Uncertainty Quantification Metrics for Deep Regression

Ziliang Xiong, Simon Kristofferson Lind, Per-Erik Forssén, Volker Krüger



LUND UNIVERSITY



Research Goals

Here we address metrics for uncertainty quantification with focus on regression tasks. We try to answer these questions:

- What do different evaluation metrics measure? lacksquare
- Are they stable and robust with limited data?

• What are their strengths and weaknesses?

Our results indicate that Calibration Error is the most stable and interpretable metric, but AUSE and NLL also have their respective use cases. We discourage the usage of Spearman's Rank Correlation and recommend replacing it with AUSE.



MSE

Uncertainty Evaluation Metrics

Area Under Sparsification Error Curve

It assesses how well the predicted uncertainty coincides with the prediction error on a test set.



Calibration Error

The difference between cumulative density function and empirical frequency.



Selected Results

Stability on varying test set sizes

Experiment 1: How quickly does each metric converge to its expected value? Experiment 2: Are the estimates unbiased for small test set sizes?



Conclusions



- 1. Stability: CE > AUSE > NLL > Spearman
- 2. Spearman converges to zero, should not be used with dense samples
- 3. All metrics converge beyond a dataset size of 1024
- 4. No metric exhibits any meaningful bias beyond a dataset size of 2⁶



Conclusions

1. Homo and Hetero: Perfect Predicted distribution returns good CE, NLL but not good ranking.



Two Typical Uncertainty-aware Regression Models



- 2. Multimodal: AUSE for DE model is good, but calibration and NLL suggest that the model
 - fails to capture generating distribution.
- Perfect Calibration is not equaivalent to perfect Ranking. \Rightarrow

Interpretability

- **CE:** lower bound of 0, and is highly interpretable; It requires the least amount of samples to be stable
- **AUSE:**unbounded and thus it lacks interpretability; robust ranking metrics, • suitable for practical tasks
- **NLL**:less stable with small test set
- **Spearman**: Unsuitable for dense samples •

