

RADIOMIC FEATURE STABILITY ANALYSIS BASED ON PROBABILISTIC SEGMENTATIONS

Christoph Haarbuerger¹, Justus Schock^{2,1}, Daniel Truhn^{3,1}, Philippe Weitz⁴,
Gustav Mueller-Franzes¹, Leon Weninger¹, Dorit Merhof¹

¹Institute of Imaging and Computer Vision, RWTH Aachen University, Germany

²Department of Diagnostic and Interventional Radiology, University Hospital Düsseldorf, Germany

³Department of Diagnostic and Interventional Radiology, University Hospital Aachen, Germany

⁴Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

ABSTRACT

Identifying image features that are robust with respect to segmentation variability and domain shift is a tough challenge in radiomics. So far, this problem has mainly been tackled in test-retest analyses. In this work we analyze radiomics feature stability based on probabilistic segmentations. Based on a public lung cancer dataset, we generate an arbitrary number of plausible segmentations using a Probabilistic U-Net. From these segmentations, we extract a high number of plausible feature vectors for each lung tumor and analyze feature variance with respect to the segmentations. Our results suggest that there are groups of radiomic features that are more (e.g. statistics features) and less (e.g. gray-level size zone matrix features) robust against segmentation variability. Finally, we demonstrate that segmentation variance impacts the performance of a prognostic lung cancer survival model and propose a new and potentially more robust radiomics feature selection workflow.

Index Terms— Radiomics, Segmentation, Feature Stability

1. INTRODUCTION

Radiomics has been increasingly applied to radio- and oncological image data [1, 2, 3]. One of the main problems with the use of radiomics lies in the curse of dimensionality: Most studies are conducted on several hundred images, while thousands of image features are extracted. Several problems arise from this setup: Firstly, by extracting such high numbers of features, despite the use of feature selection algorithms, multicollinearity and overfitting may lead to limited reproducibility of radiomic signatures [4, 5]. While image feature definitions are increasingly standardized [6, 7], other problems lie in differences in image reconstruction [8] and in how the segmentations based on which the features are extracted, are generated. In contrast to many claims [9, 10], most current radiomics pipelines are not truly quantitative because segmen-

tations are performed by humans and are thus subject to inter- and intra-rater variability [11]. We hypothesize that this limits the validity and reproducibility of radiomic signatures, even when they are based on expert segmentations.

Zwanenburg et al. [12] have already assessed robustness of features with respect to image perturbations such as translation, rotation and noise addition. Aerts et al. [1] and Peerlings et al. [13] investigated feature stability in test-retest studies. In [14], robustness of radiomics features with respect to domain shift was assessed across several institutions. Owens et al. [15] evaluated uncertainty of radiomics features for manual versus semi-automatic segmentations. In a recent review article [16], Traverso et al. concluded that there is currently no consensus regarding the question which features are optimal in terms of reproducibility. In order to quantify the degree to which radiomics features depend on the particular segmentation, we propose to perform a *probabilistic* automated segmentation that generates a *set* of plausible segmentations rather than one or few manual segmentations by experts. Feature vectors are then extracted for each segmentation individually, in order to assess robustness of individual features with respect to segmentation variability.

The set of plausible segmentations is generated by an extension of the U-Net [17], the Probabilistic U-Net (PU-Net) [18], which combines the U-Net with a conditional variational autoencoder (CVAE). With this architecture, plausible segmentations can be sampled from the latent space of the CVAE.

Based on these, we extract a high number of radiomics features and assess feature variance with respect to a set of plausible segmentation masks. We identify groups of features that are invariant with respect to the particular segmentation and others that depend on it more heavily. Furthermore, we show that segmentation variance influences the performance of a prognostic survival model on a public lung cancer dataset.

