Common limitations of performance metrics in biomedical image analysis

A. Reinke^{1,2,*}, M. Eisenmann, M.D. Tizabi, C.H. Sudre, T. Rädsch, M.J. Cardoso, V. Cheplygina, K. Farahani, B. Glocker, P. Godau, D. Heckmann-Nötzel, F. Isensee, P. Jannin, C.E. Kahn, J. Kleesiek, M. Kozubek, T. Kurc, B.A. Landman, G. Litjens, K. Maier-Hein, A.L. Martel, B. Menze, H. Müller, J. Petersen, M. Reyes, M. Riegler, N. Rieke, B. Stieltjes, R.M. Summers, S.A. Tsaftaris, B. van Ginneken, A. Kopp-Schneider, P. Jäger, L. Maier-Hein

Div. Computer Assisted Medical Interventions and HIP Helmholtz Imaging Platform, German Cancer Research Center (DKFZ) ² Faculty of Mathematics and Computer Science, Heidelberg University

* Full list of affiliations: https://arxiv.org/abs/2104.05642

Numerous pitfalls related to metrics

Example of failure:

Effect of small structures in segmentation tasks

Reference

 _	_	 	 	 	_



Prediction										



DSC = 0.79

>>

Single-pixel differences can have huge effects on the metric scores, especially relevant given high inter-rater variability and non-determinism of AI.



The DSC is strongly biased against single objects, therefore not appropriate to measure the detection of multiple objects.

DSC = 0.92

Full paper: Reinke et al. Common Limitations of Image Processing Metrics: A Picture Story. <u>https://arxiv.org/abs/2104.05642</u> You would like to contribute? Contact us!









Common limitations of performance metrics in biomedical image analysis

A. Reinke^{1,2,*}, M. Eisenmann, M.D. Tizabi, C.H. Sudre, T. Rädsch, M.J. Cardoso, V. Cheplygina, K. Farahani, B. Glocker, P. Godau, D. Heckmann-Nötzel, F. Isensee, P. Jannin, C.E. Kahn, J. Kleesiek, M. Kozubek, T. Kurc, B.A. Landman, G. Litjens, K. Maier-Hein, A.L. Martel, B. Menze, H. Müller, J. Petersen, M. Reyes, M. Riegler, N. Rieke, B. Stieltjes, R.M. Summers, S.A. Tsaftaris, B. van Ginneken, A. Kopp-Schneider, P. Jäger, L. Maier-Hein

Div. Computer Assisted Medical Interventions and HIP Helmholtz Imaging Platform, German Cancer Research Center (DKFZ) ² Faculty of Mathematics and Computer Science, Heidelberg University *Full list of affiliations: https://arxiv.org/abs/2104.05642*



Full paper: Reinke et al. Common Limitations of Image Processing Metrics: A Picture Story. <u>https://arxiv.org/abs/2104.05642</u> You would like to contribute? Contact us!

Common limitations of segmentation metrics



Metric aggregation

To produce an aggregated metric value over many images, multiple merging strategies may be applied. Special care has to be given to missing values.



Example of failure: For distance-based measures without lower/upper bounds, the strategy of how to deal with missing values is not trivial. One may choose the maximum distance of the image or normalize the metric values to [0,1] and use the worst possible value (here: 1).

Im	age	I,						
Н	D	.3						
	Igno	re NAs						
Mean HD: 6.76								

HD

Crucially, however, every choice will produce a different aggregated value, thus potentially affecting the ranking.

Metric combination

A single metric typically does not reflect all important aspects that are essential for algorithm validation. Hence, multiple metrics with different properties should be combined. **Raw metric values**

Image	Algorithm	DSC	loU	HD							
I,	AI	0.91	0.82	11.31		Ranking J: Ranking 2: Ranking 3:					
	A2	0.94	0.89	1.10							
	A3	0.95	0.90	14.49		Rank			HD		
							A3	A3	A2	lea	
		•••	•••	•••	\rightarrow	2	ΔΙ	ΔΙ			
I _n	AI	0.77	0.56	6.54			2	۸ <u>٦</u>	A2	A 2	A
	A2	0.75	0.53	2.81			3	A2	AZ	AJ	me
	A3	0.92	0.67	9.22		•••	•••	•••	***	l	



Example of failure: Ignoring missing values leads to a substantially higher DSC compared to setting missing values to the worst possible value (here: 0).



Example of failure: Mutually ependent metrics (DSC and IoU) will nd to the same ranking and should be sed interchanging, whereas metrics easuring different properties (HD) will lead to a different ranking.





Common limitations of performance metrics in biomedical image analysis

A. Reinke^{1,2,*}, M. Eisenmann, M.D. Tizabi, C.H. Sudre, T. Rädsch, M.J. Cardoso, V. Cheplygina, K. Farahani, B. Glocker, P. Godau, D. Heckmann-Nötzel, F. Isensee, P. Jannin, C.E. Kahn, J. Kleesiek, M. Kozubek, T. Kurc, B.A. Landman, G. Litjens, K. Maier-Hein, A.L. Martel, B. Menze, H. Müller, J. Petersen, M. Reyes, M. Riegler, N. Rieke, B. Stieltjes, R.M. Summers, S.A. Tsaftaris, B. van Ginneken, A. Kopp-Schneider, P. Jäger, L. Maier-Hein

Div. Computer Assisted Medical Interventions and HIP Helmholtz Imaging Platform, German Cancer Research Center (DKFZ) ² Faculty of Mathematics and Computer Science, Heidelberg University *Full list of affiliations: https://arxiv.org/abs/2104.05642*



* Thanks to Bernhard Kainz for sharing the figure.

Full paper: Reinke et al. Common Limitations of Image Processing Metrics: A Picture Story. https://arxiv.org/abs/2104.05642 You would like to contribute? Contact us!



Common limitations of detection metrics

Definition of True Positives

Definition of True Positives (TP) in object detection tasks is typically done by measuring the overlap of bounding boxes with the Intersection over Union (IoU). Depending on a threshold, the predicted object will be interpreted as TP or False Positive (FP). This definition is the pre-requisite for metric computation.

Example of failure:

Effect of threshold choice

IoU > 0.50: True positive (TP) $IoU \leq 0.50$: False positive (FP)





The threshold chosen to define TP and FP highly influences the metric values computed from them. Especially for small, diagonal structures, the size of bounding boxes changes quickly, leading to FP although the visual agreement would indicate TP (Prediction 2).





Research for a Life without Cancer





