

Let's agree to disagree: learning highly debatable multirater labelling.

Carole H. Sudre^{1,2,3}, Beatriz Gomez Anson⁴, Silvia Ingala⁵, Chris D Lane², Daniel Jimenez², Lukas Haider⁶, Thomas Varsavsky^{1,3}, Ryutaro Tanno³, Lorna Smith⁷, Sébastien Ourselin¹, Rolf H Jäger⁸, and M. Jorge Cardoso^{1,2,3}

¹ School of Biomedical Engineering and Imaging Sciences, KCL, UK

² Dementia Research Centre, UCL Institute of Neurology, UK

³ Department of Medical Physics and Biomedical Engineering, UCL, UK

⁴ Santa Creu i Sant Pau Hospital, Universitat Autònoma de Barcelona, Spain

⁵ Vrije University Medical Centre Amsterdam, The Netherlands

⁶ Queen Square Multiple Sclerosis Centre, UCL Institute of Neurology, London, UK

⁷ Cardiometabolic Phenotyping Group, Institute of Cardiovascular Science, UCL, UK

⁸ Brain Repair and Rehabilitation Group, Institute of Neurology, UCL, London, UK

Abstract. Classification and differentiation of small pathological objects may greatly vary among human raters due to differences in training, expertise and their consistency over time. In a radiological setting, objects commonly have high within-class appearance variability whilst sharing certain characteristics across different classes, making their distinction even more difficult. As an example, markers of cerebral small vessel disease, such as enlarged perivascular spaces (EPVS) and lacunes, can be very varied in their appearance while exhibiting high inter-class similarity, making this task highly challenging for human raters. In this work, we investigate joint models of individual rater behaviour and multirater consensus in a deep learning setting, and apply it to a brain lesion object-detection task. Results show that jointly modelling both individual and consensus estimates leads to significant improvements in performance when compared to directly predicting consensus labels, while also allowing the characterization of human-rater consistency.

Keywords: Deep learning · Noisy labels · Classification

1 Introduction

Detection and differentiation between types of pathological objects is a core problem of medical image analysis, generally requiring costly expert labelling. Disagreement between raters can be a result of differences in radiological training schools, rater competence, and sample appearance, among others. This problem is often exacerbated by changes in rater performance caused by retraining or observational bias.

Due to the variability in shape and intensity signatures observed across the full spectrum of lesions, even the most trained raters can present a high inter-rater variability. In such cases, finding a majority voting consensus classification is the most common strategy.

When classifying objects into multiple classes, it is often more complex to separate all object types directly, than it is to first detect all pathological objects followed by their classification, as some class decision boundaries are easier than others. This sequential detection/classification problem is, for instance, present in the context of age-related vascular changes in which macroscopic alterations can be observed on structural MR images. Among these observed changes, small elements such as enlarged perivascular spaces (EPVS) and lacunes are observed on similar image sequences [9]. EPVS, often associated with concomitant neuropathology and deleterious clinical outcome [5], appear as fluid filled structures with a linear shape. However, because of their limited size ($<10\text{mm}^3$) and highly variable appearance, EPVS instances are often confused with other concomitant lesions, such as lacunes. Because of this intrinsic uncertainty, the labelling of these small vascular lesions can be seen as a four-class problem, with classes ‘EPVS’, ‘Lacune’, ‘Undetermined’, and ‘Nothing’. This classification problem suffers from two concomitant issues: a) class imbalance with a 100:1 ratio between EPVS and lacunes, and b) noisy labelling as a result of rater disagreement. Consensus labels are also problematic in this setting, as rater behaviour is non-random and samples are not truly independent.

In this work, we build on a previously described 3-dimensional multirater Regional Convolutional Neural Network (RCNN) model, used here as a lacune and EPVS object detection system. However, rather than only learning the consensus value or a single rater, we propose to jointly learn the consensus majority voting, the associated probability for each class, and each individual rater decision, so as to appropriately model highly debatable label predictions.

2 Related work and problem specificities

Many of the recent publications on classification with noisy labels assume independence between samples and noise, and a constant mislabelling probability [4], which does not hold in the case of difficulty induced variability and rater shift. Strategies for classification in the presence of noise include sample re-weighting (importance reweighting) or curriculum-based sample selection [1, 3]. Other approaches, normally classed as label/classifier fusion, disentangle rater and label uncertainty either by iteratively favouring raters that agree with the consensus [10], or by reducing sample correlation to construct a balanced classifier [2]. Lastly, the relationships between rater behaviour can also be learned through their confusion matrices [7]. Notably, most of these works focus on the problem of classification, where balanced sampling strategies can be employed, something that is not possible in a joint detection/classification model. In this work, we argue that combining majority voting predictions while learning individual rater behaviour allows for a better model of rater consistency and sample uncertainty.

