNN)

Future

RFF in DNNs!!! See references Figure 1 Liu et al. (2021).

Next week: Neural Network?

1. **Nystroem Method vs Random Fourier Features: A Theoretical and Empirical Comparison, Advances in Neural Information Processing Systems 2012**
2. **Random features for kernel approximation: A survey on algorithms, theory, and beyond**
3. **Williams, C.K.I. and Seeger, M. "Using the Nystroem method to speed up kernel machines", Advances in neural information processing systems 2001**
T. Yang, Y. Li, M. Mahdavi, R. Jin and Z. Zhou
4. **Weighted Sums of Random Kitchen Sinks: Replacing minimization with randomization in learning.**
5. **Randomness in neural networks: an overview**
6. **Fast and scalable polynomial kernels via explicit feature maps**
7. **On the error of random Fourier features**
8. **A survey on large-scale machine learning**
9. **Sharp analysis of low-rank kernel matrix approximations**

Previous meets

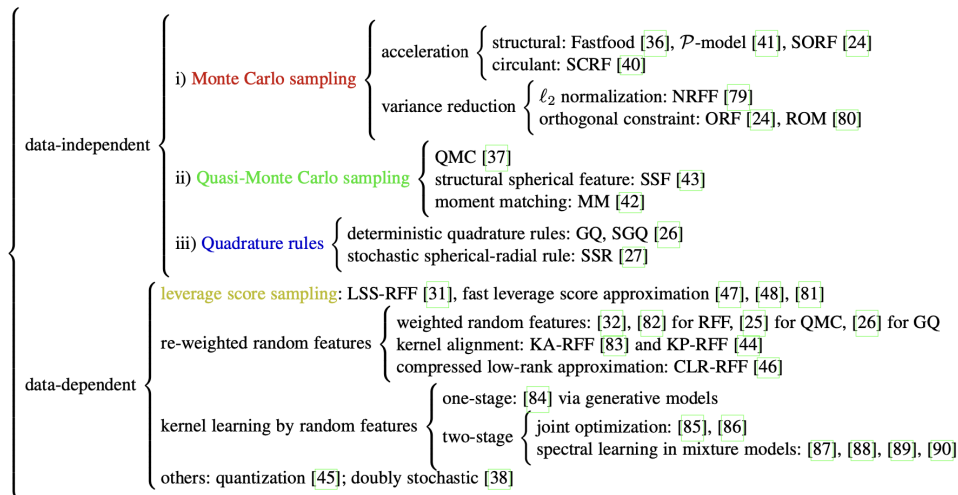
1. Nyström vs Random Fourier Features.
 - ▶ When eigenvalues are sparse Nyström is better.
 - ▶ Article hypothesis: training data dependence.
2. Random Features for kernel approximation survey.
 - ▶ Random features general approach.
 - ▶ Survey of plenty techniques: **data-independent and data-dependent**.

Questions

- ▶ Why are there no comparisons articles?
- ▶ The hypothesis was wrong?
- ▶ Why there are no comparative between algorithm and spectral sparsity.
- ▶ How relevant is spectral sparsity?
- ▶ If it is important why not a lineal combination models and adjust the coefficient by some spectral approximation or heuristic?

This survey would indirectly answer some of the questions

Taxonomy



Space and time complexities: No clear difference

Table 2
Comparison of different kernel approximation methods on space and time complexities to obtain Wx .

| Method | Kernels (in theory) | Extra Memory | Time | Lower variance than RFF |
|--|--------------------------------------|-------------------|-------------------------|-------------------------|
| Random Fourier Features (RFF) [9] | shift-invariant kernels | $\mathcal{O}(sd)$ | $\mathcal{O}(sd)$ | - |
| Quasi-Monte Carlo (QMC) [37] | shift-invariant kernels | $\mathcal{O}(sd)$ | $\mathcal{O}(sd)$ | Yes |
| Normalized RFF (NRFF) [79] | Gaussian kernel | $\mathcal{O}(sd)$ | $\mathcal{O}(sd)$ | Yes |
| Moment matching (MM) [42] | shift-invariant kernels | $\mathcal{O}(sd)$ | $\mathcal{O}(sd)$ | Yes |
| Orthogonal Random Feature (ORF) [24] | Gaussian kernel | $\mathcal{O}(sd)$ | $\mathcal{O}(sd)$ | Yes |
| Fastfood [36] | Gaussian kernel | $\mathcal{O}(s)$ | $\mathcal{O}(s \log d)$ | No |
| Spherical Structured Features (SSF) [43] | shift and rotation-invariant kernels | $\mathcal{O}(s)$ | $\mathcal{O}(s \log d)$ | Yes |
| Structured ORF (SORF) [24], [91] | shift and rotation-invariant kernels | $\mathcal{O}(s)$ | $\mathcal{O}(s \log d)$ | Unknown |
| Signed Circulant (SCRF) [40] | shift-invariant kernels | $\mathcal{O}(s)$ | $\mathcal{O}(s \log d)$ | The same |
| \mathcal{P} -model [41] | shift and rotation-invariant kernels | $\mathcal{O}(s)$ | $\mathcal{O}(s \log d)$ | No |
| Random Orthogonal Embeddings (ROM) [80] | rotation-invariant kernels | $\mathcal{O}(d)$ | $\mathcal{O}(d \log d)$ | Yes |
| Gaussian Quadrature (GQ), Sparse Grids Quadrature (SGQ) [26] | shift invariant kernels | $\mathcal{O}(d)$ | $\mathcal{O}(d \log d)$ | Yes |
| Stochastic Spherical-Radial rules (SSR) [27] | shift and rotation-invariant kernels | $\mathcal{O}(d)$ | $\mathcal{O}(d \log d)$ | Yes |

