

1 We thank all reviewers for their time and expertise. We are particularly happy that R1 and R3 are positive  
2 about our paper. R1’s summary of our contributions resonates strongly with what we wanted to achieve.

3 We are worried that R2 may have misunderstood main points of the paper. Some of their comments seem to  
4 suggest that we endorse the EF. The exact opposite is the case. The EF, which doesn’t exist in the Stats literature,  
5 is widely used in ML (see our citations). We argue that it should *not* be used. Let us clarify some of your points:

6 > (R2) “for applying the exact natural gradient [...] there is no reason to use the EF over the Fisher.”

7 We agree! However, a large portion of the literature applying natural gradient ideas to machine learning uses the EF and  
8 many researchers are under the impression that the EF and the Fisher are interchangeable. We cite 15 papers that use  
9 the EF, 9 of them from the past 3 years, and we know of at least 4 additional papers, accepted or under review at ICML,  
10 UAI and NeurIPS 2019, perpetuating this confusion. Basically, our paper makes the points that you levy against us in  
11 the review. Our goal is to point out the differences between the EF and the Fisher and clear the confusion, in the hope  
12 that future applications of natural gradient methods to machine learning will be more impactful.

13 > (R2) “the authors used the term “empirical Fisher” [... did the authors coin] the term?”

14 We did not coin this term and, as noted in L265-267, we believe it is a misnomer. To the best of our knowledge, the  
15 term “empirical Fisher” first appears in the manuscript of Martens [2014]. Although the origin of the concept might be  
16 Le Roux et al. [2007], who confused the covariance and the Fisher information, and Graves [2011] who proposed it as  
17 an approximation of the Hessian. The idea was heavily popularized by the Adam optimizer [Kingma and Ba, 2015].

18 > (R2) “This should be connected to observed Fisher information which is a well-defined concept in statistics.”

19 This is a fantastic point that we definitely need to address. While this terminology is widely used in the literature we  
20 cite, it is true that the concept might sound strange to statisticians as the terms “observed Fisher” and, even more  
21 confusing, “empirical Fisher information”, are used, in statistics, to refer to what the community we are trying to reach  
22 calls the Fisher information. We will use the additional page to make the context clearer to a larger audience and give  
23 an overview of how the EF approximation came to be.

24 > (R2) “the experiments are mainly performed on toy examples. It is not clear how [to generalize] to real networks”

25 Note again that we are arguing *against* the use of the EF concept. If it doesn’t even work on toy problems, why rely on  
26 it in “real networks” that are much harder to understand?

### 27 Questions raised by R3 & R1 (thanks for these!):

28 > (R3) “the plots in fig 3 show that the cosine [...] is close to 1 near the end of the training, which seem to contradict  
29 the claims in subsection 3.1 and 3.2: in practice it seems that the cases discussed in 3.1 and 3.2 are not met.”

30 The plots for §3.1 and 3.2 are in Fig.2, showing strong differences at the minimum. Fig. 3 showcases §3.3, about  
31 preconditioning, showing that updates can be opposite. While the cosine gets close to 1 in Fig.3, this metric captures  
32 only a small part of the relationship and the two matrices can still be very different (see figure at end of this rebuttal).

33 > (R3) “how do you obtain eq. 12? [...] line 184 the Fisher does not match eq. 2” Eq.12 and L184 show the equations  
34 for the empirical Fisher (Eq.3) and the Fisher (Eq.2) applied to the linear regression example of L178-179.

35 > (R3) “line 192 if  $b$  is arbitrarily close to 1 then  $\nabla_b \log b$  is arbitrarily close to 0 [...] I think you wanted to write that  
36 using the definition line 179 with a variance 1 normal distribution,  $b_n(\theta)$  can not get arbitrarily close to 1.” and

37 > (R1) “For classification [when] all labels are predicted correctly with probability 1, wouldn’t  $[F = EF = 0]$ ?”

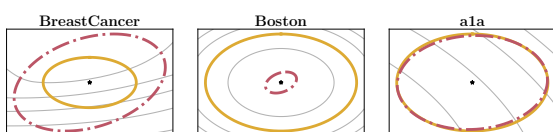
38 The subtlety of the proposed change to the definition of the GGN (Prop.3) relies on this key point: for the guarantee in  
39 Prop.2 to hold,  $\nabla_b a(b)$  needs to go to 0, which can not happen for  $a(b) = \log b$  as  $\nabla_b \log b = 1/b$  (for  $b \in [0, 1]$ ). As  
40 both reviewers noted, if  $b$  goes to 1 and is correct, the EF is 0 (the gradients are 0), the Fisher is 0 (the model is no  
41 longer probabilistic, meaning new samples can not give more information) and the Hessian is also 0 (some parameters,  
42 the weights for classification or the precision—if learned—for regression, need to reach  $\infty$  and the landscape is infinitely  
43 flat). We will add a paragraph to explain this special case.

44 > (R3) “in eq 15 I think the multiplicative  $N$  should be in front of the covariance matrix if we follow [...] eq 3”

45 We defined  $\Sigma$  as the covariance a uniformly distributed gradient scaled by  $N$  (L184) to be consistent with the definition  
46 of the loss as a sum in Eq.1. Sorry for the confusion; we will try to make it clearer.

47 **Typos found by R3: Will all be addressed, sincere thanks for the close reading.** (Eq.5, L.183-185 and Prop.2)

48 Sorry about the lost pointer to the proof of Prop.2 - it is in the supplementary §C.3.



49 Quadratic approximations (as in Fig.2) using the projection of F (yellow) and EF (red) on the two largest eigenvectors of F, at the end of training with the EF, using the settings of Fig.3. While the cosines are close to 1, the matrices can still be very different in terms of directions and scaling.