2. MATERIAL AND METHODS

2.1. Image Data

All evaluations are performed based on two publicly available lung cancer datasets: Specifically, the *LIDC-IDRI* dataset [19, 20] is used to train the PU-Net, whereas feature stability is assessed on the *Maastrro Lung1* dataset [1, 21], which are both publicly available at *The Cancer Imaging Archive (TCIA)* [22]. Expert segmentations for all 422 cases of the *Lung1* dataset are also publicly available¹. An example image with a corresponding expert segmentation is depicted in Fig. 1a. Given the CT scans, expert annotations and right-censored survival data, the *Maastrro Lung1* dataset has been used for radiomics-based lung cancer *survival analysis* [1, 4].

2.2. Segmentation using Probabilistic U-Net

The PU-Net [18] is an extension of the popular U-Net architecture [17] that models the distribution of plausible segmentations using a CVAE. By sampling from a low-dimensional latent space vector, plausible segmentation hypotheses can be generated arbitrarily which can be interpreted as an equivalent of asking a human experts to perform a manual segmentations. We trained the PU-Net as implemented in ² using a latent vector with $N = 6$ as suggested in the original PU-Net publication [18]. The network operated on 2D axial slices.

After training the PU-Net on the *LIDC-IDRI* dataset, 1000 plausible 2D segmentation masks were generated for every axial slice from the *Maastrro Lung1* dataset. Unfortunately, for 73 cases, the probabilistic segmentation failed on several slices producing no segmentation at all. By visual inspection we were not able to identify any pattern signifying the failure of generating a segmentation. These cases were excluded from this study.

2.3. Feature Extraction

Using the PU-Net, we generated up to 1000 plausible segmentations for each slice. This initial set of segmentations was then reduced to a set of *unique* segmentations resulting in dozens to few hundred segmentations. For this reduced set, we calculated the Dice score D between each probabilistic segmentation and the ground truth and sampled 25 segmentations with a uniform distribution with respect to the Dice scores related to the ground truth. This step is needed to ensure that the segmentations actually differ because even in the reduced set of segmentations generated by the PU-Net, many are almost identical. The particular choice of 25 segmentations was chosen based on visual inspection considering segmentation diversity and was not optimized further.

Next, for each nodule, 25 feature vectors, based on the 25 segmentations, were extracted using the PyRadiomics framework [7] (Version 2.2.0). Specifically, we extracted:

- 18 statistics features
- 15 shape features
- 22 gray level co-occurrence matrix (GLCM) features
- 16 gray level size zone matrix (GLSZM) features
- 16 gray level run length matrix (GLRLM) features
- 5 neighborhood difference gray tone matrix (NDGTM) features

For feature extraction, all images were resampled to a voxel spacing of $1 \times 1 \times 1 \text{ mm}^3$ and a bin width of 25 was used for grey value binning. Moreover, we extracted not only the image features based on original CT images but also on wavelet-transformed using the standard wavelets transforms implemented in PyRadiomics.

3. RESULTS

For 73 out of the 422 patients, the PU-Net did not produce segmentation hypotheses as explained in Section 2.3 but produced empty masks, which lead to an exclusion of these cases from our study.

Feature stability was assessed by calculating the intraclass correlation coefficient (ICC) for each feature across all 25 segmentations. ICCs grouped by feature categories are provided in Fig. 2. In Fig 3, we provide the same ICCs clustered by image transform (no transform, wavelet transforms).

In Fig. 2 we can see that most statistics, shape, GLCM and NGTDM features have a very high ICC > 0.9 . Many GLRLM and especially GLSZM features have a considerably lower ICCs that can be as low as 0.2. If the features are grouped by image transform as in Fig. 2 it becomes obvious that wavelet features are generally subject to higher segmentation dependence than features extracted from raw CT images. Overall, 28.7% of all features had an ICC < 0.9 .

In order to demonstrate the relation between the prognostic value and the stability a feature, in Fig 4, each for each feature the univariate concordance index (cindex) is plotted against the stability rank. Cindex is a widely-used performance measure in survival analysis and defined as

$$\text{cindex} = \frac{\# \text{ concordant pairs}}{\# \text{ possible pairs}} \in [0, 1]. \quad (1)$$

Stability rank is defined as a descending ranking of all features based on ICC such that the feature with the highest ICC has a stability rank of 1. The last experiment is based on the radiomic signature by Aerts et al. based on a Cox model and the four features identified in [1]. We extracted radiomic signatures based on this model for all 25 segmentations and calculated the corresponding cindices, which are depicted in Fig. 5 as a histogram. Cindices vary between 0.569 and 0.577. In comparison, using the expert segmentation, a cindex of 0.574 is achieved.