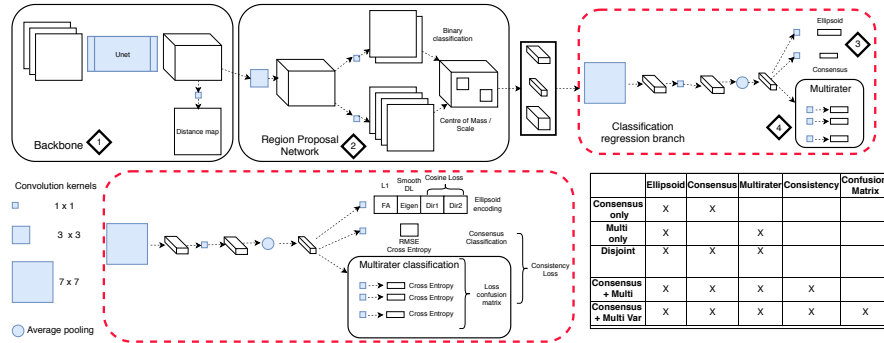


Fig. 1. Detection and multirater classification architecture framework. In this work focus is put on the classification branch (red dashed) further detailed on the second row along with the description of the different training regimes.

3 Methods

3.1 Network architecture

The multirater 3D RCNN framework presented in [6] is composed of four stages: 1) a backbone network learning the features using as target a distance map to the objects of interest; 2) a region proposal network (RPN) regressing the location of candidate centres of mass of target objects together with their spatial scale; 3) patches from the RPN representation are fed into a two layer network in order to regress the average object classification and object shape; 4) a multibranch fully connected layer is used to model the behaviour of each rater. (see Fig. 1.)

Shape encoding Instead of modelling each object by its encompassing cuboid [6], the shape of each candidate object is encoded as a four-parameter simplified encompassing ellipsoid, namely using the largest eigenvalue, the two first components of the associated eigenvector and the value of fractional anisotropy of the associated tensor.

3.2 Multi-rater classification

The classification of the candidate objects is defined as a four-class problem, i.e. EPVS, Lacune, Undetermined, Nothing. From a human-rater point of view, the classification can be seen in two different ways: a multi-rater consensus, here modeled as the probability of a class to be chosen among the six raters computed as the average rating, and a rater-specific categorical label.

Consensus average classification When modelling the consensus/average rater, the training can be performed using either a hard or a probabilistic classification, or a combination of both. Note, here, the hard classification corresponds simply to the majority voting categorical consensus, while the continuous probability

encodes the uncertainty over the final classification. As a consequence, a cross-entropy loss is used to learn the consensus, while a root mean square error loss is used over the resulting class probabilities.

Independent rater modelling In the last stage of classification, a cross-entropy loss is used to learn each rater label independently. Inter-rater behaviour can be enforced through a variability loss (L_{var}) penalizing the difference between the effective and predicted probabilistic confusion matrices. Noting C (resp. \widehat{C}) the observed (resp. predicted) confusion matrix, $L_{var} = \sum_{(i,j)} |C_{i,j} - \widehat{C}_{i,j}|$

Ideally, we would also like to have consistency between the predicted group consensus and the consensus of individual prediction. In order to achieve this, the following consistency loss L_{cons} is introduced:

$$L_{cons} = \sqrt{\sum_{k=1}^K \widehat{p}_k - \frac{1}{R} \sum_{r=1}^R \widehat{p}_{kr}}$$

with \widehat{p}_k denoting the predicted consensus probability, and \widehat{p}_{kr} denoting the predicted probability given by rater r for class k .

Compensating for inter-rater variability and enhancing individual rater characteristics The EPVS labelling problem is highly variable in terms of rater agreement; sometimes all raters agree with each other, while other times raters converge to completely different decisions. As a consequence, when predicting the group consensus, we have enforced consensus learning from samples of high agreement. To this effect, sample importances were downweighted according to their observed variability, here expressed as $var = 1 - \sum_{k=1}^K p_k^2$, where p_k is the observed classification probability for class k . The sample is then weighted by $\exp(-var)$. Conversely, when modelling individual raters, and in order to learn rater-specific behaviours, we promote samples for which the individual rater disagrees with the consensus. This is achieved by weighting each rater-sample combination by the inverse of its contribution to the consensus ($1/p_{kr}$), where p_{kr} is the observed probability for the sample to be classified as k if rater r labels it as k .

