

Homework 5 Selected Solutions

APPM 5650 Fall 2021

Randomized Algorithms

Due date: Monday, Sept 27 2021 at 10:20 AM

Theme: Randomized Sketches

Instructor: Prof. Becker

Problem 1: [READING and MATH] Read exercise 4.1.4 and Lemma 4.1.5 about approximate identities, from Vershynin's 2018 "High-dimensional probability" book.

Deliverable: Do exercise 4.1.6: if $A \in \mathbb{R}^{m \times n}$ and all the singular values of A are in the set $[1 - \delta, 1 + \delta]$ for some $\delta > 0$, then $\|A^T A - I_n\| \leq 3 \max(\delta, \delta^2)$. Note: you may assume $m \geq n$, otherwise the hypothesis cannot be true unless $\delta \geq 1$ (why?) and the result is not interesting.

Solution:

Let the thin SVD of A be $A = U\Sigma V^T$. Then since $U^T U = I_m$ and $V V^T = V^T V = I_n$,

$$\begin{aligned} \|A^T A - I_n\| &= \|V \Sigma^T U^T U \Sigma V^T - I_n\| \\ &= \|V \Sigma^2 V^T - I_n\| \\ &= \|V \Sigma^2 V^T - V V^T\| \\ &= \|V(\Sigma^2 - I)V^T\| \\ &= \|\Sigma^2 - I\| \\ &= \|\sigma^2 - \mathbf{1}\|_\infty \end{aligned}$$

where we switched from the spectral norm (on matrices) to the ℓ_∞ norm on vectors at the last step, and used the unitary invariance of the spectral norm. Then because $\sigma_1 \geq \sigma_2 \dots$, we have

$$\begin{aligned} \|\sigma^2 - \mathbf{1}\|_\infty &= \max_i |\sigma_i^2 - 1| \\ &= \max_i \max(\sigma_i^2 - 1, 1 - \sigma_i^2) \\ &= \max(\sigma_1^2 - 1, 1 - \sigma_n^2) \\ &\leq \max((1 + \delta)^2 - 1, 1 - (1 - \delta)^2) \\ &= \max(\delta + 2\delta^2, \delta - 2\delta^2) \\ &= \delta + 2\delta^2 \\ &\leq \max(\delta, \delta^2) + 2 \max(\delta, \delta^2) = 3 \max(\delta, \delta^2) \end{aligned}$$

If $m < n$ then $A^T A$ has zero singular values and we know that $\|A^T A - I_n\| \geq 1$.

Problem 2: [PROGRAMMING] Try various sketches, do they converge to the identity as number of repeated trials N increases to ∞ ? Each draw of a sketch S_i is $m \times n$, but unlike the above exercise, here $m < n$, and instead of looking at $S_i^T S_i \approx I_n$, we ask about

$$\left\| \frac{1}{N} \sum_{i=1}^N S_i^T S_i - I_n \right\| \tag{1}$$

as $N \rightarrow \infty$ (in either spectral or Frobenius norm).

Also consider a histogram of outcomes of sketching

- a) $x \in \mathbb{R}^n$ where x is very sparse, and
- b) $x \in \mathbb{R}^n$ where all the entries of x are similar in value.

Look at $\|Sx\|^2/\|x\|^2$.

Generate S in all of the following manners:

