

1 First of all, we would like to thank all reviewers for their suggestions to improve our paper submission.

2 Reviewers **#1** and **#2** suggest experiments to measure if UMAL yields calibrated outputs. We have performed an

3 additional empirical study to assess this point. We highlight that our system predicts an output distribution $p(y|x, \mathbf{w})$

4 (not a confidence value). In particular, we have computed the % of actual test data that falls into different thresholds of

5 predicted probability. Ideally, given a certain threshold $\theta \in [0, 1]$, the amount of data points with a predicted probability

6 above or equal to $1 - \theta$ should be similar to θ . In the left part of Figure 1 we plot these measures for different methods (in

7 green, our model) when considering the BCN RPF dataset. Furthermore, (following reviewers **#1** and **#2** suggestions)

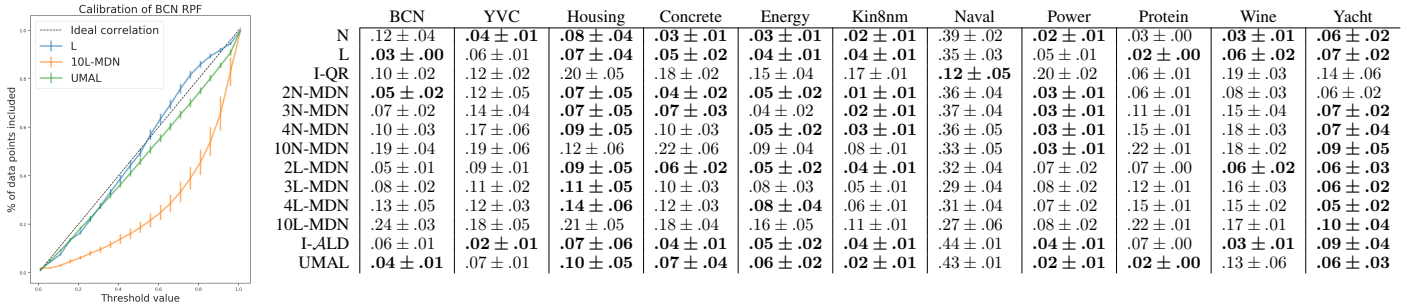
8 we have evaluated calibration quality on the UCI datasets with the same architectures as [2,3]. In the right part of Figure

9 1 we report the mean absolute error between the empirical measures and the ideal ones for all datasets. As it is shown,

10 the conditional distribution predicted by UMAL has low error values. **Therefore, we can state that UMAL produces**

11 **proper and calibrated conditional distributions that are especially suitable for heterogeneous problems.**

Figure 1: Plot with the performance of three different models in terms of calibration. The mean and standard deviation for all folds of the mean absolute error between the predicted calibration and the perfect ideal calibration is represented in the table.



12 For the sake of completeness, we also have computed UMAL negative log-likelihood for UCI datasets (see Table 1)

13 following [2]. These results restate that UMAL is always in the best positions. However, it should be noted that most of

14 these databases have a **small sample size** and that in this regime, aleatoric uncertainty cannot be reliably estimated. We

15 hypothesize that a better solution would be to simultaneously estimate epistemic (as [1,2,3]) and aleatoric uncertainty.

16 Minor comments:

17 - Quantile Regression allows us to approximate a desired quantile, unlike the classical regression that only estimates the

18 mean or the median. This is useful since we could capture confidence intervals without strong assumptions about the

19 distribution function to approximate (**Rev. #3. Q5.1**).

20 - IQR is a function, $\phi(\mathbf{x}, \tau)$, that can be evaluated for any real value of $\tau \in (0, 1)$ to give us a point-wise estimation

21 of an infinite number of quantiles $(\mu_\tau)_{\tau \in (0, 1)}$ (**Rev. #3. Q5.3**). On the other hand, I-ALD is a function, $\phi(\mathbf{x}, \tau)$, that

22 predicts the $(\mu_\tau, b_\tau)_{\tau \in (0, 1)}$ parameters of each ALD that is identified with a real asymmetry value, τ . Thus, now each

23 ALD tries to estimate, in a non-point-wise manner, their corresponding quantile (**Rev. #3.Q5.2**).

24 - If we consider each ALD as a component of a single mixture model we arrive to UMAL and, in turn, we solve the

25 crossing quantiles problem because all ALDs are optimized jointly to produce the output distribution (**Rev. #3. Q5.4**).

