

Performance Metrics for Probabilistic Ordinal Classifiers

Adrian Galdran^{1,2,✉}

¹ BCN Medtech, Universitat Pompeu Fabra, Barcelona, Spain,
adrian.galdran@upf.edu

² University of Adelaide, Adelaide, Australia

Abstract. Ordinal classification models assign higher penalties to predictions further away from the true class. As a result, they are appropriate for relevant diagnostic tasks like disease progression prediction or medical image grading. The consensus for assessing their categorical predictions dictates the use of distance-sensitive metrics like the Quadratic-Weighted Kappa score or the Expected Cost. However, there has been little discussion regarding how to measure performance of probabilistic predictions for ordinal classifiers. In conventional classification, common measures for probabilistic predictions are Proper Scoring Rules (PSR) like the Brier score, or Calibration Errors like the ECE, yet these are not optimal choices for ordinal classification. A PSR named Ranked Probability Score (RPS), widely popular in the forecasting field, is more suitable for this task, but it has received no attention in the image analysis community. This paper advocates the use of the RPS for image grading tasks. In addition, we demonstrate a counter-intuitive and questionable behavior of this score, and propose a simple fix for it. Comprehensive experiments on four large-scale biomedical image grading problems over three different datasets show that the RPS is a more suitable performance metric for probabilistic ordinal predictions. Code to reproduce our experiments can be found at github.com/agaldran/prob_ord_metrics.

Keywords: Ordinal Classification · Proper Scoring Rules · Model Calibration · Uncertainty Quantification

1 Introduction and Related Work

The output of predictive machine learning models is often presented as categorical values, *i.e.* “hard” class membership decisions. Nonetheless, understanding the faithfulness of the underlying probabilistic predictions giving rise to such hard class decisions can be essential in some critical applications. Meaningful probabilities enable not only high model accuracy, but also more reliable decisions: a doctor may choose to order further diagnostic tests if a binary classifier gives a $p = 45\%$ probability of disease, even if the hard prediction is “healthy” [2]. This is particularly true for ordinal classification problems, *e.g.* disease severity staging [6,7] or medical image grading [14,21]. In these problems, predictions

should be *as close as possible to the actual category*; further away predictions must incur in heavier penalties, as they have increasingly worse consequences.

There is a large body of research around performance metrics for medical image analysis [20]. Most existing measures, like accuracy or the F1-score, focus on assessing hard predictions in specific ways that capture different aspects of a problem. In ordinal classification, the recommended metrics are Quadratic-Weighted Kappa and the Expected Cost [5,16]. In recent years, measuring the performance of “soft” probabilistic predictions has attracted an increasing research interest [12,19]. For this purpose, the current consensus is to employ Calibration Errors like the ECE and Proper Scoring Rules like the Brier score [16]. We will show that other metrics can instead be a better choice for assessing probabilistic predictions in the particular case of ordinal classification problems.

How to measure the correctness of probabilistic predictions is a decades-old question, naturally connected to forecasting, *i.e.* predicting the future state of a complex system [9]. A key aspect of forecasting is that, contrary to classifiers, forecasters do not output hard decisions, but probability distributions over possible outcomes. Weather forecasts do not tell us whether it will rain tomorrow or not, they give us a probability estimate about the likelihood of raining, leaving to us the decision of taking or not an umbrella, considering the personal cost of making such decision. The same applies for financial investments or sports betting, where it is also the final user who judges risks and makes decisions based on probabilistic forecasts. In this context, Proper Scoring Rules (PSRs) have been long used by the forecasting community to measure predictive performance [10]. PSRs are the focus of this paper, and will be formally defined in section 2.1.

Relation to Calibration: A popular approach to assess the quality of probabilistic predictions is measuring calibration. A model is well calibrated if its probabilistic predictions are aligned with its accuracy on average. PSRs and calibration are intertwined concepts: PSRs can be decomposed into a calibration and a resolution component [8]. Therefore, a model needs to be both calibrated and resolved (*i.e.* having *sharp*, or *concentrated* probabilities) in order to have a good PSR value. For example, if a disease appears in 60% of the population, and our model is just “`return p=0.6`”, in the long run the model is correct 60% of the time, and it is perfectly calibrated, as its confidence is fully aligned with its accuracy, despite having zero predictive ability. If the model predicted in a “resolved” manner with $p = 0.99$ the presence of the disease, but being correct only 70% of the time, then it would be overconfident, which is a form of miscalibration. Only when the model is simultaneously confident and correct can it attain a good PSR value.