- a) Gaussian
- b) Haar (cf. HW 2). Note: this is really part of a Haar matrix – it should *not* be a full $n \times n$ orthogonal matrix. It should be $m \times n$ like all the other sketches, but with orthonormal rows.
- c) Fast Johnson-Lindenstrauss, $S = RHD$ where R samples m of n rows uniformly at random, H is the discrete cosine transform, and D is diagonal with Rademacher random variables on the diagonal.
- d) Simple sampling of m of n rows (e.g., $S = R$ in the above notation)
- e) *at least one* of the following variants
 - i. the approach from “How to Fake Multiply by a Gaussian Matrix” (Kaprалov, Potluru, Woodruff, ICML 2016, <https://arxiv.org/abs/1606.05732>)
 - ii. the FJLT approach from the original FJLT paper, where R is not simple sub-sampling but rather the matrix R is such that each entry is 0 with a certain probability, otherwise drawn from an appropriately scaled Gaussian (see “Approximate nearest neighbors and the fast Johnson-Lindenstrauss transform”, Ailon and Chazelle, 2006, STOC, <https://dl.acm.org/citation.cfm?id=1132597>)
 - iii. The low-density parity check (LDPC) codes used for the R term in the FJLT, from “Low Rank Approximation using Error Correcting Coding Matrices” (Ubaru, Mazumdar, Saad, ICML 2015, <http://proceedings.mlr.press/v37/ubaru15.html>).
 - iv. CountSketch (“Simple and deterministic matrix sketching”, Edo Liberty, KDD 2013, <https://arxiv.org/abs/1206.0594>; or Clarkson and Woodruff, <https://doi.org/10.1145/3019134>)
 - v. the approach from “Very Sparse Random Projections” (Li, Hastie, Church, KDD 2006, <https://dl.acm.org/citation.cfm?id=1150436>), where each entry S_{ij} is iid with value $\sqrt{3}$ with probability 1/6, value $-\sqrt{3}$ with probability 1/6, and value 0 with probability 2/3.
 - vi. the $S = \text{const} \cdot THG\Pi HB$ approach in eq. (7) from “Fastfood – Approximating Kernel Expansions in Loglinear Time” (Le, Sarlós, Smola, ICML 2013, <https://arxiv.org/abs/1408.3060>)
 - vii. Fast Johnson-Lindenstrauss, where H is now a Hadamard transform instead of the DCT
 - viii. Anything else you can find in the literature, e.g., see Fig. 1

For each of the (at least 5) methods for S , plot the quantity in Eq. (1) as a function of N , and check that it behaves as you think it should.

Try with several choices of n and m to make sure your scaling is correct; most of the sketch matrices need a $\sqrt{n/m}$ or $\sqrt{1/m}$ scaling.

Deliverables: A plot of the quantity in Eq. (1) as a function of N for each type of sketch, and histograms of $\|Sx\|_2^2/\|x\|_2^2$ for the different types of sketches and the two types of x mentioned above. Include at least code snippets for how you generated each type of matrix.

Are these good metrics to look at? Can you think of other metrics to look at?

Hint: do the scaling slightly wrong on purpose, and make sure that your plots detect this error, so that you can be sure your true scaling really is correct. What kind of graph is best

type	sketch	complexity
ℓ_1	Dense Cauchy [Sohler and Woodruff, 2011]	$\mathcal{O}(nd^2 \log n + d^3 \log d + d^{\frac{11}{2}} \log^{\frac{3}{2}} d / \epsilon^2)$
ℓ_1	Fast Cauchy [Clarkson et al., 2013]	$\mathcal{O}(nd \log n + d^3 \log^5 d + d^{\frac{17}{2}} \log^{\frac{3}{2}} d / \epsilon^2)$
ℓ_1	Sparse Cauchy [Meng and Mahoney, 2013a]	$\mathcal{O}(\text{nnz}(A) \log n + d^7 \log^5 d + d^{\frac{19}{2}} \log^{\frac{3}{2}} d / \epsilon^2)$
ℓ_1	Reciprocal Exponential [Woodruff and Zhang, 2013]	$\mathcal{O}(\text{nnz}(A) \log n + d^3 \log d + d^{\frac{13}{2}} \log^{\frac{3}{2}} d / \epsilon^2)$
ℓ_1	Lewis Weights [Cohen and Peng, 2015]	$\mathcal{O}(\text{nnz}(A) \log n + d^3 \log d + d^{\frac{9}{2}} \log^{\frac{3}{2}} d / \epsilon^2)$
ℓ_2	Gaussian Transform	$\mathcal{O}(nd^2 + d^3 \log(1/\epsilon)/\epsilon)$
ℓ_2	SRHT [Tropp, 2011]	$\mathcal{O}(nd \log n + d^3 \log n \log d + d^3 \log(1/\epsilon)/\epsilon)$
ℓ_2	Sparse ℓ_2 embedding [Cohen, 2016]	$\mathcal{O}(\text{nnz}(A) \log n + d^3 \log d + d^3 \log(1/\epsilon)/\epsilon)$
ℓ_2	Refinement Sampling [Cohen et al., 2015a]	$\mathcal{O}(\text{nnz}(A) \log(n/d) \log d + d^3 \log(n/d) \log d + d^3 \log(1/\epsilon)/\epsilon)$

