1 We would like to thank the reviewers for the careful and thorough reading of our submission. We are appreciative of the
2 many suggestions for improvements and insightful questions. In the limited space below, we respond to some of the
3 main concerns raised by the three reviewers.

4 **The novelty of our analysis in the case $p = 2$ [Reviewer 1 and Reviewer 2]** As the reviewers note, for the case
5 $p = 2$, the tight approximation bound of $CSS$ is already known in the classic work [14]. However, our work has a
6 different goal: to provide a unified way of analyzing approximation bounds for all different values of $p$. With this goal
7 in mind, we introduced the Riesz-Thorin theorem as a general framework, and found that the existing analysis in [14]
8 unfortunately does not fit in this framework, due to the technical differences stated below.

9 The main technical difference can be found in line 219-222: the Riesz-Thorin theorem requires that Equation (4)
10 should hold for *all* $\{a_t\}_{t \in \binom{[m]}{1}} \in \mathbb{C}^{\binom{[m]}{1}}, \{b_J\}_{J \in \binom{[m]}{k}} \in \mathbb{C}^{\binom{[m]}{k}}$. In comparison, to prove that $CSS$ is a $\sqrt{k+1}$
11 approximation, one only needs to show that the equation holds for $\{a_t\}_{t \in \binom{[m]}{1}} \in \mathbb{C}^{\binom{[m]}{1}}, \{b_J\}_{J \in \binom{[m]}{k}} \in K$, where $K$
12 is a subset of $\mathbb{C}^{\binom{[m]}{k}}$ defined as $K = \{\{b_J\}_{J \in \binom{[m]}{k}} | b_J = \det(S_J), S \in \mathbb{C}^{k \times m}\}$. It is easy to see that the requirement
13 of Riesz-Thorin is significantly stronger, since the set $K$ is determined by only $km$ parameters (the matrix $S$), while
14 $\mathbb{C}^{\binom{[m]}{k}}$ is a $\binom{m}{k}$-dimensional space. In other words, we need a much stronger inequality in order to apply Riesz-Thorin
15 theorem, hence we provided a brand new proof for the $p = 2$ setting (line 421-423, as mentioned by reviewer 1) which
16 is completely different from [14].

17 **The choice of $p$ in $l_p$ low rank approximation [Reviewer 1 and Reviewer 2]** While our paper is mostly theoretical,
18 we do believe that choosing an appropriate value of $p$ can make a difference in practice. The $\ell_p$ low rank approximation
19 problem has attracted interest relatively recently, see for instance ICML 2017 paper [13]. A related problem $\ell_p$ linear
20 regression however has been studied extensively in the statistics community, and these two problems share similar
21 motivation. In particular, if we assume a statistical model $A_{ij} = A_{ij}^\star + \varepsilon_{ij}$, where $A^\star$ is a low rank matrix and $\varepsilon_{ij}$ are
22 i.i.d. noise, the different values of $p$ correspond to the MLE of different noise distributions, say $p = 1$ for Laplacian
23 noise and $p = 2$ for Gaussian noise. To capture a broader range of realistic noises in complex datasets, it is beneficial to
24 expand our choices beyond the standard ones $(1, 2, \infty)$.

25 **Sampling Based Algorithms [Reviewer 2]** This is definitely an important direction for the future works that we
26 intend to explore. Our determinantal weights are indeed inspired by [14] and we will include the comparison in our final
27 version. The sampling methods for the $p = 2$ case, e.g. volume sampling, are closely related to the determinantal point
28 process. To generalize this approach to $\ell_p$ setting, we need an efficient way to implement the exponentiated variant of
29 volume sampling (i.e. sample a subset $S$ with probability proportional to $\det(V_S)^\alpha$). To the best of our knowledge, this
30 problem has not been resolved yet. In fact, even computing the normalizing constant of the distribution is difficult - it
31 was stated as an open problem in Section 7.2 of the survey "Determinantal point processes for machine learning". We
32 refer the reviewer to the NeurIPS 2018 paper "Exponentiated Strongly Rayleigh Distributions" for recent progress on
33 this problem.

34 **Clarification on Theorem 1.2 [Reviewer 3]** We thank the reviewer for pointing out the ambiguity of the informal
35 statement of the theorem. We proved that there are infinitely many different values of $k$, such that for each $k$, there
36 exists a matrix $A$ such that CSS cannot do better than $(k + 1)^{1 - \frac{1}{p}}$ approximation. The dimensions of the matrix $A$ ($m$
37 and $n$) are not fixed for these different $k$'s.

38 **Optimization landscape of low rank approximation [Reviewer 3]** We refer the reviewer to the paper "Neural
39 networks and principal component analysis: Learning from examples without local minima" where the paper shows
40 that for low-rank approximation problem even as easy as PCA ($p = 2$), saddle points exist. We are happy to modify the
41 relevant sentence in our paper as "Unfortunately the loss surface of the problem suffers from many saddle points" and
42 replace [12] with the reference above.

43 **The optimality gap of the bound when $1 < p < 2$ [Reviewer 3]** (line 83) We note that the lower bound (Theorem
44 1.2) also applies to the case $1 < p < 2$. Therefore, we have an upper bound of $(k + 1)^{\frac{1}{p}}$ and a lower bound $(k + 1)^{1 - \frac{1}{p}}$,
45 hence we were correct on the $(k + 1)^{\frac{2}{p} - 1}$ optimality gap.

46 **Improvements on Exposition [Reviewer 1,2,3]** We thank all three reviewers for their suggestions on the exposition
47 and will take them into account in the final version. In particular, we will clarify the difference of ordered and unordered
48 sets; add more explanation to the reduction in the proof of Lemma 2.2; improve the sentence structure in Lemma 2.3
49 and include the definition of CUR factorization in the introduction.