The two most widely adopted PSRs are the Brier and the Logarithmic Score [1,11]. Unfortunately, none of these is appropriate for the assessment of ordinal classification probabilities [3]. A third PSR, long used by forecasting researchers in this scenario, the Ranked Probability Score (RPS, [4]), appears to have been neglected so far in biomedical image grading applications. This paper first covers the definition and basic properties of PSRs, and then motivates the use the RPS for ordinal classifiers. We also illustrate a counter-intuitive behavior of the RPS, and propose a simple modification to solve it. Our experiments cover two relevant

biomedical image grading problems and illustrate how the RPS can better assess probabilistic predictions of ordinal classification models.

2 Methods

2.1 Scoring Rules - Notation, Properties, Examples

We consider a K -class classification problem, and a classifier that takes an image \mathbf{x} and maps it into a vector of probabilities $\mathbf{p} \in [0, 1]^K$. Typically, \mathbf{p} is the result of applying a softmax operation on the output of a neural network. Suppose \mathbf{x} belongs to class $y \in \{1, \dots, K\}$, and denote by \mathbf{y} its one-hot representation. A Scoring Rule (SR) \mathcal{S} is any function taking the probabilistic prediction \mathbf{p} and the label \mathbf{y} and producing a number $\mathcal{S}(\mathbf{p}, \mathbf{y}) \in \mathbb{R}$ (a score). Here we consider negatively oriented SRs, which assign lower values to *better predictions*.

Of course, the above is an extremely generic definition, to which we must now attach additional properties in order to encode our understanding of what *better predictions* means for a particular problem.

Property 1: A Scoring Rule (SR) is *proper* if its value is minimal when the probabilistic prediction coincides with the ground-truth in expectation.

Example: The Brier Score [1] is defined as the sum of the squared differences between probabilities and labels:

$$\text{Brier}(\mathbf{p}, \mathbf{y}) = \|\mathbf{p} - \mathbf{y}\|_2^2 = \sum_{i=1}^K (p_i - y_i)^2. \quad (1)$$

Since its value is always non-negative, and it decreases to 0 when $\mathbf{p} = \mathbf{y}$, we conclude that the Brier Score is indeed proper.

Property 2: A Proper Scoring Rule (PSR) is *local* if its value only depends on the probability assigned to the correct category.

Example: The Brier Score is non-local, as its value depends on the probability placed by the model on all classes. The Logarithmic Score [11], given by:

$$\mathcal{L}(\mathbf{p}, \mathbf{y}) = -\log(p_c) \quad (2)$$

where c is the correct category of \mathbf{x} , rewards the model by placing as much probability mass as possible in c , regardless of how the remaining probability is distributed. It is, therefore, a local PSR. The Logarithmic Score is also known, when taken on average over a dataset, as the Negative Log-Likelihood.

Property 3: A PSR is *sensitive to distance* if its value takes into account the order of the categories, in such a way that probability placed in categories further away from the correct class is more heavily penalized.

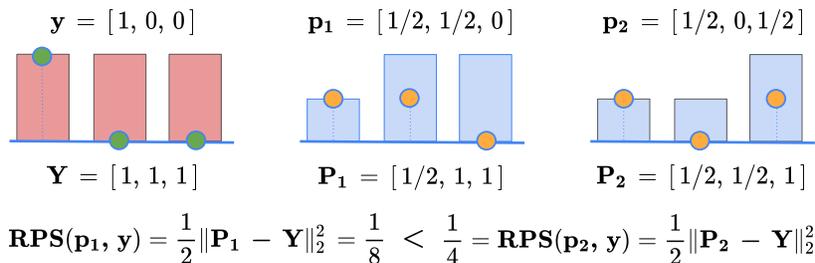


Fig. 1: The RPS is sensitive to distance, suitable for assessing probabilistic predictions on biomedical image grading problems. It is the difference between the cumulative probability distributions of the label and a probabilistic prediction.

Example: Both the Brier and the Logarithmic scores are insensitive to distance (shuffling \mathbf{p} and \mathbf{y} won't affect the score). Sensitivity to distance is essential for assessing ordinal classifiers. Below we define the Ranked Probability Score (RPS) [4,18], which has this property, and is therefore more suitable for our purposes.