Models:

Here we test various random features based algorithms:

| Data-independent | Data-dependent |
|---|----------------|
| RFF ORF SORF ROM Fastfood QMC SSF GQ | LS-RFF |

Biased tests!

Datasets:

Table 6
Dataset statistics.

| datasets | d | #traing | #test | random split | scaling |
|-----------------|------|-----------|---------|--------------|---------|
| <i>ijcnn1</i> | 22 | 49,990 | 91,701 | no | - |
| <i>EEG</i> | 14 | 7,490 | 7,490 | yes | mapstd |
| <i>cod-RNA</i> | 8 | 59,535 | 157,413 | no | mapstd |
| <i>covtype</i> | 54 | 290,506 | 290,506 | yes | minmax |
| <i>magic04</i> | 10 | 9,510 | 9,510 | yes | minmax |
| <i>letter</i> | 16 | 12,000 | 6,000 | no | minmax |
| <i>skin</i> | 3 | 122,529 | 122,529 | yes | minmax |
| <i>a8a</i> | 123 | 22,696 | 9,865 | no | - |
| <i>MNIST</i> | 784 | 60,000 | 10,000 | no | minmax |
| <i>CIFAR-10</i> | 3072 | 50,000 | 10,000 | no | - |
| <i>MNIST-8M</i> | 784 | 8,100,000 | 10,000 | no | - |

Results:

| datasets | approximation | | lr | | liblinear | |
|----------------|---------------|--------------------|-----------|-----------|-----------|-----------|
| | small s | large s | small s | large s | small s | large s |
| <i>ijcnn1</i> | SSF | SORF, QMC, ORF | - | - | Fastfood | - |
| <i>EEG</i> | SSF | ORF | - | - | - | - |
| <i>cod-RNA</i> | SSF | - | - | - | - | - |
| <i>covtype</i> | ORF | - | - | - | - | - |
| <i>magic04</i> | SSF | SSF, ORF, QMC, ROM | - | - | - | - |
| <i>letter</i> | SSF | SSF, ORF | - | - | - | - |
| <i>skin</i> | SSF, ROM | QMC | - | - | - | - |
| <i>a8a</i> | - | - | - | - | SSF | - |

- Approximation: small (SSF: 8/10), large (ORF 4/10, QMC: 3/10, SSF: 2/10).
- Lineal regression: Adjusting the model kernel error are alleviated.

Having in mind taxonomy

- ▶ SSF: data independent, Quasi-Monte Carlo sampling.
- ▶ ORF: data independent, Monte Carlo sampling.
- ▶ QMC: data independent, Quasi-Monte Carlo sampling.

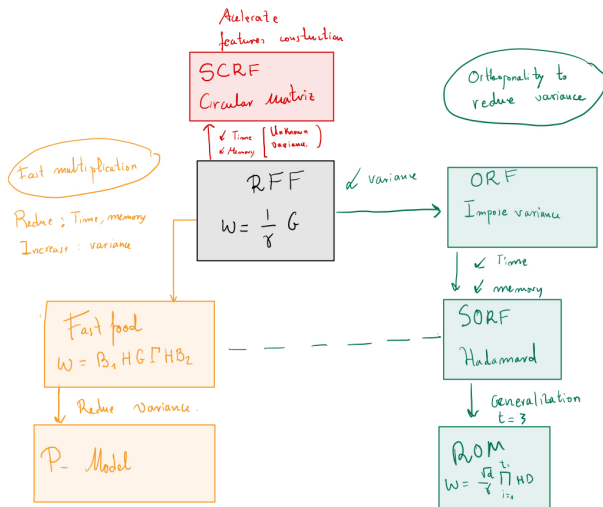
Where is LS-RFF? (See appendix B)

- ▶ Highest (except ijcn1) approximation error.
- ▶ Most expensive time cost.