4 Data and experiments

4.1 Data

16 subjects that were part of a large tri-ethnic cohort investigating the relationship between cardiovascular risk factors and brain health [8] were chosen due to their elevated vascular burden. 4147 EPVS and lacunes were manually segmented using jointly 1mm3 structural MR sequences (T1, T2, FLAIR) using ITKSnap⁹. Individual segmented lesions, defined using connected components,

⁹ <http://www.itksnap.org/pmwiki/pmwiki.php?n=Main.HomePage>

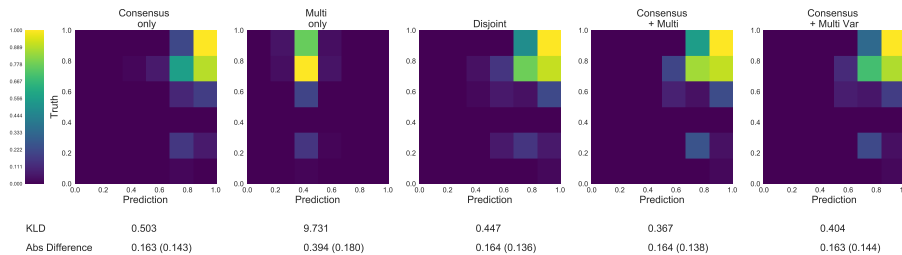


Fig. 2. Comparison of EPVS probability distributions quantitatively evaluated in terms of KLD and absolute error (mean sd).

were then classified by six trained raters using an in house dedicated viewer. Only objects bigger than 5 voxels were used in this study, resulting in a database of 2202 elements. 14 subjects were used for training and two subjects for testing. The test set contained 184 objects that were all classified at least by one rater as EPVS. Inter-rater accuracy ranged from 0.47 to 0.92 with a mean of 0.72.

4.2 Training modes

In order to investigate the model’s ability to handle label noise, different training regimes were adopted (see Fig. 1): 1) Training only the shape + consensus classification (Consensus only); 2) Staged training of the shape encoding followed by the independent rater multihead (Multi Only); 3) Staged training of shape and classification, followed by training the independent rater multihead (Disjoint); 4) Staged training of ‘shape and consensus only’, followed by ‘multihead only’ finishing by ‘shape, consensus and multihead’ with consistency loss (Consensus + Multi); 5) Training as in 4, with an extra loss over the confusion matrix L_{var} . All models were trained for 10000 iterations with a learning rate of 10^{-4} and using the Adam optimiser.

5 Experiments and Results

5.1 Consensus probability

As a first experiment, we investigate the ability of each training mode to appropriately predict the distribution of EPVS classification probabilities. Figure 2 presents the joint histograms of the target and predicted distributions. The resulting Kullback-Leibler Divergence (KLD) over the distributions are displayed below along with the mean absolute error in prediction. Results show that all methods explicitly learning the average consensus are able to reproduce it well. The ability of the different models to reproduce individual rater behaviour was evaluated by comparing predicted inter-rater agreement with observed inter-rater agreement. Figure 3 presents the pairwise agreement results between observations and between predictions, and measures of correlation and

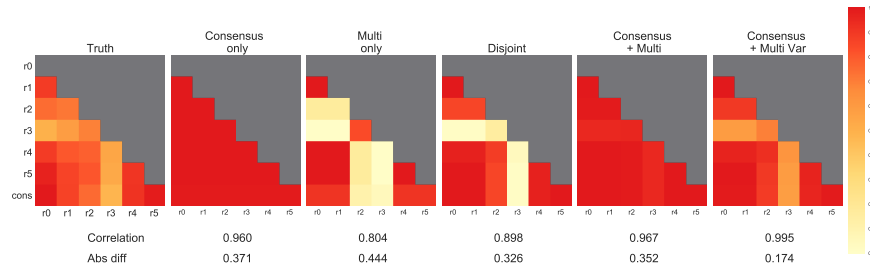


Fig. 3. Pairwise agreement scores between observed rater labels compared to the agreement scores between categorized individual rater predictions for each training mode. The left most element is the target inter-rater behaviour. Pearson correlation coefficient and mean absolute difference against the target inter-rater behaviour are presented below.

absolute difference between agreement matrices. One can note that, as expected, no inter-rater behaviour is learnt when adopting the ‘Consensus only’ framework. Furthermore, we observe that the inter-rater agreement learnt with both the ‘Multi only’ and the ‘Disjoint’ model is exacerbated compared to the truth. This rater behaviour exacerbation fades away when enforcing consistency between multi-rater consensus predictions and the consensus of individual rater predictions (i.e. ‘Consensus+Multi’ model).