2.2 The Ranked Probability Score for Ordinal Classification

Consider a test sample (\mathbf{x}, \mathbf{y}) in a 3-class classification problem, with label \mathbf{y} and two probabilistic predictions $\mathbf{p}_1, \mathbf{p}_2$:

$$\mathbf{y} = [1, 0, 0], \quad \mathbf{p}_1 = \left[\frac{1}{4}, \frac{3}{4}, 0 \right], \quad \mathbf{p}_2 = \left[\frac{1}{4}, 0, \frac{3}{4} \right] \quad (3)$$

In this scenario, both the Brier and the Logarithmic scores produce the same penalty for each prediction, whereas a user might prefer \mathbf{p}_1 over \mathbf{p}_2 due to the latter assigning more probability to the second category. Indeed, if we use the arg-max operator to generate a hard-decision for this sample, we will obtain a prediction of class 2 and class 3 respectively, which could result in the second model declaring a patient as severely unhealthy with serious consequences. In this context, we would like to have a PSR that takes into account distance to the true category, such as the Ranked Probability Score (RPS, [4]), given by:

$$\text{RPS}(\mathbf{p}, \mathbf{y}) = \frac{1}{K-1} \sum_{i=1}^{K-1} \left[\sum_{j=1}^i (p_j - y_j) \right]^2 = \frac{1}{K-1} \|\mathbf{P} - \mathbf{Y}\|_2^2. \quad (4)$$

The RPS is the squared ℓ_2 distance between the cumulative distributions \mathbf{Y} of the target label \mathbf{y} and \mathbf{P} of the probabilistic prediction \mathbf{p} , discounting their last component (as they are both always one) and normalizing so that it varies in the unit interval. In the above example, the RPS would give for each prediction a penalty of $\text{RPS}(\mathbf{p}_1, \mathbf{y}) = 1/8$, $\text{RPS}(\mathbf{p}_2, \mathbf{y}) = 1/4$, as shown in Figure 1.

Among many interesting properties, one can show that the RPS is proper [17], and reduces to the Brier score for $K = 2$. Despite the RPS dating back

more than 50 years [4], and enjoying great popularity in the weather forecasting community, it appears to be much less known in the image analysis and computer vision areas, where we could not find any trace of it. The **first goal** of this paper is to bring to the attention of computer vision researchers this tool for measuring the performance of probabilistic predictions in ordinal classification.

2.3 The Squared Absolute RPS

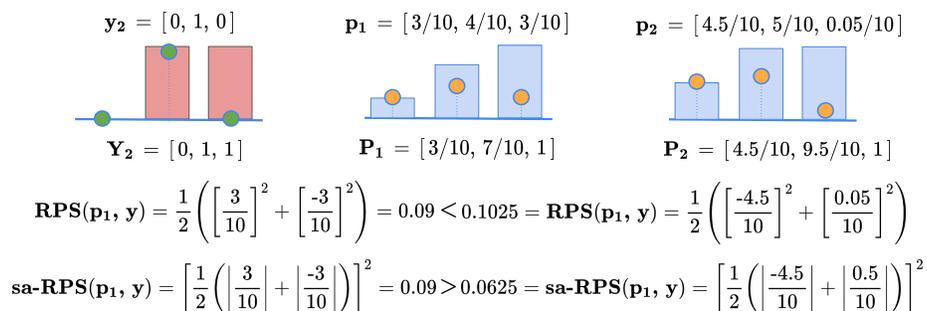
Our **second goal** in this paper is to identify and then fix certain failure modes of the RPS that might lead to counter-intuitive behaviors. First, in disease grading and other ordinal classification problems it is customary to assign penalties to mistakes that grow quadratically with the distance to the correct category. This is the reason why most works utilize the Quadratic-Weighted Kappa Score (QWK) instead of the linearly weighted version of this metric. However, the RPS increases the penalty linearly, as can be quickly seen with a simple 3-class problem and an example $(\mathbf{x}_1, \mathbf{y}_1)$ of class 1 ($\mathbf{y}_1 = [1, 0, 0]$):

$$\text{RPS}([1, 0, 0], \mathbf{y}_1) = 0, \quad \text{RPS}([0, 1, 0], \mathbf{y}_1) = 1/2, \quad \text{RPS}([0, 0, 1], \mathbf{y}_1) = 1. \quad (5)$$

