Multi-Scale Evaluation of Uncertainty Quantification Techniques for Deep Learning based MRI Segmentation

Benjamin Lambert^{1,3}, Florence Forbes², Senan Doyle³, Alan Tucholka³, Michel Dojat¹

¹Univ. Grenoble Alpes, Inserm U1216, Grenoble Institut Neurosciences, GIN, 38000, Grenoble, France ²Univ. Grenoble Alpes, Inria, CNRS, Grenoble INP, LJK, 38000, Grenoble, France ³Pixyl, Research and Development Laboratory, 39000, Grenoble, France

Abstract

Deep learning systems are powerful models that tend to be considered as black-boxes. To facilitate their acceptance in clinical routine, models that *know when they don't know* are desired. A multitude of methods have been proposed to quantify uncertainty of Deep Learning predictions. Here, we propose to compare three of them : Monte Carlo Dropout, Deep Ensemble, and Heteroscedastic models. We illustrate this analysis on a task of segmentation of White-Matter Hyperintensities in FLAIR MRI sequences of Multiple-Sclerosis patients. Evaluation is carried-out at 3 different scales : the voxel, the lesion, and the whole image. Results show the superiority of the Heteroscedastic approach, outperforming competing methods at the voxel and imagelevel, and ranking second in the lesion task. This technique is also the fastest and the most memory-efficient. Alternative methods achieved lower uncertainty estimates quality, while also requesting longer inference time or larger computational requirements.

Introduction

Deep Learning (DL) techniques have become the gold standard for biomedical image segmentation, but they tend to be considered as black-boxes. This is partly due to the inability of neural networks to express the uncertainty in their predictions. In recent years, massive efforts have been carried out to develop models that can express their confidence¹.

In this work, we propose an in-depth evaluation of 3 state-ofthe-art approaches to quantify uncertainty attached to DL predictions : Monte Carlo Dropout² (M1), Deep Ensemble³ (M2), and Heteroscedastic network⁴ (M3). Uncertainty estimates are evaluated at the voxel, lesion and image levels. We illustrate this comparison on an automatic segmentation task to detect White-Matter Hyperintensities (WMH) from T2-weighted FLAIR MRI sequences of Multiple-Sclerosis (MS) patients.



MC-Dropout

Figure 1 : Illustration of segmentation masks and uncertainty maps obtained with each technique.



lesion probability, u is the uncertainty attached to the prediction.

Experiments

We used a brain dataset composed of 238 T2-weighted FLAIR MRI sequences of MS patients, with ground truth segmentations of WMH. The dataset was split into 187 scans for training and 51 for testing. Example of segmentation masks and uncertainty maps are illustrated in Figure 1.

At the voxel-level, we assessed the quality of uncertainty with Area Under Confidence-Classification estimates Characteristic curves (AUCCC)⁵ (Figure 3). This metric is agnostic of the model segmentation performance, which is desired for a fair comparison. An AUCCC of 0.5 indicates that the uncertainties of correct and incorrect voxels are confounded, hence meaningless. Alternatively, an AUCCC of 1 indicates that correct voxels are systematically assigned with a lower uncertainty than incorrect voxels.

Image uncertainty is calculated as the mean of uncertainties of all voxels predicted as lesions. We plotted correlation curves between images uncertainties and their Dice scores (segmentation quality) (Figure 5). We used the Pearson correlation coefficient (PCC) between both quantities to assess the quality of image uncertainties, and the Dice scores to estimate segmentation performance.





Results and Conclusion

Results are presented in Figure 6, with top-scoring techniques highlighted in green. The Heteroscedastic approach (M3) outperformed competing approaches on 2 of the 3 evaluation scales regarding uncertainty estimates as well as for segmentation performance. This technique is also the fastest, as uncertainty and segmentation are simultaneously obtained in a single-step, contrary to MC-Dropout (M1) and Deep Ensemble (M2) that require the aggregation of multiple predictions.

Method Scale & Metrics	MC-Dropout (M1)	Deep Ensemble (M2)	Heteroscedastic (M3)
Voxel - AUCCC (E1)	0.731	0.705	0.745
Lesion - AULSC (E2)	0.791	0.697	0.766
Image - PCC (E3)	-0.685	-0.497	-0.826
Image - Dice (E3)	0.777	0.780	0.784

1.	Abd
	learı
	(202
2.	Yariı
	appr
	Proc
	Cont
3.	Laks
	prec
	of th
	Syst
4.	Rich
	Refi
	with
5.	Hua
	unce
	arXiv
6.	Roy
	qual
	Cont



Figure 5 : Correlation curves between segmentation performance and uncertainty estimates.

Figure 6 : Evaluation results of the presented methods. Top-performing method is highlighted in green.

References

lar, Moloud, et al. A review of uncertainty quantification in deep rning: Techniques, applications and challenges. Information Fusion 21):243-297.

in Gal and Zoubin Ghahramani. 2016. Dropout as a Bayesian roximation: representing model uncertainty in deep learning. In ceedings of the 33rd International Conference on International ference on Machine Learning - Volume 48 (ICML'16). 1050–1059.

shminarayanan B, Pritzel A, Blundell C. Simple and scalable dictive uncertainty estimation using deep ensembles. In Proceedings he 31st International Conference on Neural Information Processing tems (NIPS'17): 6402-6413.

hard McKinley, Michael Rebsamen, et al. Uncertainty-Driven inement of Tumor-Core Segmentation Using 3D-to-2D Networks Label Uncertainty. BrainLes at MICCAI (1) 2020: 401-411.

ang X, Yang J, Li L, Deng H, Ni B, Xu Y. Evaluating and boosting ertainty quantification in classification. arXiv preprint v:1909.06030.2019.

AG, Conjeti S, Navab N, Wachinger C. Inherent brain segmentation lity control from fully convnet monte carlo sampling. In International ference on Medical Image Computing and Computer-Assisted Intervention 2018 Sep 16 (pp. 664-672). Springer, Cham.