Why they did not test data indepenence techniques?

- Not worthy? (A bit weird)

Monte Carlo sampling based approaches



Hadamard matrix (wikipedia)

The Hadamard matrix can be constructed using the Sylvester construction, which involves recursively combining smaller Hadamard matrices to form larger ones.

Start with the 1×1 Hadamard matrix $H_1 = [1]$. To construct the $2n \times 2n$ Hadamard matrix H_{2n} , form the matrix

$$H_{2n} = \begin{bmatrix} H_n & H_n \\ H_n & -H_n \end{bmatrix}.$$

Using this construction, we can generate Hadamard matrices of any power of 2 size. For example, we can construct the 4×4 Hadamard matrix as follows:

$$H_4 = \begin{bmatrix} H_2 & H_2 \\ H_2 & -H_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}.$$

Walsh-Hadamard matrix

Rearrange the rows of the matrix according to the number of sign change of each row.

$$W_4 = \frac{1}{\sqrt{4}} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix}$$

Note that the Walsh matrix of size 4 is equivalent to the second-order Hadamard matrix that I gave earlier.

Quasi-Monte Carlo Sampling

Hypothesis: p factorizes with respect to the dimension

$$p(x) = \prod_{j=1}^d p_j(x_j)$$

QMC transform the integral to

$$k(x - x') = \int_{[0,1]^d} \exp(i(x - x')^T \phi^{-1}(t)) dt,$$

where $\phi^{-1}(t) = (\phi_1^{-1}(t_1), \dots, \phi_d^{-1}(t_d))$ where $\{t\}$ is a low discrepancy sequence.

$$W = [\phi_1^{-1}(t_1), \dots, \phi_d^{-1}(t_s)]$$

SSF: Spherical structured feature maps for kernel approximation

It improves the space and time complexities of QMC for approximating shift- and rotation-invariant kernels

SSF generates points

$$\{v_1, v_2, \dots, v_s\}$$

asymptotically uniformly distributed on the sphere S_{d-1} and constructs the transformation matrix as

$$W_{\text{SSF}} = [\Phi^{-1}(t)v_1, \Phi^{-1}(t)v_2, \dots, \Phi^{-1}(t)v_s]^T \in \mathbb{R}^{s \times d},$$

where $\Phi^{-1}(t)$ is the inverse of the cumulative distribution function of the standard normal distribution evaluated at t .

$$V = \frac{1}{\sqrt{d/s}} \begin{bmatrix} \Re F_{\Lambda} & -\Im F_{\Lambda} \\ \Im F_{\Lambda} & \Re F_{\Lambda} \end{bmatrix}$$

where F_{Λ} consists of the row of the Fourier discrete matrix.

The discrete Fourier transform (DFT) matrix is defined as:

$$F_N = \frac{1}{\sqrt{N}} \begin{pmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & w_N & w_N^2 & \cdots & w_N^{N-1} \\ 1 & w_N^2 & w_N^4 & \cdots & w_N^{2(N-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & w_N^{N-1} & w_N^{2(N-1)} & \cdots & w_N^{(N-1)^2} \end{pmatrix},$$

where $w_N = e^{-2\pi i/N}$ is a primitive N -th root of unity.

The matrix F_N is a unitary matrix and its inverse is given by the Hermitian transpose, i.e., $F_N^{-1} = F_N^H$.

$$F = \frac{1}{2} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & i & -1 & -i \\ 1 & -1 & 1 & -1 \\ 1 & -i & -1 & i \end{pmatrix}$$

The selection of d^2 rows from F is done by minimizing the discrete Riesz 0-energy such that the points spread as evenly as possible on the sphere.

Leverage score based sampling: main idea

Generate samples w from a distribution q .

Feature mapping:

$$\varphi_q(x) = \frac{1}{\sqrt{s}} \left(\sqrt{\frac{p(w_1)}{q(w_1)}} e^{-iw_1^T x}, \dots, \sqrt{\frac{p(w_s)}{q(w_s)}} e^{-iw_s^T x} \right)^T \quad (15)$$

How to design q

Let define ridge leverage function:

$$l_\lambda(w_i) = p(w_i) z_{p,w_i}^T(X) (K + n\lambda I)^{-1} z_{p,w_i}(X), \quad (16)$$

where λ is the KRR regularization parameter.