Also, the RPS has a hidden preference for symmetric predictions. To see this, consider a second example $(\mathbf{x}_2, \mathbf{y}_2)$ in which the correct category is now the middle one ($\mathbf{y}_2 = [0, 1, 0]$), and two probabilistic predictions: $p_{\text{sym}} = [3/10, 4/10, 3/10]$, $p_{\text{asym}} = [1/10, 5/10, 9/10]$. In principle, there is no reason to prefer p_{sym} over p_{asym} , unless certain prior/domain knowledge tells us that symmetry is a desirable property. In this particular case, p_{asym} is actually more confident on the correct class than p_{sym} , which is however the preferred prediction for the RPS:

$$\text{RPS}([0.30, 0.40, 0.30], \mathbf{y}_2) = 0.09 < 0.1025 = \text{RPS}([0.45, 0.50, 0.05], \mathbf{y}_2). \quad (6)$$

Fig. 2: The Ranked Probability Score displays some counter-intuitive behavior that the proposed sa-RPS can fix. Here, \mathbf{p}_2 places more probability on the correct class but \mathbf{p}_1 is preferred due to its symmetry.



In order to address these aspects of the conventional RPS, we propose to implement instead the Squared Absolute RPS (sa-RPS), given by:

$$\text{sa-RPS}(\mathbf{p}, \mathbf{y}) = \frac{1}{K-1} \left[\sum_{i=1}^K \left| \sum_{j=1}^i (p_j - y_j) \right| \right]^2 \quad (7)$$

Replacing the inner square in eq. (4) by an absolute value, we manage to break the preference for symmetry of the RPS, and squaring the overall result we build a metric that still varies in $[0,1]$ but gives a quadratic penalty to further away predictions. This is illustrated in Fig. 2 above.

2.4 Evaluating Evaluation Metrics

Our **third goal** is to demonstrate how the (sa-)RPS is useful for evaluating probabilistic ordinal predictions. In the next section we will show some illustrative examples that qualitatively demonstrate its superiority over the Brier and logarithmic score. However, it is hard to quantitatively make the case for one performance metric over another, since metrics themselves are what quantify modeling success. We proceed as follows: we first train a neural network to solve a biomedical image grading problem. We generate probabilistic predictions on the test set and apply distance sensitive metrics to (arg-maxed) hard predictions (QWK and EC, as recommended in [16]), verifying model convergence.

Here it is important to stress that, contrary to conventional metrics (like accuracy, QWK, or ECE) PSRs can act on an individual datum, without averaging over sets of samples. We exploit this property to design the following experiment: we sort the probabilistic predictions of the test set according to a score \mathcal{S} , and then progressively remove samples that are of worst quality according to \mathcal{S} . We take the arg-max on the remaining probabilistic predictions and compute QWK and EC. If \mathcal{S} prefers better ordinal predictions, we must see a performance increase on that subset. We repeat this process, each time removing more of the worse samples, and graph the evolution of QWK and EC for different scores \mathcal{S} : a better score should result in a faster QWK/EC-improving trend.

Lastly, in order to derive a single number to measure performance, we compute the area under the remaining samples vs QWK/EC curve, which we call Area under the Retained Samples Curve (AURSC). In summary:

What we expect to see:

As we remove test set samples considered as worse classified by RPS, we expect to more quickly improve QWK/EC on the resulting subsets. We measure this with the Area under the Retained Samples Curve (AURSC)

3 Experimental Results

We now give a description of the data we used for experimentation, analyze performance for each considered problem, and close with a discussion of results.

3.1 Datasets and Architecture

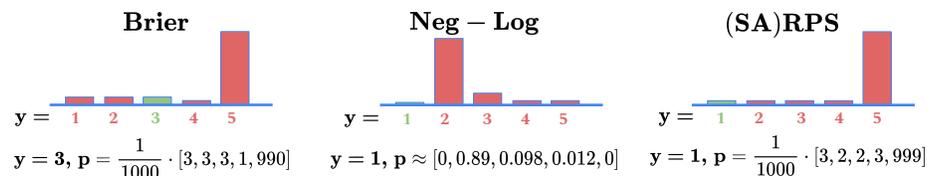
Our experiments are on two different medical image grading tasks: **1)** the **TMED-v2** dataset ([13], link) contains 17,270 images from 577 patients, with an aortic stenosis (AS) diagnostic label from three categories (none, early AS, or significant AS). The authors provide an official train/test distribution of the data that we use here. **2)** **Eyepacs** (link) contains retinal images and labels for grading Diabetic Retinopathy (DR) stage into five categories, ranging from healthy to proliferative DR. It has 35,126 images for training and 53,576 in the test set.