5.2 Consistency between consensus and multirater average

This experiment aims to test the efficacy of the loss function introduced in Section 3.2 with the aim of promoting the agreement between the multi-rater consensus labelling and the consensus of individual predictions. Figure 4 left presents the boxplots of the difference between the average of the predicted

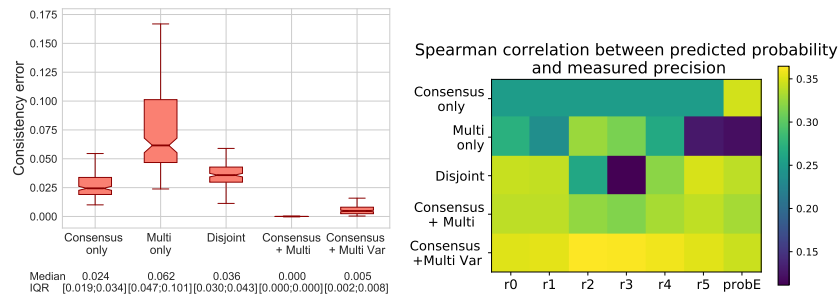


Fig. 4. Left: Boxplot of the consistency error between predicted consensus probability and average of individual raters predictions. Right: Spearman coefficient between predicted probability of classifying the element as an EPVS and measured precision ($1/var$) over multiple models. ProbE refers to the predicted probability of an EPVS for the consensus of raters.

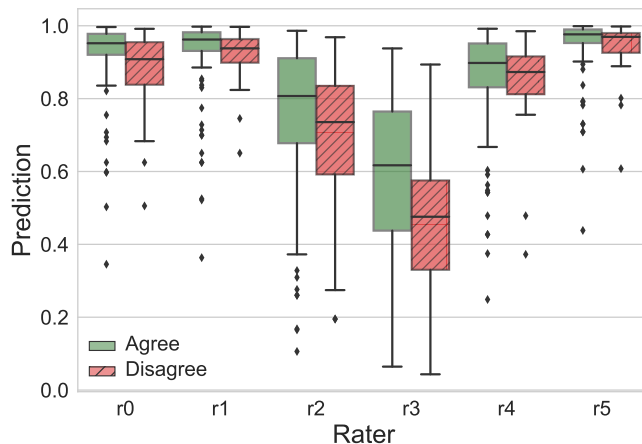


Fig. 5. Boxplot of the predicted probabilities for the agreeing and disagreeing cases for each rater with the best training regime (Consensus + Multi Var).

individual raters and the consensus prediction. Numerical results of median and interquartile range are presented below the graph. Training regimes that promote consistency between the consensus prediction and the average of independent predictions both reach, as expected, a very high level of consistency. Conversely, simpler models only optimising for independent rater predictions do not achieve a good consensus estimation.

5.3 Variability, disagreements and individual rater quality

In this experiment, we would like to assess if the probabilistic predictions of each individual rater provide a good proxy for sample uncertainty (defined as the variability of individual ratings). To this end, we estimate the Spearman correlation coefficient between each individual prediction and the measured precision defined as $1/var$, displayed on Figure 4 (right). As already noted from Figure 3, no rater-specific information can be modeled using only the consensus. Individualized rater predictions were found to be strongly associated with overall variability, primarily when consistency losses were applied.

5.4 Labelling introspection

We used the best overall model ('Consensus+Multi Var') to study the distribution of probabilistic predictions for objects whose rater classification was in agreement with the consensus versus objects where there was rater disagreement, plotted in Figure 5. High prediction probabilities for individual raters were found to be a good surrogate marker rater agreement and rater consistency

over training samples. Results suggest that rater 3, and to a lesser degree rater 2, displayed inconsistent labelling behaviour. Interestingly, when asked about their rating practice, both raters 2 and 3 indicated having undergone clinical retraining, possibly explaining the observed shift in their labelling. Retraining the model without these two raters resulted in an improvement in the consensus prediction, with a KLD reduced to 0.202 and a mean absolute error over the predicted probability of 0.15. This experiment suggests that one can use the proposed framework to identify not only inter-rater disagreement but also intra-rater inconsistency, and potentially correct for it.

6 Discussion and conclusion

In this work, we investigated different training regimes in presence of noisy labelling with the aim of predicting both rater consensus and individualized predictions. We found that promoting agreement between predicted multi-rater consensus and the consensus of individualized predictions can provide good model accuracy together with the ability to introspect rater behaviour, thus not only allowing the identification of noisy labels/subjects but also assessing rater skills so as to prevent bias in large scale studies and enforce appropriate radiological training. Future work will explore the use of this information in an active learning setting and develop the accuracy of the multi-rater model estimates.

