
On the Importance of High-Rank Kernel Approximation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We investigate the problem of training kernel approximation models under memory
2 constraints. Our first step toward this goal is to compare the performance of
3 two leading kernel approximation methods, the Nyström method [3] and random
4 Fourier features (RFFs) [2], in this memory constrained setting. Although the
5 Nyström method is generally known to outperform RFFs for a fixed number of
6 features [4], we find that RFFs are able to outperform Nystrom features in this
7 constrained setting. This holds true in spite of the fact that the Nyström method
8 often has smaller kernel approximation error than the RFF approximation, in terms
9 of both the spectral norm and the Frobenius norm, under the same budget. In
10 order to explain this phenomenon, we leverage the recent work of Avron et al.[1];
11 in this work, they consider an approximation \tilde{K} to a kernel matrix K , which
12 satisfies $(1 - \Delta)(K + \lambda \mathbb{1}) \preceq \tilde{K} + \lambda \mathbb{1} \preceq (1 + \Delta)(K + \lambda \mathbb{1})$, and they bound
13 the generalization performance of the regression model trained using \tilde{K} , with
14 regularization constant λ , in terms of Δ . We demonstrate across a number of tasks
15 that this metric Δ correlates much more strongly with generalization performance
16 than other measures of kernel approximation error, and show that under a fixed
17 memory budget, random Fourier features systematically attains lower Δ values.
18 We explain that this is due to the fact that the RFF approximation is higher rank
19 than the Nyström approximation, when under the same memory budget. Taking
20 inspiration from these results, we propose using *low-precision* random Fourier
21 features (LP-RFFs) in order to generate an even higher rank approximation under
22 the same memory budget. We build on the theory of Avron et al. to analyze
23 the generalization performance of LP-RFFs, proving that performance decays as
24 **...make formal statement**. Empirically, we show the we are able to get within **xx%**
25 of the performance of full-precision RFF models, using **xx** times less memory
26 during training, on the TIMIT speech recognition dataset.

27 1 Introduction

- 28 • Kernel methods are powerful, but don't scale well. Need kernel approximation methods.
- 29 • For kernel approximation methods to do well, generally need lots of features. When
30 performing GPU training (limited memory), or when considering model deployment, the
31 memory requirement for these models becomes a large bottleneck. An example of this is the
32 recent work in large-scale kernel methods for speech recognition, which are constrained in
33 terms of the number of features they can fit on GPU.
- 34 • As a result, we perform a thorough empirical comparison of Nystrom and RFF, two lead-
35 ing kernel approximation methods, in the memory constrained setting. Surprisingly, we

find that although the Nystrom method often has lower kernel approximation error (Frobenius/Spectral norm), it performs worse than RFF in this setting. We note that this runs counter to the common belief that Nystrom is better than RFF.

- We explain these observations using the recent theoretical results of Avron et al. We note that the bounds in this paper differ in important ways from prior bounds which are generally based on the kernel approximation error (and thus useless to explain the observed phenomenon).
- We note that the primary insight provided by these recent bounds is that it is crucial to have a high-rank decomposition whose spectrum roughly aligns with the exact spectrum, even if it has high kernel approximation error. Note that this runs contrary to a lot of the literature kernel approximation literature, which is often narrowly concerned with reducing the kernel approximation error. Furthermore, we consider the computational cost, in terms of memory, of computing these features, another aspect which is generally not considered when comparing Nystrom and RFF.
- We leverage this insight to propose LP-RFF.
- We analyze LP-RFF theoretically, showing that you can upper bound the value Δ for LP-RFF in terms of the number of bits per feature used.
- We show that LP-RFF are able to perform on par with full-precision features, at a fraction of the memory.
- This opens the door for training kernel approximation models with many more features than are currently being used. To fully leverage these insights to attain faster, lower-memory, and lower-energy training, it will be necessary to implement efficient version of our algorithm (mention TPUs, GPUs, FPGAs, etc). **AM: Are TPUs already being used for neural net training in low-precision?**
- Include sexy plot in intro summarizing the main results of the paper? Generalization performance vs. memory usage on TIMIT, for Nyström, RFF, and LP-RFF?