We train a ConvNeXt [15], minimizing the CE loss with the adam algorithm for 10 epochs starting with a learning rate of $l = 1e-4$, decaying to zero over the training. We report average Area under the Retained Samples Curve (AURSC) for 50 bootstrap iterations in each dataset below, and also plot the evolution of performance as we remove more samples considered to be worse by four PSRs: the Brier score, the Logarithmic score (Neg-Log), RPS and sa-RPS.

3.2 How is RPS useful? Qualitative Error Analysis

The obvious application of RPS would be to train better ordinal classification models. But beyond this, RPS also enables improved, fine-grained error analysis. Let us see this through a simple experiment. Since PSRs assess samples individually, we can sort our test set using RPS, NLL, and Brier score. The worst-scored items are what the model considers the wrongest probabilistic predictions. The result of sorting predictions on the Eyepacs test set with the Brier, Neg-Log and RPS rules is shown on Fig. 3. We can see that the prediction identified as worst by the RPS does indeed violate more heavily the order of categories, placing more probability on class 5 for a sample of class 1. On the other hand, for the same test set and predictions, the Brier score finds worst a prediction with 99% of the probability on class 3 and a label of class 5, and the Neg-Log score identifies a sample of class 1 for which the model wrongly predicts class 2.

Fig. 3: For the same test set and predictions, the RPS finds wrong samples that are more incorrect from the point of view of ordinal classification.



3.3 Quantitative Experimental Analysis

Quantitative results of the experiment described in section 2.4, computing AURSC values for all PSRs, are shown in Table 2, with dispersion measures obtained from 50 bootstrapped performance measurements.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	13.46 \pm 0.35	3.76 \pm 0.21	17.36 \pm 0.04	2.84 \pm 0.07
Neg-Log	13.56 \pm 0.35	3.62 \pm 0.2	17.44 \pm 0.04	2.67 \pm 0.07
RPS	14.76 \pm 0.28	2.68 \pm 0.14	17.81 \pm 0.03	1.99 \pm 0.04
sa-RPS	14.95 \pm 0.25	2.53 \pm 0.12	17.86 \pm 0.03	1.88 \pm 0.04

Table 1: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **ConvNeXt**, for each PSR; **best** and **second best** values are marked.

We see that for the considered ordinal classification problems, distance-sensitive scores consistently outperform the Brier and Neg-Log scores. Also, the Square-Absolute Ranked Probability Score always outperforms the conventional Ranked Probability Score. It is worth stressing that when observing bootstrapped performance intervals, neither the Brier nor the Logarithmic scores manage to overlap the SA-RPS interval in any of the two datasets, and in the Eyepacs dataset not even the best RPS result reaches the performance of worst SA-RPS result.

For a visual analysis, Fig. 4 shows the full Sample Retention Curves from which AURSC-QWK values in Table 1 were computed. These curves show how

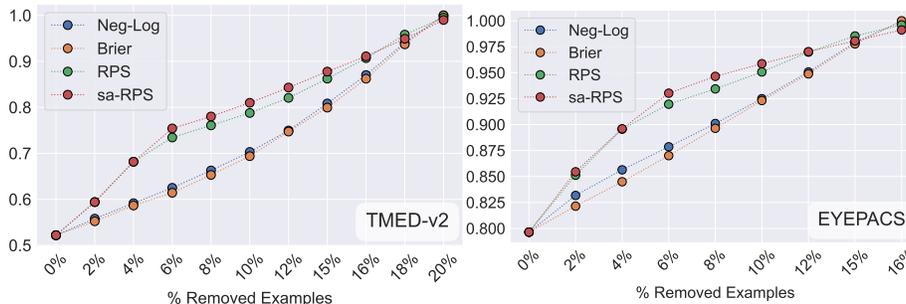


Fig. 4: We sort probabilistic predictions in each test set using several PSRs: Brier, Neg-Log, RPS, sa-RPS. We progressively discard worse-scored samples, improving the metric of interest (only QWK shown). Removing worse samples according to RPS and sa-RPS leads to better QWK, implying that they both capture better ordinal classification performance at the probabilistic level.