¹<https://xnat.bmia.nl>

²<https://github.com/stefanknegt/Probabilistic-Unet-Pytorch>

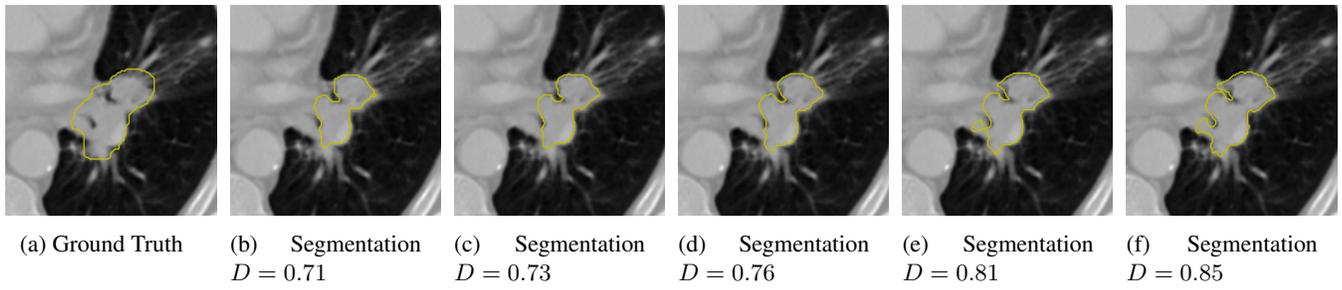


Fig. 1: Ground truth (a) and examples for corresponding probabilistic segmentations (b – f). The dice scores D refer to the ground truth.

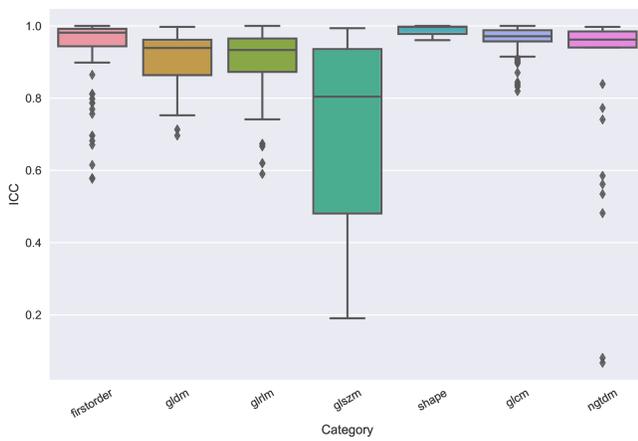


Fig. 2: ICCs for all features clustered by feature category.

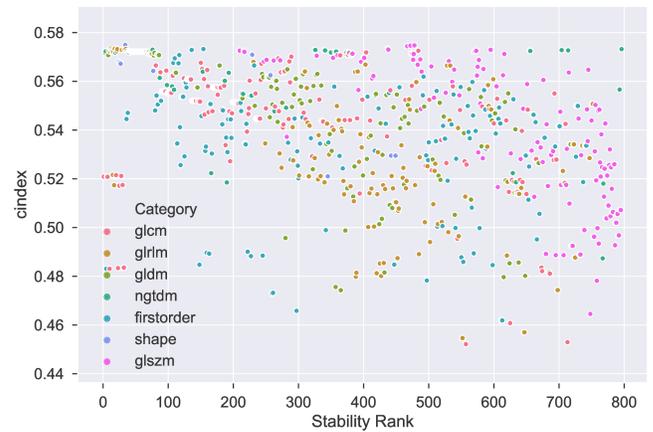


Fig. 4: Stability rank (x-axis) and corresponding concordance index (cindex) for all features.