Define

$$d_K^\lambda = \int_{\mathbb{R}^d} l_\lambda(w) dw = \text{tr}[K(K + n\lambda I)^{-1}]. \quad (17)$$

Number of effective degrees of freedom (independent parameters in a learning problem)

$$d^\lambda \ll n$$

q is designed as

$$q(w) = \frac{l_\lambda(w)}{d_K^\lambda}. \quad (18)$$

Improvements

- ▶ Cost K approximation and inverse ($O(n^3)$)
- ▶ LS-RFF uses a subset of data to approximate K . ($O(ns^2 + s^3)$).
- ▶ SLS-RFF (Surrogate Leverage Score-RFF) to avoid inverting an $s \times s$ matrix. $O(ns^2)$ same complexity as RFF.

Re-weighted random features

- ▶ KA-RFF (Kernel Alignment-RFF) : It pre-computes a large number of random features that are generated by RFF, and then select a subset of them by solving a simple optimization problem based on kernel alignment. In particular, the optimization problem is

$$\max_{a \in P_J} \sum_n^{i,j=1} y_i y_j \sum_{t=1}^J a_t z_p(x_i, w_t) z_p(x_j, w_t). \quad (19)$$

- ▶ KP-RFF (Kernel Polarization-RFF) and CLR-RFF (Compression Low Rank-RFF): It first generates a large number of random features by RFF and then selects a subset from them using an energy-based scheme or solving and optimization problem.

Hilbert cores

The problem of building a small, weighted subset of the data that approximates the full dataset, is known as the **Hilbert cores**. One solution:

- ▶ Greedy iterative geodesic ascent,
- ▶ Frank-Wolfe based methods,
- ▶ Johnson-Lindenstrauss random projection.

Kernel learning by random features

Construct random features using sophisticated learning techniques, e.g., by learning the spectral distribution of kernel from the data.

- ▶ Representative approaches in this class often involve a one- stage or two-stage process.
 1. Learns the random features,
 2. and then incorporates them into kernel methods for prediction.
- ▶ Example: leverage sampling and random features selection.
- ▶ One-stage algorithms aim to simultaneously learn the spectral distribution of a kernel and the prediction model by solving a single joint optimization problem or using a spectral inference scheme.

computational cost

Random Fourier Features Li (2021): reducing the computational complexity (roughly from $O(n^3)$ in time and $O(n^2)$ in space to $O(n^2)$ and $O(n\sqrt{n})$ respectively) without having to trade-off the statistical prediction accuracy.

| Step | Task | Theory | Cost | Memory |
|------|--|---|--------|---------|
| 1 | Sampling Fourier components u_1, \dots, u_m from $p(u)$. | Inverse transform sampling, Accept or reject, Montecarlo. | $O(1)$ | $O(m)$ |
| 2 | Compute $z_f(x) = (\sin(u_1^T x), \cos(u_1^T x), \dots, \sin(u_m^T x), \cos(u_m^T x))$ | | $O(m)$ | $O(2m)$ |

Nyström

| Step | Task | Theory | Cost | Memory |
|------|---|-------------------------------|----------|----------|
| 1 | Sampling | SVD and matrix multiplication | $O(m)$ | $O(m)$ |
| 2 | constructing a low rank matrix by $\hat{K}_r = K_b \hat{K}_b^\dagger K_b^T$. | | $O(n^3)$ | $O(m^2)$ |

| Training model | Cost | Memory |
|---------------------------|--|-------------|
| Ridge regression | $O(m^2(N + m))$ | $O(mn + n)$ |
| KRR | $O(n^3)$ | $O(n^2)$ |
| SVM | $O(n^2)$ or $O(n^3)$ | $O(n^2)$ |
| Backpropagation one layer | $O(nod)$ (o output, d: mini batch size) | $O(no)$ |

References I

- Amit Daniely, Roy Frostig, and Yoram Singer. Toward deeper understanding of neural networks: The power of initialization and a dual view on expressivity, 2017.
- Quoc Viet Le, Tamas Sarlos, and Alexander Johannes Smola. Fastfood: Approximate kernel expansions in loglinear time, 2014.
- Jaehoon Lee, Yasaman Bahri, Roman Novak, Samuel S. Schoenholz, Jeffrey Pennington, and Jascha Sohl-Dickstein. Deep neural networks as gaussian processes, 2018.
- Zhu Li. Sharp analysis of random fourier features in classification, 2021.
- Fanghui Liu, Xiaolin Huang, Yudong Chen, and Johan A. K. Suykens. Random features for kernel approximation: A survey on algorithms, theory, and beyond, 2021.

References II

- Ali Rahimi and Benjamin Recht. Weighted sums of random kitchen sinks: Replacing minimization with randomization in learning. In D. Koller, D. Schuurmans, Y. Bengio, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 21. Curran Associates, Inc., 2008. URL https://proceedings.neurips.cc/paper_files/paper/2008/file/0efe32849d230d7f53049ddc4a4b0c60-Paper.pdf.
- Danica J. Sutherland and Jeff Schneider. On the error of random fourier features, 2015.