PSRs can indeed take a single probabilistic prediction and return a score that is correlated to QWK, which is computed over sets of samples. This is because as we remove samples according to any PSR, performance in the remaining test set improves in all cases. The curves in Fig. 4 also tell a more complete story of how the two distance-sensitive scores outperform the Brier and Neg-Log scores, particularly for TMED and Eyepacs. Just by removing a 5%-6% of samples with worse (higher) RPS, we manage to improve QWL and EC to a greater extent.

4 Conclusion and Future Work

We have shown that Proper Scoring Rules are useful tools for diagnosing probabilistic predictions, but the standard Brier and Logarithmic scores should not be preferred in ordinal classification problems like medical image grading. Instead, the Ranked Probability Score, popular in the forecasting community, should be favoured. We have also proposed sa-RPS, an extension of the RPS that can better handle some pathological cases. Future work will involve using the RPS to learn ordinal classifiers, and investigating its impact in calibration problems.

Acknowledgments

This work was supported by a Marie Skłodowska-Curie Fellowship (No 892297).

References

1. Brier, G.W.: Verification of Forecasts Expressed in Terms of Probability. *Monthly Weather Review* **78**(1), 1–3 (Jan 1950). [https://doi.org/10.1175/1520-0493\(1950\)078<0001:VOFEIT>2.0.CO;2](https://doi.org/10.1175/1520-0493(1950)078<0001:VOFEIT>2.0.CO;2), publisher: American Meteorological Society Section: *Monthly Weather Review* 2, 3
2. Cahan, A., Gilon, D., Manor, O., Paltiel, O.: Probabilistic reasoning and clinical decision-making: do doctors overestimate diagnostic probabilities? *QJM: An International Journal of Medicine* **96**(10), 763–769 (Oct 2003). <https://doi.org/10.1093/qjmed/hcg122>, <https://doi.org/10.1093/qjmed/hcg122> 1
3. Constantinou, A.C., Fenton, N.E.: Solving the Problem of Inadequate Scoring Rules for Assessing Probabilistic Football Forecast Models. *Journal of Quantitative Analysis in Sports* **8**(1) (Mar 2012). <https://doi.org/10.1515/1559-0410.1418>, publisher: De Gruyter 2
4. Epstein, E.S.: A Scoring System for Probability Forecasts of Ranked Categories. *Journal of Applied Meteorology and Climatology* **8**(6), 985–987 (Dec 1969). [https://doi.org/10.1175/1520-0450\(1969\)008<0985:ASSFPF>2.0.CO;2](https://doi.org/10.1175/1520-0450(1969)008<0985:ASSFPF>2.0.CO;2) 2, 4, 5
5. Ferrer, L.: Analysis and Comparison of Classification Metrics (Nov 2022). <https://doi.org/10.48550/arXiv.2209.05355>, <http://arxiv.org/abs/2209.05355>, arXiv:2209.05355 [cs] 2
6. Galdran, A., Chelbi, J., Kobi, R., Dolz, J., Lombaert, H., ben Ayed, I., Chakor, H.: Non-uniform Label Smoothing for Diabetic Retinopathy Grading from Retinal Fundus Images with Deep Neural Networks. *Translational Vision Science & Technology* **9**(2), 34 (Jun 2020). <https://doi.org/10.1167/tvst.9.2.34> 1