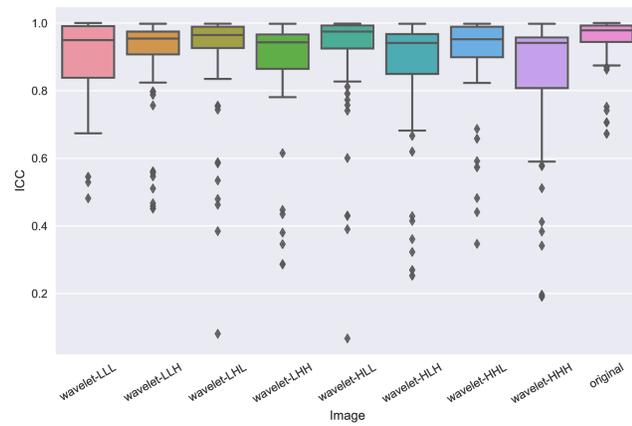


Fig. 3: ICCs for all features clustered by image transform.

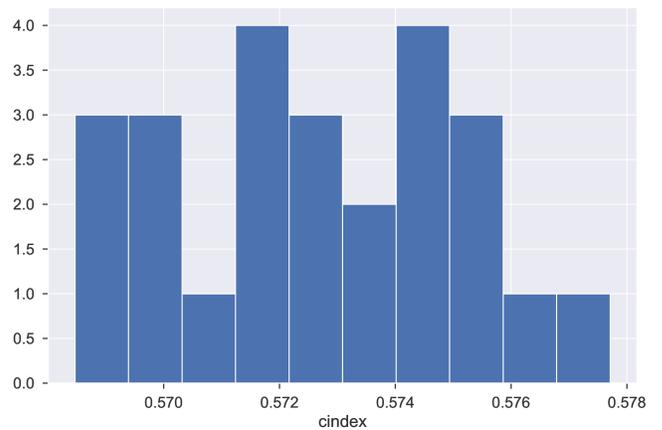


Fig. 5: Histogram of cindices using a Cox model based on the Aerts signature [1] over 25 probabilistic segmentations.

4. DISCUSSION

We have trained a probabilistic segmentation algorithm that provides segmentation hypotheses that can be used for extracting radiomics features. Based on the feature vectors, we have analyzed feature stability with respect to varying segmentation masks. The results show that there are groups of radiomics features that are subject to higher and lower variance across segmentations, respectively. This is in line with other works in which stability of feature vectors originating from multiple manual segmentations by experts were evaluated [1, 13]. Moreover, we were able to show in Fig 4 that variance in segmentations carries over to a prognostic model. However, the variance with respect to the cindex across segmentations is relatively small, indicating that the signature by Aerts et al. is relatively robust against segmentation variance.

There is currently no consensus on which ICC cutoff should be used to exclude features from further analyses [16], however, assuming a cutoff of 0.9, about one third of all features in our analysis could be discarded. Thus, in a standard pipeline for radiomic signature development, the curse of dimensionality and multicollinearity could be considerably alleviated, which may lead to radiomic signatures that are more robust and reproducible. Moreover, feature scores could be averaged over segmentation mask hypotheses to further improve robustness.

Based on these findings, rather than “just” extracting radiomics features from a single expert segmentation, we envision that a future radiomic signature development pipeline could be composed of the following steps:

1. Train a probabilistic segmentation algorithm based on the expert segmentation
2. Generate N plausible segmentations for each case
3. For each feature, calculate ICC with respect to the N segmentations
4. Discard features that are subject to low ICC
5. For the remaining features, average feature scores over N probabilistic segmentations
6. Run “standard” radiomics pipeline (feature selection, model fit, etc.)