Figure 1: Types of sketches as listed in “Weighted SGD for ℓ_p Regression with Randomized Preconditioning” by Yang, Chow, Ré and Mahoney, SODA 2016, <https://arxiv.org/abs/1502.03571>

(e.g., log or linear scaling on the axes)? What do you expect the histogram to be centered around? Try changing m, n to catch bugs. Usually we want $n \gtrsim 100$ otherwise if it is very small you can get some weird effects.

Solution:

See the figures in Fig. 2 and the source code as well. The take-aways: Eq. (1) is not a good “catch-all” metric to describe how good a sketch is; for example, the simple sub-sampling sketch R does great by this metric (partly because $R^T R$ is always diagonal). Looking at $\|Sx\|$ for several types of x , we can see that simple sub-sampling has a huge variance, much worse than all the other types. Other than R , all other types of sketches have roughly similar performance. The bias in the histograms is probably because I have asked Matlab to fit the histogram with a Gaussian distribution, but that’s not a perfect fit since the variable $\|Sx\|^2$ is non-negative, so it cannot be symmetric, and hence Matlab is compensating by estimating a slightly larger mean.

```

1  rng(0); % Make it reproducible
2
3  % Dimensions:
4  % n    = 10; m = 1; m = 9; % try a variety...
5  n     = 256; m = 20;
6
7  totalN = 5e3; % Number of trials
8
9  ALGO_NAMES = {'Gaussian', 'Haar', 'FJLT', 'Simple', ...
10               'CountSketch', 'Very sparse RP'};
11  nAlgos     = length( ALGO_NAMES );
12  SUMS       = cell( nAlgos , 1 );
13  [SUMS{:}]  = deal( zeros(n) );
14  [errs, err_x, err_y, condList] = deal( zeros( totalN, nAlgos ) );
15  errFcn      = @(S) norm( S - eye(n), 'fro' ); % faster than spectral norm
16
17  % — Pick two different types of vectors to observe
18  x = zeros(n,1); % very sparse test vector
19  x( randperm(n,5) ) = randn(5,1);
20  x = x/norm(x);
21
22  y = randn(n,1); % more uniform arbitrary test vector
23  y = y/norm(y);
24
25  fprintf('Starting... ');
26  for N = 1:totalN
27      if ~mod(totalN,100)
28          fprintf('\b\b\b\b\b\b%4.1f%%', N/totalN*100 );