7. Galdran, A., Dolz, J., Chakor, H., Lombaert, H., Ben Ayed, I.: Cost-Sensitive Regularization for Diabetic Retinopathy Grading from Eye Fundus Images. In: Medical Image Computing and Computer Assisted Intervention – MICCAI 2020. pp. 665–674. Lecture Notes in Computer Science, Springer International Publishing, Cham (2020). https://doi.org/10.1007/978-3-030-59722-1_64 1
8. Gneiting, T., Balabdaoui, F., Raftery, A.E.: Probabilistic forecasts, calibration and sharpness. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **69**(2), 243–268 (2007). <https://doi.org/10.1111/j.1467-9868.2007.00587.x>, [_eprint: https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9868.2007.00587.x](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1467-9868.2007.00587.x) 2
9. Gneiting, T., Katzfuss, M.: Probabilistic Forecasting. *Annual Review of Statistics and Its Application* **1**(1), 125–151 (2014). <https://doi.org/10.1146/annurev-statistics-062713-085831> 2
10. Gneiting, T., Raftery, A.E.: Strictly Proper Scoring Rules, Prediction, and Estimation. *Journal of the American Statistical Association* **102**(477), 359–378 (Mar 2007). <https://doi.org/10.1198/01621450600001437> 2
11. Good, I.J.: Rational Decisions. *Journal of the Royal Statistical Society: Series B (Methodological)* **14**(1), 107–114 (1952). <https://doi.org/10.1111/j.2517-6161.1952.tb00104.x> 2, 3
12. Gruber, S.G., Buettner, F.: Better Uncertainty Calibration via Proper Scores for Classification and Beyond (Oct 2022), <https://openreview.net/forum?id=PikKk21F6P> 2
13. Huang, Z., Long, G., Wessler, B., Hughes, M.C.: TMED 2: A Dataset for Semi-Supervised Classification of Echocardiograms. In: Unpublished Technical Report (2022), https://tmed.cs.tufts.edu/papers/HuangEtAl_TMED2_DataPerf_2022.pdf 7
14. Jaroensri, R., et al.: Deep learning models for histologic grading of breast cancer and association with disease prognosis. *npj Breast Cancer* **8**(1), 1–12 (Oct 2022). <https://doi.org/10.1038/s41523-022-00478-y> 1
15. Liu, Z., Mao, H., Wu, C.Y., Feichtenhofer, C., Darrell, T., Xie, S.: A ConvNet for the 2020s. pp. 11976–11986 (2022) 7
16. Maier-Hein, L., et al.: Metrics reloaded: Pitfalls and recommendations for image analysis validation (Feb 2023). <https://doi.org/10.48550/arXiv.2206.01653>, <http://arxiv.org/abs/2206.01653>, arXiv:2206.01653 [cs] 2, 6
17. Murphy, A.H.: On the “Ranked Probability Score”. *Journal of Applied Meteorology and Climatology* **8**(6), 988–989 (Dec 1969). [https://doi.org/10.1175/1520-0450\(1969\)008<0988:OTPS>2.0.CO;2](https://doi.org/10.1175/1520-0450(1969)008<0988:OTPS>2.0.CO;2) 4
18. Murphy, A.H.: A Note on the Ranked Probability Score. *Journal of Applied Meteorology and Climatology* **10**(1), 155–156 (Feb 1971). [https://doi.org/10.1175/1520-0450\(1971\)010<0155:ANOTRP>2.0.CO;2](https://doi.org/10.1175/1520-0450(1971)010<0155:ANOTRP>2.0.CO;2) 4
19. Perez-Lebel, A., Morvan, M.L., Varoquaux, G.: Beyond calibration: estimating the grouping loss of modern neural networks (Jan 2023). <https://doi.org/10.48550/arXiv.2210.16315>, <http://arxiv.org/abs/2210.16315>, arXiv:2210.16315 [cs, stat] 2
20. Reinke, A., et al.: Understanding metric-related pitfalls in image analysis validation (Feb 2023). <https://doi.org/10.48550/arXiv.2302.01790> 2
21. Silva-Rodríguez, J., Colomer, A., Sales, M.A., Molina, R., Naranjo, V.: Going deeper through the Gleason scoring scale: An automatic end-to-end system for histology prostate grading and cribriform pattern detection. *Computer Methods and Programs in Biomedicine* **195**, 105637 (Oct 2020). <https://doi.org/10.1016/j.cmpb.2020.105637> 1

Appendix A Further Experimental Results

In the main paper we reported results without dispersion measures due space constraints. Below we show expanded tables with standard deviation, and we include several more architectures (plus ConvNeXt in Tab. 2, here also with dispersion measures). Specifically, we include results obtained with a Swin-Transformer (Tab. 3), a Densnet121 (Tab. 4), a Resnet50 (Tab. 5), a Resnet18 (Tab. 7, and a Mobilent V2 (Tab. 8). In general, we observe similar trends as in the main paper. For all six considered neural networks, distance-sensitive PSRs always achieve greater performance than the Brier and the Logarithmic scores, and also the Square-Absolute Ranked Probability Score always outperforms the conventional Ranked Probability Score. It is worth stressing that when observing bootstrapped performance intervals, the neither the Brier nor the Logarithmic scores manage to overlap the SA-RPS interval in any of the two datasets, and in the Eyepacs dataset not even the best RPS result reaches the performance of worst SA-RPS result.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	13.46 \pm 0.35	3.76 \pm 0.21	17.36 \pm 0.04	2.84 \pm 0.07
Neg-Log	13.56 \pm 0.35	3.62 \pm 0.2	17.44 \pm 0.04	2.67 \pm 0.07
RPS	14.76 \pm 0.28	2.68 \pm 0.14	17.81 \pm 0.03	1.99 \pm 0.04
sa-RPS	14.95 \pm 0.25	2.53 \pm 0.12	17.86 \pm 0.03	1.88 \pm 0.04

