1 We thank the reviewers for their detailed and constructive feedback.

2 **R1:** *It is not clear... how the assumption of bounded gradients...:* We would like to point out that the assumption of
3 bounded gradient is only employed in the theoretical analysis of the bias. Such assumptions are quite common in the
4 analysis of ML algorithms (see e.g., [Hazan and Kale, 2014]). As mentioned in our submission, in many settings, this
5 holds because of the clipping of the gradients (see e.g., [Goodfellow, Bengio and Courville, 2016]).

6 *Same as above for the normalization of the class embeddings and input embeddings...:* As discussed in lines 167-174
7 of our submission, it is correct that we assume the embeddings to be normalized for RF-softmax, however, such
8 normalization is widely used in practice (e.g., see references [26], [27], and [28]). Furthermore, we have empirically
9 shown (on both NLP and extreme classification datasets) that with proper setting of $\tau$, the normalized embeddings do
10 not degrade the final performance of the model.

11 **R2:** *It feels that there is a skip between Sections 2 and 3 ...:* We thank the reviewer for the constructive comment. We
12 will enhance the presentation by better motivating the kernel-based sampling to approximate exponential families (i.e.,
13 RF-softmax method and its analysis) and smoothing the transition from section 2 to section 3. We have motivated the
14 use of random fourier features by comparing its performance against two most natural candidates in Table 1. We will
15 improve the relevant discussion in the final version.

16 *The paper perhaps lacks a discussion on approaches ... such as minimization of Fisher divergence ... an approximation*
17 *scheme for log-partition function was considered in [2, Section 2.3]...:* We thank the reviewer for the suggestion.
18 In [Hyvarinen 2005], the partition function $Z$ is just the function of model parameter and thus disappears in scoring
19 function. However, in our case, the partition function depends on the input $h$ (which changes during the training).
20 Therefore, while calculating the score function (taking derivative of $Z$ with respect to $(h, c)$), the partition function has
21 a non-trivial contribution. As for [Vembu et al., 2009], given ways to generate uniform samples for the set of classes,
22 they propose a MCMC approach to sample a class with a distribution that is close to full softmax distribution. Such
23 methods do not come with precise sample/computational complexity guarantees. We will include a discussion in the
24 final version.

25 *The notation is not very clear, especially in the appendix...:* We will highlight the distribution/random variables with
26 respect to which we take the expectations. Also, we will try to eliminate any inconsistency/ambiguity regarding the
27 notation elsewhere in the paper.

28 *In Eq. (5), I was not clear for the motivation of $m$ until checking the appendix and reading the whole paper...:* Adjusting
29 the logits for negative classes using their expected number of occurrence [Bengio Senecal, 2008] is critical for the
30 unbiasedness of sampled softmax loss, e.g., it ensures that $Z'$ is an unbiased estimator of $Z$. We will add a comment to
31 clarify the process of adjusting the logits in (5).

32 *Related work can also be improved and relevance of the work in the context of extreme classification...:* We will
33 highlight relevant papers in extreme classification literature in the final version.

34 **R3:** *What do you get by applying Theorem 1 to RF-softmax ...:* As pointed out in the discussion following Theorem 1,
35 the result provides a guidance for selecting a sampling distribution with low bias (by highlighting the requirement of
36 tight multiplicative approximation). We will combine Theorem 1 and 2 (at least in the setting with large D) in the final
37 version to obtain the bounds for RFF.

38 *Theorem 2 and Remark 2: Where is the dependence on D coming from in $o_D(1)$? ...:* Thanks for the comment. We
39 will include a comment on how $\gamma_2$ depends on $\gamma_1$. As it's clear from the proof of Theorem 2 and the statement of
40 Remark 2, one can choose $\gamma_1 = const\sqrt{(d \log D)/D}$. Now $\gamma_1$ (and thus $\gamma_2$) scales as $o_D(1)$ (while keeping other
41 parameters fixed).

42 *What happens if the inner product $\phi(c_i)^T \phi(h)$ is negative?:* With normalized embeddings, for finite $\nu$, $e^{(\nu h^T c_i)}$ is
43 strictly positive. Therefore, for reasonably large $D$, with high probability, $\phi(c_i)^T \phi(h)$ should be non-negative. In
44 general, one can replace $\phi(c_i)^T \phi(h)$ with $max(0, \phi(c_i)^T \phi(h))$ (with minor modifications in the proofs).

45 *How does RF-softmax compare to hierarchical softmax...:* As discussed in [Balnc Rendle, 2018] and references therein,
46 in many tasks, the final solution of hierarchical softmax is worse than those of both full-softmax and sampled-softmax.

47 *"$q_j$ should provide a tight uniform multiplicative approximation of $e^{o_j}$." How tight is the bound...:* We believe that the
48 bounds presented in Theorem 1 are fairly tight. In particular, these bounds recover the unbiasedness of the gradient for
49 full softmax distribution.

50 *It would be good to show the calculation that verifies (15):* (15) follows from the existing literature on RFF. We plan to
51 include a citation to [Yu et al. 2016, Lemma 1].

52 We will fix all the typos pointed out by the reviewers and eliminate other remaining typos in the final version.