```

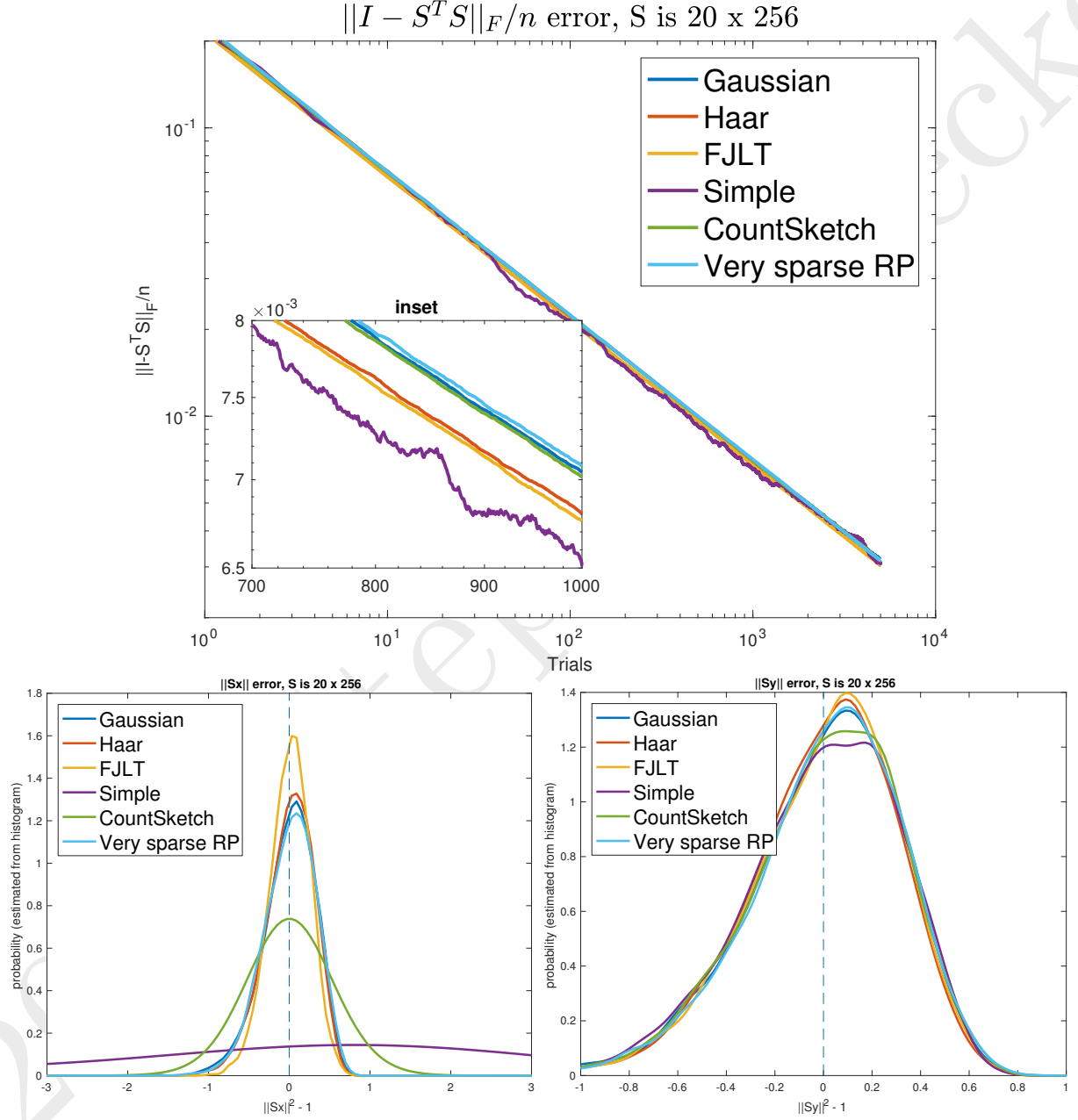


Figure 2: Top: error in Eq. (1) as a function of number of samples N ; Bottom, left: error $\|Sx\|^2 - 1$ for x very sparse and $\|x\|=1$; right, error $\|Sy\|^2 - 1$ for y with Gaussian entries and normalized so $\|y\|=1$.

```

29     end
30
31     for ALGO = 1:nAlgos
32         switch ALGO
33             case 1 % Gaussian
34                 S = 1/sqrt(m)*randn(m,n);
35             case 2 % Haar
36                 [S,R] = qr(randn(n,m),0);
37                 S = (S*diag(sign(diag(R))))'; % debias
38                 S = sqrt(n/m)*S;
39             case 3 % FJLT with DCT
40                 ind = randperm(n,m);
41                 HD = dct( diag(sign(randn(n,1))) );
42                 S = sqrt(n/m)*HD(ind,:);
43             case 4 % Simple sub-sampling
44                 ind = randperm(n,m);
45                 I = eye(n);
46                 S = sqrt(n/m)*I(ind,:);
47             case 5 % Count sketch
48                 D = diag( sign(randn(n,1)) );
49                 targetRows = randi(m,n,1);
50                 S = zeros(m,n);
51                 for j = 1:n
52                     i = targetRows(j);
53                     S(i,:) = S(i,:) + D(j,:);
54                 end
55             case 6 % Very sparse random projections
56                 S = zeros(m,n);
57                 ind = rand(m,n);
58                 S( ind < 1/6 ) = -1;
59                 S( ind > 5/6 ) = +1;
60                 S = sqrt(3/m)*S;
61         end
62
63         SUMS{ ALGO } = SUMS{ ALGO } + S'*S;
64         errs(N,ALGO) = errFcn( SUMS{ALGO}/N );
65         err_x(N,ALGO) = norm(S*x);
66         err_y(N,ALGO) = norm(S*y);
67     end
68 end
69 fprintf('\nFinished.\n');
70 % Make it norm^2
71 err_x = err_x.^2;
72 err_y = err_y.^2;
73 %% Plot Frobenius norm error of deviation from I
74 fh = figure(1); clf;
75 loglog( errs/n, 'linewidth',2 )