Acknowledgments We are extremely grateful to all the participants of the SABRE study, and past and present members of the SABRE team. This work was supported by an Alzheimer’s Society Junior Fellowship (AS-JF-17-011), the Wellcome/EPSRC Centre for Medical Engineering [WT 203148/Z/16/Z], IMI2 grant AMYPAD [115952], the MSCA-ITN-Demo [721820], and the Wellcome Flagship Programme in High-Dimensional Neurology. The SABRE study was funded at baseline by the Medical Research Council, Diabetes UK, and the British Heart Foundation. At follow-up, the study was funded by the Wellcome Trust (067100, 37055891 and 086676/7/08/Z), the British Heart Foundation (PG/06/145, PG/08/103/ 26133, PG/12/29/29497 and CS/13/1/30327) and Diabetes UK (13/0004774). We gratefully acknowledge NVIDIA corporation for the donation of a GPU Tesla K40 that was used in the preparation of this work.

References

1. Bouguelia, M.R., Nowaczyk, S., Santosh, K.C., Verikas, A.: Agreeing to disagree: active learning with noisy labels without crowdsourcing. *International Journal of Machine Learning and Cybernetics* **9**(8), 1307–1319 (aug 2018)
2. Hongzhi Wang, Suh, J.W., Das, S.R., Pluta, J.B., Craige, C., Yushkevich, P.A.: Multi-Atlas Segmentation with Joint Label Fusion. *IEEE TPAMI* **35**(3), 611–623 (mar 2013). <https://doi.org/10.1109/TPAMI.2012.143>
3. Jiang, L., Zhou, Z., Leung, T., Li, L.J., Fei-Fei, L.: MentorNet: Learning Data-Driven Curriculum for Very Deep Neural Networks on Corrupted Labels (dec 2017), <http://arxiv.org/abs/1712.05055>

4. Li, Y., Yang, J., Song, Y., Cao, L., Luo, J., Li, L.J.: Learning From Noisy Labels With Distillation (2017)
5. Ramirez, J., Berezuk, C., McNeely, A.A., Gao, F., McLaurin, J., Black, S.E.: Imaging the Perivascular Space as a Potential Biomarker of Neurovascular and Neurodegenerative Diseases. *Cellular and molecular neurobiology* (mar 2016)
6. Sudre, C., Gomez Anson, B., Ingala, S., Lane, C., Jimenez, D., Haider, L., Varsavsky, T., Smith, L., Ourselin, S., Jäger, R., Cardoso, M.: 3d multirater rnn for multimodal multiclass detection and characterisation of extremely small objects. In: Proc. of the 2nd Int. MIDL conf. Proceedings of Machine Learning Research, vol. 102, pp. 447–456. PMLR, London, United Kingdom (08–10 Jul 2019)
7. Tanno, R., Saeedi, A., Sankaranarayanan, S., Alexander, D.C., Silberman, N.: Learning from noisy labels by regularized estimation of annotator confusion. In: Conference on Computer Vision and Pattern Recognition (2019)
8. Tillin, T., Forouhi, N.G., McKeigue, P.M., group Chatuverdi, N.f.t.S., Chaturvedi, N.: Southall And Brent REvisited: cohort profile of SABRE, a UK population-based comparison of cardiovascular disease and diabetes in people of European, Indian Asian and African Caribbean origins. *International Journal of Epidemiology* **41**(1), 33–42 (feb 2012). <https://doi.org/10.1093/ije/dyq175>
9. Wardlaw, J.M., Smith, E.E., Biessels, G.J., Cordonnier, C., Fazekas, F., Frayne, R., Lindley, R.I., O’Brien, J.T., Barkhof, F., Benavente, O.R., Black, S.E., Brayne, C., Breteler, M.M.B., Chabriat, H., DeCarli, C., de Leeuw, F.E., Doubal, F., Duering, M., Fox, N.C., Greenberg, S., Hachinski, V., Kilimann, I., Mok, V., van Oostenbrugge, R., Pantoni, L., Speck, O., Stephan, B.C.M., Teipel, S., Viswanathan Anand, Werring, D., Chen, C., Smith, C., van Buchem, M.A., Norrving, B., Gorelick, P.B., Dichgans, M.: Neuroimaging standards for research into small vessel disease and its contribution to ageing and neurodegeneration. *Lancet Neurology* **12**, 822–838 (2013)
10. Warfield, S.K., Zou, K.H., Wells, W.M.: Simultaneous truth and performance level estimation (STAPLE): an algorithm for the validation of image segmentation. *IEEE TMI* **23**(7), 903–21 (jul 2004)