2 Related Work

- Give an overview of kernel approximation generalization bounds. Note that the central object they generally use to bound the generalization error is $\|K - \tilde{K}\|$. Discuss the more recent results which approach the problem differently **AM: We should read recent papers closely.**
- Discuss Nystrom vs. RFF paper. Our main differentiators: 1) we consider larger/more challenging datasets, 2) we consider much larger numbers of features, 3) we consider memory budget.
- Discuss low-precision, and recent interest in doing low-precision training.
- Discuss fastfood/circulant RFF, and other attempts at reducing memory footprint of projection matrix. Note that this doesn't address the memory of the features themselves, which can take a ton of space for moderate size mini-batches. Furthermore, our method can be combined with these methods in order to reduce the memory footprint even further (in fact, that is what we do!).
- Discuss other attempts at improving RFF (orthogonal RFF, Spherical Structured Feature Maps, Gaussian quadrature, quasi-monte carlo, etc.), and why we are different. Discuss that while these methods often improve the kernel approximation error for a fixed number of features, they do not increase the rank for a fixed amount of memory.
- Perhaps discuss methods of improving Nystrom? Sampling distributions, etc.

3 Preliminaries and Notation

3.1 Background: Kernel approximation

- Give a very brief overview of kernel approximation, Nystrom and RFF.
- Discuss the amount of memory each of these methods uses during training. It is important to be very clear about how we are measuring memory.

86 3.2 Avron Generalization Bound

- Present their bound, and discuss in moderate detail.

$$R(\tilde{f}) \leq \frac{1}{1 - \Delta} \hat{R}_K(f) + \frac{\Delta}{1 + \Delta} \cdot \frac{\text{rank}(\tilde{K})}{n} \cdot \sigma_v^2$$

- 87 • Mention that we can understand Δ in terms of $D = \max_i |\log(\sigma_i + \lambda) - \log(\tilde{\sigma}_i + \lambda)|$
88 ($1 + \Delta \geq \exp(D)$), as a way of visualizing the bound; if you plot both spectra $+\lambda$ in
89 log-scale, the Δ is the maximum gap between the two lines (if eigenvectors are the same.
90 So even if the eigenvectors are different, this us gives lower-bound for the true value of Δ).
91 • Discuss their result which bounds the number of random Fourier features needed to get a
92 specific Δ value.

93 3.3 Notation

- 94 • Present notation for Nystrom and RFF features, kernel matrices, etc.
- 95 • Specify that n denotes number of datapoints, d is data dimension, m is number of random
96 features, etc.

97 4 Nyström vs. Random Fourier Features

- 98 • In this section, we will present a number of experiments comparing the performance of
99 Nyström and RFF. We show that under a fixed memory budget, RFF outperforms Nyström,
100 and this is in spite of the fact that Nyström often attains lower kernel approximation error
101 (Frobenius/spectral).
102 • we will explain this using fixed design theory, and the Δ parameter from [1].

103 4.1 Nyström vs. RFF: Empirical comparison

- 104 • Across various classification/regression datasets (TIMIT, yearpred, CovType, Census, Adult),
105 we compare Nyström vs. RFF, for a wide range of numbers of features.
- 106 • We compare these two methods in various ways: 1) For a fixed # of features, 2) For a fixed
107 amount of memory, 3) for a fixed kernel approximation error (x-axis = spectral/Frobenius
108 norm). For the y-axis in these plots, we consider kernel approximation error, and downstream
109 performance.
- 110 • We observe that RFF outperforms Nyström for a fixed memory budget, and for a fixed kernel
111 approximation error.

