

1 **ImageNet results (R2, R3 & R4):** As suggested by most reviewers,  
 2 we perform our experiments on the ImageNet dataset by certifying  
 3 a pre-trained ResNet-50 from Cohen et al. in [6] and produce the  
 4 following plots similar to Figure 2 in our paper. Just like the CIFAR-  
 5 10 results, the CDF-based method significantly outperforms the naive  
 6 baseline for different levels of smoothing noise  $\sigma$ . Due to space  
 7 constraints, we show plots for  $\sigma = 0.25$  noting that plots for  $\sigma = 0.5$   
 8 demonstrate similar trends.

9 **Calibration of confidence information (R1):** We agree that it is  
 10 improper to interpret confidence scores as the “probability” of the  
 11 classification being correct – indeed, this is the reason why we prefer  
 12 the term “class score” over “class probability.” However, we think it  
 13 is widely accepted that scores are often strongly correlated with the  
 14 chance of correct classification, even if this is not an exact or linear  
 15 relationship because of poor calibration.

16 **Correlation between prediction score and certified radius (R1):**  
 17 We do not mean to imply no correlation exists, but rather that the  
 18 correlation in Fig 4 is *weak*; there are many images with large radius  
 19 but low confidence, and visa versa.

20 **Theorem 1 is an immediate application of Lemma 2 in Salman**  
 21 **et al. (R1):** While the two theorems might be related, the connection  
 22 does not seem so obvious for us. Theorem 1 in our paper is regarding  
 23 the *baseline* certificate against which we compare our CDF-based certificate (Theorem 2) which is our main technical  
 24 contribution. We will make a reference to this lemma in the paper.

25 **Parameters  $s_1, \dots, s_n$  are not discussed in text (R2):** In our experiments, we set  $s_1, \dots, s_n$  such that the number  
 26 of confidence score values in every interval  $(a, s_1), (s_1, s_2), \dots, (s_n, b)$  is equal. We chose this approach instead of  
 27 setting the parameters at regular intervals in  $(a, b)$  because doing so would make a lot of intervals empty. It made sense  
 28 to split the range in such a way that the intervals are well-balanced. There could be other reasonable ways to set these  
 29 parameters which could have an impact on the performance. We will add this discussion to the paper.

30 **In page 3, how to determine the parameter  $a, b$  and  $c$ ? (R3):** Parameters  $a$  and  $b$  depend on the confidence measure  
 31 used and denote the range  $(a, b)$  in which the confidence scores lie. For instance, if you are using average prediction  
 32 score, then  $a = 0$  and  $b = 1$  (see page 6 under section 4: Confidence measures) and if you are considering margin of  
 33 average prediction score, then  $a = -1$  and  $b = 1$ . Parameter  $c$  is the confidence threshold you would like to certify,  
 34 i.e., the confidence score is guaranteed to be above  $c$  within the certified radius. This parameter can be set by the user  
 35 to calculate the corresponding certified radius. In our experiments, to keep computing requirements within reach, we  
 36 compute the threshold  $c$  for a chosen radius instead. But, the original task of computing the radius for a given  $c$  can be  
 37 accomplished using a simple binary-search. (See page 5, lines 172-175).

38 **The rationale behind confidence notions (R3):** The *average* score is a very conventional measure of classification  
 39 confidence, and is the more direct analog of classical confidence scores in the setting of a smoothed classifier. The  
 40 classifier *margin* is less conventional; it measures how much more certain the top class is over the second-place class.  
 41 We find that the distribution of margins is often more concentrated than average scores, and we can therefore produce  
 42 certificates with stronger bounds. There is a tradeoff here that the user can choose – the average is more interpretable but  
 43 provides weaker bounds, the margin is less conventional (although large margin is still a good indicator of confidence)  
 44 and provides tighter certificates.

45 **Relationship between radii  $\sigma$  and confidence threshold (R3):** If you consider two Gaussians separated by  $\epsilon$  units,  
 46 their overlap will be large when  $\sigma$  is large relative to  $\epsilon$ , and the overlap is very small when  $\sigma$  is small. For this reason,  
 47 when  $\sigma$  is large, the confidence bound decreases more slowly with increasing  $\epsilon$ , allowing us to certify larger radii.  
 48 However, when  $\epsilon$  is small (e.g., for certifying small radii), the overlap between Gaussians is large, even for small  $\sigma$ , and  
 49 so we can produce certificates with small  $\sigma$ . In this case, we can take advantage of the fact that the sampled scores are  
 50 more homogenous when  $\sigma$  is small to get tighter bounds. Our experiments show numerous values of  $\sigma$  to show that our  
 51 results hold for a range of choices. However in practice the optimal  $\sigma$  is proportional to the radius you want to certify.  
 52 We will include this discussion when we update our paper.

53 **Related work (R2 & R3):** We will be sure to add a section on related work in randomized smoothing and certifying  
 54 confidences.

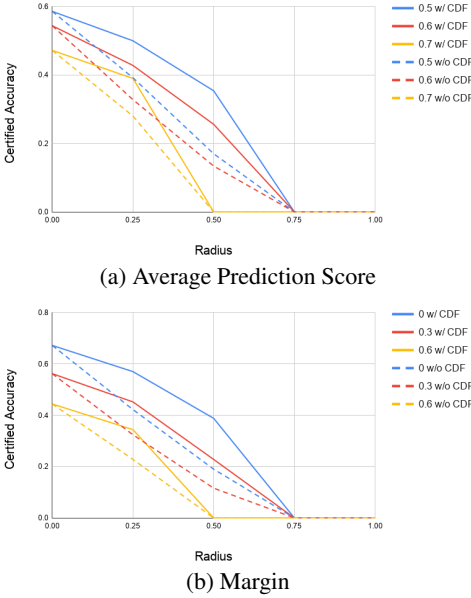


Figure 1: ImageNet Results ( $\sigma = 0.25$ ).