26 - To obtain the distribution predicted by UMAL we will evaluate the learnt function $\psi(\mathbf{x}, \tau)$ for each τ in any

27 discretization of its interval of definition, $(0, 1)$. For instance, by considering the partition $sel_\tau = [0.01, 0.02, \dots, 0.99]$

28 as parameter of the function defined in Algorithm 3 of the paper (**Rev. #1. Q2.2**).

29 We are seriously considering the suggestions made by reviewers **#1** and **#2** regarding typos and polishing. Thus, the

30 paper is being revised by several native English speakers to improve its flow.

Table 1: Comparison of the Negative Log-Likelihood of the test set over different train-test folds proposed in [3].

	Housing	Concrete	Energy	Kin8nm	Naval	Power	Protein	Wine	Yacht
Normal distribution	2.76 ± .34	3.20 ± .16	2.13 ± .24	-1.15 ± .03	-3.67 ± .01	2.83 ± .03	2.84 ± .03	1.05 ± .14	1.86 ± .31
Laplace distribution	2.59 ± .20	3.21 ± .13	2.06 ± .20	-1.08 ± .04	-3.73 ± .04	2.87 ± .03	2.74 ± .01	1.00 ± .08	1.54 ± .37
Independent QR	10.96 ± 2.4	10.19 ± .95	9.45 ± 1.3	9.22 ± .66	5.14 ± .89	8.39 ± .45	8.14 ± .52	12.30 ± .91	10.32 ± 2.9
2 comp. Normal MDN	2.74 ± .30	3.25 ± .21	2.02 ± .30	-1.15 ± .05	-3.66 ± .02	2.85 ± .05	2.56 ± .03	1.33 ± .61	1.55 ± .32
3 comp. Normal MDN	2.68 ± .28	3.64 ± .28	2.30 ± .43	-1.15 ± .05	-3.66 ± .01	2.85 ± .04	2.90 ± .15	0.69 ± 1.0	1.54 ± .52
4 comp. Normal MDN	2.87 ± .46	3.74 ± .28	2.46 ± .39	-1.12 ± .04	-3.66 ± .03	2.86 ± .05	3.32 ± .11	0.52 ± .90	1.43 ± .36
10 comp. Normal MDN	3.10 ± .46	5.64 ± 1.1	3.03 ± .71	-0.99 ± .06	-3.64 ± .03	2.86 ± .04	4.94 ± .75	0.75 ± .95	1.75 ± .49
2 comp. Laplace MDN	2.61 ± .23	3.28 ± .14	2.06 ± .30	-1.10 ± .04	-3.70 ± .06	2.91 ± .05	2.50 ± .03	0.59 ± .63	1.37 ± .42
3 comp. Laplace MDN	2.65 ± .25	3.45 ± .16	2.30 ± .21	-1.09 ± .03	-3.66 ± .06	2.95 ± .04	2.65 ± .06	-0.81 ± .70	1.39 ± .35
4 comp. Laplace MDN	2.76 ± .42	3.57 ± .14	2.31 ± .35	-1.10 ± .05	-3.68 ± .06	2.93 ± .04	2.79 ± .08	-0.65 ± .96	1.45 ± .35
10 comp. Laplace MDN	3.17 ± .46	3.95 ± .34	2.80 ± .49	-0.98 ± .07	-3.62 ± .10	2.96 ± .05	3.46 ± .12	0.52 ± .74	1.63 ± .34
Independent ALD	2.79 ± .56	3.87 ± .12	2.28 ± .11	-1.00 ± .05	-2.82 ± .01	2.89 ± .02	2.68 ± .01	1.01 ± .07	1.78 ± .41
UMAL model	2.59 ± .26	3.74 ± .15	2.13 ± .14	-1.09 ± .03	-2.81 ± .01	2.85 ± .03	2.40 ± .01	0.14 ± .70	1.41 ± .38

31 [1] Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. ICML, 2016.

32 [2] Lakshminarayanan, Balaji, et al. Simple and scalable predictive uncertainty estimation using deep ensembles. NIPS, 2017.

33 [3] J. M. Hernández-Lobato, et al. Probabilistic backpropagation for scalable learning of Bayesian neural networks. ICML, 2015.