112 4.2 Nyström vs. RFF: “ Δ -analysis”

- 113 • Across a few small datasets (Census, Adult), we measure the Δ parameter, and plot the
114 performance of Nyström vs. RFF in terms of this parameter.
- 115 • We observe that Nyström and RFF models that have similar Δ values, have similar general-
116 ization performance. We show that this correlation is much stronger than the correlation
117 between kernel approximation error and generalization performance.
- 118 • We plot the value of Δ attained by Nyström and RFF as a function of memory and number
119 of features. We show that for a fixed memory budget, RFF attains much smaller values of Δ .

120 5 Low-Precision Random Fourier Features (LP-RFF)

- 121 • We introduce low-precision RFF **AM: Should we do this earlier, in a “Methods” section?**
- 122 • We present our theoretical results about LP-RFF, bounding it’s approximation error (maybe
123 push this to appendix), and bounding it’s Δ , as a function of the number of quantized
124 features, and the precision being used.
- 125 • We show our experiments, which have three parts: (1) Explore the performance on big tasks,
126 (2) Explain trends on big experiments by doing detailed analysis on small tasks, (3) show
127 low-precision training doesn’t degrade performance much relative to FP training.

128 5.1 LP-RFF: Method

- 129 • Present LP-RFF method in detail.

130 5.2 LP-RFF: Theory

- 131 • Discuss bounds on kernel approximation performance (if we have space).
- 132 • Discuss bounds on Δ .

133 5.3 LP-RFF: Experiments

- 134 • We show performance on the “big tasks”.
- 135 • We explain these results in terms of the “small tasks”, where we measure Δ , showing that
- 136 below a certain number of bits, Δ can increase a lot. If we have space, we explain the
- 137 trade-offs here a bit, showing in the fixed design setting that higher noise means higher λ^* ,
- 138 which means fewer bits can be used until Δ increases a lot.
- 139 • We briefly present low-precision training, showing that performance doesn’t degrade much
- 140 relative to full-precision training.

141 6 Conclusion

- 142 • We have shown the practical significance of the recent theory of [1]. This came in two parts:
- 143 (1) It explains why RFF outperforms Nyström, even when they both have the same kernel
- 144 approximation error; a direct consequence is that under a fixed memory budget, RFF is
- 145 better. (2) It inspired us to define, analyze, and test the performance of LP-RFF.
- 146 • We showed that we can do much better, for a fixed memory budget, using LP-RFF, because
- 147 they are a much higher-rank decomposition of the kernel matrix.
- 148 • This opens the door for larger-scale kernel approximation experiments, which can leverage
- 149 recent advances in chips, etc, for fast training. For future work, we plan to implement
- 150 efficient version of our LP-RFF training, to be able to apply this to large/challenging tasks
- 151 (e.g., speech recognition, computer vision).

152 References

- 153 [1] Haim Avron, Michael Kapralov, Cameron Musco, Christopher Musco, Ameya Velingker, and
- 154 Amir Zandieh. Random fourier features for kernel ridge regression: Approximation bounds and
- 155 statistical guarantees. In *Proceedings of the 34th International Conference on Machine Learning,*
- 156 *ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, pages 253–262, 2017.
- 157 [2] Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances*
- 158 *in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference*
- 159 *on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December*
- 160 *3-6, 2007*, pages 1177–1184, 2007.
- 161 [3] C.K.I. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In T.K.
- 162 Leen, T.G. Dietterich, and V. Tresp, editors, *Advances in Neural Information Processing Systems*
- 163 *13*, pages 682–688. MIT Press, 2001.
- 164 [4] Tianbao Yang, Yu-Feng Li, Mehrdad Mahdavi, Rong Jin, and Zhi-Hua Zhou. Nyström method
- 165 vs random fourier features: A theoretical and empirical comparison. In *Advances in Neural*
- 166 *Information Processing Systems 25: 26th Annual Conference on Neural Information Processing*
- 167 *Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United*
- 168 *States.*, pages 485–493, 2012.