1  We thank all reviewers for the constructive comments. Below, we label each comment by Reviewer#.Comment#.

2  R1.1 [ascent when $p$ exceeds the range in Thm 2, descent-region] As $p$ becomes larger, eventually the model error
3  will increase again. The reason is that there will be some columns of $\mathbf{X}_{\text{train}}$ very similar to those of true features.
4  Then, BP will pick those very similar but wrong features, and the error will approach null risk. As a result, the whole
5  double-descent behavior should indeed be <descent-ascent(p<n)-descent(p>n)-ascent> (as shown in Fig. 1, the same is
6  also true for $\ell_2$-minimization). By the "descent region", we mean the second "descent" part. Thus, Fig. 1 shows that
7  model error is insensitive to $\|\beta\|_2$ before it ascends again. We will make these clear in the final version.
8  R1.2 [constant term on the RHS of (9)] Currently we can reduce, but cannot completely get rid of, the constant term in
9  (9). It does not appear in Prop. 4 due to the different path of proof (see (88) of Supplemental Material).
10 R1.3 [upper-bound of BP is not affected by $\|\beta\|_2$] This insensitivity to $\|\beta\|_2$ is a special feature of BP. This can be seen
11 by considering the following (perhaps over-simplified) example. Given one training datum ($x \in \mathbb{R}^2, y \in \mathbb{R}$) generated
12 by the true model $y = x^T \beta + 0.1$ (here 0.1 is the noise), where $\beta = [\beta_1, \ \beta_2]^T$, and $x = [1, \ 0.5]^T$. Assume $\beta_2 = 0$
13 (for sparsity). If we minimize $\|\hat{\beta}\|_1$ subject to $y = x^T \hat{\beta}$, we get $\hat{\beta} = [\beta_1 + 0.1, \ 0]^T$ (and thus the model error is always
14 $\|\beta - \hat{\beta}\|_2 = 0.1$) for any $\beta_1$. In a higher dimensional space, due to a similar reason (under certain conditions), the
15 generalization error of BP could be insensitive to $\|\beta\|_2$.

16 R2.1 [upper bound may be larger than the null risk] We agree. On the other hand, one important goal is to understand
17 in what settings min-$\ell_1$ is better than min-$\ell_2$. Eq. (9) can already provide new and useful answers, e.g., when $p$ is
18 exponentially large and when noise is low, the model error of min-$\ell_1$ is much lower than that of min-$\ell_2$ (null risk).
19 R2.2 [bounds are loose, range of $p$ does not include $p$ near $n$, squared loss and isotropic Gaussian data] We acknowledge
20 these limitations. We believe the results can be extended to independent non-Gaussian data, while other generalization
21 may require more work. However, we also would like to point out that, even under our model, the analysis for
22 BP is significantly more challenging than min-$\ell_2$ solutions because, unlike min-$\ell_2$ solutions that can be written as
23 pseudo-inverses, min-$\ell_1$ solutions have no closed form. Thus, we have to resort to large-$p$ asymptotes to capture the
24 model error of BP, which is why the characterization of $p \approx n$ remains unknown. We believe that this difficulty is the
25 reason why there are so few results on overfitting BP solutions. Despite these limitations, our results still successfully
26 reveal key insights on the difference between min-$\ell_1$ and min-$\ell_2$ solutions, which we believe are of significant interest
27 to the community (as the other reviewer commented).
28 R2.3 [required sparsity relates to $n$] Such requirements are not uncommon, e.g., Thm. 3.1 of [16] also requires the
29 sparsity to be below some function of $M$ (incoherence of $\mathbf{X}_{\text{train}}$), which is related to $n$ according to Prop. 9 in our paper.
30 R2.4 [establish double descent by the upper-bound] Our result already implies double descent. The reason is that, when
31 $p \le n$, the MSE solution is independent of $\ell_1$ vs. $\ell_2$. From known results for min-$\ell_2$, the model error has a peak when
32 $p \to n$ from below. Thus, our upper bound implies that the model error of BP has to descend from that peak.
33 R2.5 [related work: arXiv preprint arXiv:2002.01586 (2020)] Thank you for the pointer and we will cite this reference.
34 However, the reference studies classification, while we study regression. The notion of "fitting the data" is quite
35 different between the two models. Hence, we feel that the conclusions are not directly comparable.
36 R2.6 [(minor) Fig. 2 starts from different positions (x axis)] That is because we focus on the interpolating regime $p > n$
37 and $n$ varies. We will add back the regime $p < n$ for the final version. Thank you for the suggestion!

38 R2.7 / R3.1 [connection to DNN, motivation of studying $\ell_1$ minimization] We agree that the connection between min-$\ell_1$
39 and DNN is not clear yet. Similar to other related work, our approach is to use linear models as a starting point to
40 understand which types of overfitting solutions approximate the generalization power of DNN better. Our hope is that
41 such analysis would eventually lead to better training methods than SGD (which is closely related to $\ell_2$ minimization).

42 R3.2 [Gaussian features assumption, limited range of $p$ compared with $\ell_2$, constants in Theorem 2 are not sharp] We
43 acknowledge these limitations. Please also refer to our response to R2.2.

44 R4.1 [$n$ and $p$ are large for theory] One reason for large $p$ and $n$ is that we aim for with-high-probability results, which
45 leads to stricter conditions than average-case analysis. Another reason is that we have not optimized the constants.
46 Nonetheless, the numerical results suggest that the main predicted trends hold for much smaller $n$ and $p$. We will clarify
47 this in the final version. Thank you for the suggestion!
48 R4.2 [model error independent of the signal strength] Thanks for sharing your code! There seems to be a small
49 yet consequential bug on Line 39, which should be `print(np.linalg.norm(x[:p] - beta_exact))` instead of
50 `print(np.linalg.norm(x[:p] - beta_exact[:, 0]))`. (The original code returns the norm of a $p \times p$ matrix
51 instead of a $p \times 1$ vector.) After correcting this line, the result becomes: 0.019, 0.020, 0.021, and 0.021 for $\|\beta\|_2$
52 equals to 0.1, 1, 10, and 100, respectively, which verifies our conclusion "model error can be independent of the signal
53 strength". Please also refer to our response to R1.3 for the intuition. We will post our code for the final version.