

1 We warmly thank the reviewers for their time and for sharing valuable feedback, we will use it to improve our paper.
2 Overall, we are encouraged by the reviewers’ reactions. Next we summarize our contributions at a high level, and then
3 we address each reviewer’s comments separately. We hope that our responses clear out any remaining concerns.

4 This paper should be regarded as a theoretical contribution that takes a critical stand on the “usual assumptions” on
5 which PAC-Bayes inequalities are based (i.e. (a) bounded loss, (b) data-free prior, and (c) i.i.d. data observations),
6 clarifying their role and illustrating how to obtain PAC-Bayes inequalities in cases where these assumptions are
7 removed. Importantly, our work enables new PAC-Bayes inequalities with data-dependent priors. Furthermore, our
8 paper contributes a unified approach to understanding PAC-Bayes inequalities, and their distinctions. Our paper also
9 makes a case for the usefulness of formalizing data-dependent distributions as stochastic kernels, which may bring a
10 new perspective to the PAC-Bayes literature, in spite of the fact that stochastic kernels are well-known in other literature,
11 namely, on stochastic processes and their applications. We did also aim to write an informative and easy-to-read paper
12 that may help the larger machine learning community to connect with the PAC-Bayes literature.

13 **Response to Reviewer 1**

14 We are grateful for your comments on the notation, we gladly will address this in the revised paper. Also we appreciate
15 your feedback about the need to elaborate more on the novelty and relevance of our paper. We will add discussions
16 about the new kinds of results one can obtain now that were not possible before, including rates where relevant. We will
17 add a remark that there is indeed a direct relationship between moments of the loss and concentration of eigenvalues of
18 the covariance matrix in the least-squares example. Consider a simple noise-free linear regression $Y = \mathbf{X}^\top \mathbf{w}^*$, a “ridge”
19 prior $p^0(\mathbf{w}) \propto e^{-\lambda \|\mathbf{w}\|^2}$, and let $\hat{\Sigma}_0 = \frac{1}{n} \sum_i \mathbf{X}_i \mathbf{X}_i^\top - \mathbb{E}[\mathbf{X}_1 \mathbf{X}_1^\top]$. Then, for the log-exponential moment of the loss
20 we have the identity $\log \mathbb{E} \int e^{L(\mathbf{w}) - \hat{L}_S(\mathbf{w})} p^0(d\mathbf{w}) = \frac{d}{2} \log(2\pi) + \log \mathbb{E} |\det(\hat{\Sigma}_0 + \lambda \mathbf{I})|^{-\frac{1}{2}}$ which shows equivalence of
21 concentration of eigenvalues of $\hat{\Sigma}_0$ and concentration of the loss since $\lambda_i(\hat{\Sigma}_0) \rightarrow 0$ as $n \rightarrow \infty$ for i.i.d. instances.

22 **Response to Reviewer 2**

23 We will do a better job at discussing the usual formulation (for all posterior distributions) versus the formulation in
24 terms of a stochastic kernel. Their equivalence can be illustrated by discussing specific stochastic kernels. Note that
25 each theorem holds for an arbitrary kernel, so in that sense it holds “for all posteriors” as per the usual formulation
26 in the literature. On the other hand, as we mentioned in our paper, we think that it is important to make explicit the
27 data-dependence of the distributions, which is what the stochastic kernel formulation does. We will add a comment that
28 the mathematical proof of our theorem is valid for (convex) $F : \mathbb{R}^k \rightarrow \mathbb{R}$ with an arbitrary k , while $k = 2$ is relevant
29 for the known PAC-Bayes bounds. However, applications with $k > 2$ might emerge in the future, which is why we
30 think it is interesting to point out our theorem’s validity for arbitrary k , as this could enable new results. We are aware
31 that the refined form of Pinsker inequality ($\text{kl}(p\|q) \geq (p - q)^2 / (2q)$, for $p < q$) is stronger than the other one only
32 when $q < 0.25$, we apologize for the imprecision, this will be clarified in the revised paper. Thanks for pointing out the
33 missing reference Blanchard & Fleuret (2007), we will make sure to include it in the revised paper.

34 **Response to Reviewer 3**

35 Many thanks for your positive reaction to our work! Responding to the four raised points regarding the related literature:
36 (1) & (3) We will discuss in the introduction the work of Catoni (2007) in connection with data-dependent priors
37 (regular conditional probability distributions), and accordingly instead of claiming “firsts” we will highlight instead
38 that our work makes a case for the usefulness of representing data-dependent distributions as stochastic kernels, with
39 attribution to Catoni (2007) and Alquier (2008) as predecessors. (2) We will comment on our main theorem being valid
40 under general data-generation assumptions (i.e. not restricted to (c) i.i.d. data). (We are aware of this, sadly we missed
41 commenting on it, but happy to fix it.) We’ll expand the coverage of literature on PAC-Bayes bounds for non-i.i.d. data,
42 including Alquier & Wittenberger (2012). The pointer to Rio’s inequality is most helpful, we really appreciate this!
43 (4) We meant that in some cases we can have concentration of the smallest eigenvalue of a sample covariance matrix
44 even for unbounded instances, our reference is Section 5.4.2 of Vershynin (2011), discussing eigenvalue concentration
45 of a sample covariance matrix with heavy-tailed observations. We will add discussion about Holland (2019) too.

46 **Response to Reviewer 4**

47 We’d like to clarify that our Theorem 2 is not restricted to convex losses, only the function F is restricted to be convex
48 (e.g. $F(x, y) = c(x - y)^2$ was used with $x = \hat{L}_S(h)$ and $y = L(h)$), but the theorem is valid e.g. for the 01 loss or for
49 the ramp loss, which are non-convex. We are glad for the question, we will emphasize this in the revised paper. Also,
50 the notation “KL” (upper-case) is the for KL divergence in the generic case (between any pair of distributions), while
51 “kl” (lower-case) is reserved for Bernoulli distributions. We think this notation is helpful as a visual aid, which has been
52 used by other authors as well, but of course we can insert reminders throughout the paper to help the reader.