Our study has several limitations: First, the dataset used for feature analysis originates from a single scanner, using a single reconstruction algorithm. Thus, variability arising from protocol and device difference as investigated in [5] is not considered here. Furthermore, the publicly available expert segmentations are prone to errors which might carry over to the generated probabilistic segmentations. Accordingly, the PU-Net segmentations have several shortcomings: First, segmentation is only performed on axial slices in 2D. Moreover, training was partly unstable, producing many similar segmentation hypotheses that had to be filtered out. Finally, for 73 cases, the automatic segmentation using the PU-Net

failed completely producing empty masks. These cases had to be excluded from our study and limit the comparability to other works on the same dataset. Finally, our findings are only based on a single dataset from a single modality. In future work, it is desirable to conduct the same experiments on other cancer types and imaging devices to assess feature robustness in a more general sense.

5. CONCLUSION

Using a set of plausible segmentation hypotheses generated by a PU-Net segmentation algorithm, we analyzed variance of radiomic features with respect varying segmentations, showing that there are groups of image features that are subject to different degrees of robustness. Furthermore, we showed that segmentation variance carries impacts a radiomics survival model on a public lung cancer dataset.

6. ACKNOWLEDGEMENTS

The authors wish to acknowledge financial support from Interreg V-A Euregio Meuse-Rhine (“Euradiomics”).

7. REFERENCES

- [1] H. J. W. L. Aerts, E. R. Velazquez, R. T. H. Leijenaar, C. Parmar, P. Grossmann, S. Cavalho, J. Bussink, R. Monshouwer, B. Haibe-Kains, D. Rietveld, F. Hoesbers, M. M. Rietbergen, C. R. Leemans, A. Dekker, J. Quackenbush, R. J. Gillies, and P. Lambin, “Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach,” *Nature Communications*, vol. 5, Jun 2014.
- [2] P. Kickingeder, S. Burth, A. Wick, M. Götz, O. Eidel, H.-P. Schlemmer, K. H. Maier-Hein, W. Wick, M. Bendzus, A. Radbruch, and D. Bonekamp, “Radiomic profiling of glioblastoma: Identifying an imaging predictor of patient survival with improved performance over established clinical and radiologic risk models,” *Radiology*, vol. 280, no. 3, pp. 880–889, sep 2016.
- [3] E. de la Rosa, D. M. Sima, T. V. Vyvere, J. S. Kirschke, and B. Menze, “A radiomics approach to traumatic brain injury prediction in CT scans,” in *2019 IEEE 16th International Symposium on Biomedical Imaging (ISBI 2019)*. Apr. 2019, IEEE.
- [4] M. L. Welch, C. McIntosh, B. Haibe-Kains, M. F. Milosevic, L. Wee, A. Dekker, S. H. Huang, T. G. Purdie, B. O’Sullivan, H. J. Aerts, and D. A. Jaffray, “Vulnerabilities of radiomic signature development: The need for safeguards,” *Radiotherapy and Oncology*, 2018.