76 legend( ALGO_NAMES,'fontsize',18);
77 ylim([2e-3,.2]);
78
79 h = gca;
80 axes2 = copyobj(h,fh);
81 axes2.Position = [0.18 0.18 0.35 0.35];
82 xlim(axes2,[7e2,1e3]);

```

```

83 ylim(axes2,[6.5e-3,8e-3]);
84 title(axes2,'inset');
85
86 title(h,sprintf('$||I-S^TS||_F/n$ error, S is %d x %d', m, n ) ,...
87     'interpreter','latex','fontsize',18);
88 xlabel(h,'Trials');
89 ylabel(h,'||I-S^TS||_F/n');
90 export_fig 'HW5_identityError' '-pdf' -transparent
91 %% Or, look at probability of exceeding a certain threshold (histogram)
92 grid = linspace( -3,3,100 );
93 figure(4); clf;
94 for ALGO = 1:nAlgos
95     data = 1 - err_x(:,ALGO);
96
97     % histogram( data , 'Normalization','Probability','DisplayStyle','stairs')
98
99     pd = fitdist(data,'kernel','Kernel','normal');
100     y = pdf( pd, grid );
101     plot( grid, y, 'linewidth', 2 )
102     hold all
103 end
104 legend( ALGO_NAMES, 'location','northwest','fontsize',18 );
105 title(sprintf('||Sx|| error, S is %d x %d', m, n ) );
106 hl = line( [0,0],[0,1.8],'linestyle','—','HandleVisibility','off');
107 xlabel('||Sx||^2 - 1');
108 ylabel('probability (estimated from histogram)');
109 export_fig 'HW5_x_error' '-pdf' -transparent
110 %%
111 figure(5); clf;
112 grid = linspace( -1,1,100 );
113 for ALGO = 1:nAlgos
114     data = 1 - err_y(:,ALGO);
115     % histogram( data , 'Normalization','pdf','DisplayStyle','stairs');
116
117     pd = fitdist(data,'kernel','Kernel','normal');
118     y = pdf( pd, grid );
119     plot( grid, y, 'linewidth', 2 )
120     hold all
121 end
122 legend( ALGO_NAMES, 'location','northwest','fontsize',18 );
123 title(sprintf('||Sy|| error, S is %d x %d', m, n ) );
124 hl = line( [0,0],[0,1.4],'linestyle','—','HandleVisibility','off');
125 xlabel('||Sy||^2 - 1');
126 ylabel('probability (estimated from histogram)');
127 export_fig 'HW5_y_error' '-pdf' -transparent

```