Table 2: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **ConvNeXt**, for each PSR; **best** and **second best** values are marked.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	13.60 \pm 0.39	3.56 \pm 0.25	17.33 \pm 0.04	2.84 \pm 0.07
Neg-Log	13.71 \pm 0.39	3.43 \pm 0.25	17.40 \pm 0.04	2.68 \pm 0.06
RPS	14.59 \pm 0.30	2.75 \pm 0.18	17.77 \pm 0.03	2.02 \pm 0.05
sa-RPS	14.79 \pm 0.26	2.58 \pm 0.15	17.82 \pm 0.03	1.91 \pm 0.04

Table 3: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **Swin-T**, for each PSR; **best** and **second best** values are marked.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	13.25 \pm 0.35	3.76 \pm 0.2	17.16 \pm 0.05	3.07 \pm 0.08
Neg-Log	13.31 \pm 0.34	3.66 \pm 0.2	17.23 \pm 0.05	2.91 \pm 0.07
RPS	14.43 \pm 0.28	2.81 \pm 0.15	17.67 \pm 0.03	2.15 \pm 0.05
sa-RPS	14.70 \pm 0.23	2.59 \pm 0.11	17.72 \pm 0.03	2.02 \pm 0.04

Table 4: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **DenseNet**, for each PSR; **best** and **second best** values are marked.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	13.14 \pm 0.44	3.93 \pm 0.27	16.91 \pm 0.05	3.5 \pm 0.08
Neg-Log	13.24 \pm 0.43	3.80 \pm 0.26	16.97 \pm 0.05	3.34 \pm 0.07
RPS	14.32 \pm 0.34	2.98 \pm 0.19	17.50 \pm 0.03	2.41 \pm 0.04
sa-RPS	14.46 \pm 0.30	2.84 \pm 0.16	17.54 \pm 0.03	2.30 \pm 0.04

Table 5: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **ResNet50**, for each PSR; **best** and **second best** values are marked.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	13.41 \pm 0.41	3.61 \pm 0.23	16.87 \pm 0.05	3.59 \pm 0.08
Neg-Log	13.47 \pm 0.41	3.51 \pm 0.22	16.93 \pm 0.05	3.44 \pm 0.08
RPS	14.49 \pm 0.33	2.73 \pm 0.16	17.48 \pm 0.04	2.46 \pm 0.05
sa-RPS	14.82 \pm 0.29	2.47 \pm 0.13	17.54 \pm 0.03	2.31 \pm 0.05

Table 6: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **ResNet34**, for each PSR; **best** and **second best** values are marked.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	13.12 \pm 0.33	4.33 \pm 0.23	16.74 \pm 0.06	3.66 \pm 0.08
Neg-Log	13.14 \pm 0.34	4.25 \pm 0.24	16.83 \pm 0.05	3.45 \pm 0.08
RPS	14.41 \pm 0.27	3.21 \pm 0.17	17.36 \pm 0.04	2.55 \pm 0.05
sa-RPS	14.73 \pm 0.23	2.93 \pm 0.14	17.41 \pm 0.04	2.42 \pm 0.05

Table 7: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **ResNet18**, for each PSR; **best** and **second best** values are marked.

	TMED		Eyepacs	
	AURSC-QWK \uparrow	AURSC-EC \downarrow	AURSC-QWK \uparrow	AURSC-EC \downarrow
Brier	11.13 \pm 0.43	4.92 \pm 0.23	16.62 \pm 0.05	3.85 \pm 0.07
Neg-Log	11.13 \pm 0.43	4.84 \pm 0.23	16.71 \pm 0.05	3.63 \pm 0.07
RPS	12.61 \pm 0.36	3.74 \pm 0.17	17.28 \pm 0.04	2.66 \pm 0.05
sa-RPS	12.86 \pm 0.33	3.46 \pm 0.15	17.34 \pm 0.03	2.52 \pm 0.04

Table 8: Areas under the Retained Samples Curve for **TMED** and **Eyepacs**, with a **Mobilenet**, for each PSR; **best** and **second best** values are marked.