- [5] R. Berenguer, M. del Rosario Pastor-Juan, J. Canales-Vázquez, M. Castro-García, M. V. Villas, F. M. Legorburo, and S. Sabater, “Radiomics of CT features may be nonreproducible and redundant: Influence of CT acquisition parameters,” *Radiology*, vol. 288, no. 2, pp. 407–415, Aug. 2018.
- [6] A. Zwanenburg, S. Leger, M. Vallières, and S. Löck, “Image biomarker standardisation initiative - feature definitions,” Dec. 2016.
- [7] J. J. van Griethuysen, A. Fedorov, C. Parmar, A. Hosny, N. Aucoin, V. Narayan, R. G. Beets-Tan, J.-C. Fillion-Robin, S. Pieper, and H. J. Aerts, “Computational radiomics system to decode the radiographic phenotype,” *Cancer Research*, vol. 77, no. 21, pp. e104–e107, oct 2017.
- [8] M. Meyer, J. Ronald, F. Vernuccio, R. C. Nelson, J. C. Ramirez-Giraldo, J. Solomon, B. N. Patel, E. Samei, and D. Marin, “Reproducibility of CT radiomic features within the same patient: Influence of radiation dose and CT reconstruction settings,” *Radiology*, Oct. 2019.
- [9] R. J. Gillies, P. E. Kinahan, and H. Hricak, “Radiomics: Images are more than pictures, they are data,” *Radiology*, vol. 278, no. 2, pp. 563–577, 2016, PMID: 26579733.
- [10] S. S. F. Yip and H. J. W. L. Aerts, “Applications and limitations of radiomics,” *Physics in Medicine and Biology*, vol. 61, no. 13, pp. R150, 2016.
- [11] L. Joskowicz, D. Cohen, N. Caplan, and J. Sosna, “Inter-observer variability of manual contour delineation of structures in ct,” *European Radiology*, vol. 29, no. 3, pp. 1391–1399, Mar 2019.
- [12] A. Zwanenburg, S. Leger, L. Agolli, K. Pilz, E. G. C. Troost, C. Richter, and S. Löck, “Assessing robustness of radiomic features by image perturbation,” *Sci Rep*, vol. 9, no. 1, pp. 614, Jan 2019.
- [13] J. Peerlings, H. C. Woodruff, J. M. Winfield, A. Ibrahim, B. E. V. Beers, A. Heerschap, A. Jackson, J. E. Wildberger, F. M. Mottaghy, N. M. DeSouza, and P. Lambin, “Stability of radiomics features in apparent diffusion coefficient maps from a multi-centre test-retest trial,” *Scientific Reports*, vol. 9, no. 1, Mar. 2019.
- [14] J. Kalpathy-Cramer, A. Mamomov, B. Zhao, L. Lu, D. Cherezov, S. Napel, S. Echegaray, D. Rubin, M. McNitt-Gray, P. Lo, J. C. Sieren, J. Uthoff, S. K. N. Dilger, B. Driscoll, I. Yeung, L. Hadjiiski, K. Cha, Y. Balagurunathan, R. Gillies, and D. Goldgof, “Radiomics of lung nodules: A multi-institutional study of robustness and agreement of quantitative imaging features,” *Tomography (Ann Arbor, Mich.)*, vol. 2, no. 4, pp. 430–437, 12 2016.
- [15] C. A. Owens, C. B. Peterson, C. Tang, E. J. Koay, W. Yu, D. S. Mackin, J. Li, M. R. Salehpour, D. T. Fuentes, L. E. Court, and J. Yang, “Lung tumor segmentation methods: Impact on the uncertainty of radiomics features for non-small cell lung cancer,” *PLOS ONE*, vol. 13, no. 10, pp. e0205003, Oct. 2018.
- [16] A. Traverso, L. Wee, A. Dekker, and R. Gillies, “Repeatability and reproducibility of radiomic features: A systematic review,” *International Journal of Radiation Oncology*, vol. 102, no. 4, pp. 1143–1158, Nov. 2018.
- [17] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *Medical Image Computing and Computer-Assisted Intervention – MICCAI*. 2015, pp. 234–241, Springer International Publishing.
- [18] S. A. A. Kohl, B. Romera-Paredes, C. Meyer, J. D. Fauw, J. R. Ledsam, K. H. Maier-Hein, S. M. A. Eslami, D. J. Rezende, and O. Ronneberger, “A Probabilistic U-Net for Segmentation of Ambiguous Images,” 2018.
- [19] S. G. Armato et al., “The lung image database consortium (LIDC) and image database resource initiative (IDRI): A completed reference database of lung nodules on CT scans,” *Medical Physics*, vol. 38, no. 2, pp. 915–931, Jan. 2011.
- [20] S. Armato et al., “Data from lidc-idri,” 2015.
- [21] H. J. W. L. Aerts, E. Rios Velazquez, R. T. H. Leijenaar, C. Parmar, P. Grossmann, S. Carvalho, and P. Lambin, “Data from NSCLC-radiomics,” in *The Cancer Imaging Archive*, 2015.
- [22] K. Clark, B. Vendt, K. Smith, J. Freymann, J. Kirby, P. Koppel, S. Moore, S. Phillips, D. Maffitt, M. Pringle, L. Tarbox, and F. Prior, “The cancer imaging archive (TCIA): Maintaining and operating a public information repository,” *Journal of Digital Imaging*, vol. 26, no. 6, pp. 1045–1057, July 2013.