

1 We thank the reviewers for the thought-invoking questions and helpful comments on improving the manuscript. We will
2 add more text to improve readability and, if necessary, remove some of the technical lemmas to make room. We also
3 have extended our empirical proof of maximal informativeness to $k = 15$.

5 **R1, R2, & R3:** The LLW hinge loss is calibrated with respect to the 0-1 loss while the WW hinge loss is not. Why
6 does the WW SVM still outperform the LLW SVM? In other words, how can calibration be used as a justification for
7 performance?

8 **A:** The LLW SVM performs worse for a reason unrelated to calibration. Doğan et al. [2016] makes the distinction
9 between *relative* and *absolute* margin losses. The WW and CS loss are both based on relative margins, while LLW is
10 based on absolute margins. Doğan et al. [2016] on their page 20 gave an explanation for the worse performance of all
11 losses based on absolute margin. Hence, the poor performance of LLW is a consequence of using absolute margin. Out
12 of the nine SVM losses considered by Doğan et al. [2016], only the CS and WW losses are relative margin based. We
13 will add this discussion to our manuscript.

15 **R2, R3 & R4:** Why is consistency with respect to the ordered partition loss desirable? What is the intuition?

16 **A:** Regarding the intuition behind the ordered partition loss, the basic idea is that we want to rank the labels,
17 where ties are allowed and each S_i is a set of labels that are tied. We want the correct label to be as high
18 up the ranking as possible. The lower the true class is ranked, the larger the loss. That is what the definition of
19 the ordered partition loss says. We should have said this in the initial draft and will add this discussion to our manuscript.

21 **R1 & R4:** How to get the surrogate decision function ψ to recover the ordered partition/buckets? Is there an excess
22 risk bound?

23 **A:** In line 125 of our manuscript, we cited Finocchiaro et al. [2019] who provided an explicit ψ given L , ℓ , and φ . We
24 refer to [Finocchiaro et al., 2019, Definition 6] for the construction of ψ . The excess risk bound for their constructed
25 ψ can be found in [Finocchiaro et al., 2019, Theorem 6]. We will make the theorem references explicit in the manuscript.

27 **R2:** What are the consequences of the maximally informative property?

28 **A:** Intuitively, a discrete loss $\ell : \mathcal{R} \rightarrow \mathbb{R}_+^k$ (where \mathcal{R} is finite) with embedding φ (an injection into the domain
29 of L) is maximally informative for a surrogate L if $\varphi(\mathcal{R})$ captures all the essential information contained in the
30 surrogate L in the most compact way. To better convey this intuition, we replace “maximally informative” with the new
31 terminology “minimally emblematic.” Let us say that a set of vectors $E \subseteq \mathbb{R}^k$ is an *emblem* of L if for all $p \in \Delta^k$, the
32 set $E \cap \operatorname{argmin}_v \langle p, L(v) \rangle$ is nonempty. Then we can equivalently define ℓ with φ to be *minimally emblematic* for L if
33 $\varphi(\mathcal{R})$ is an emblem of L of minimal cardinality. In other words, $\varphi(\mathcal{R})$ is a minimal set of minimizers of all possible
34 L -inner risks. We will update our manuscript with this discussion and the new terminology.

36 **R3:** How does performance over the ordered partition loss translate to the 0-1 loss?

37 **A:** Results from our section on the “argmax link” provide a partial answer to this. Namely, we show in two common
38 regimes, the Bayes optimal ordered partition has a top bucket consisting of a single element. When this occurs, the
39 argmax link recovers the most probable class, i.e., the unique element from the top bucket. We will modify the
40 manuscript to clarify this point.

41 References

42 U. Doğan, T. Glasmachers, and C. Igel. A unified view on multi-class support vector classification. *The Journal of*
43 *Machine Learning Research*, 17(1):1550–1831, 2016.

44 J. Finocchiaro, R. Frongillo, and B. Waggoner. An embedding framework for consistent polyhedral surrogates. In
45 *Advances in neural information processing systems*, pages 10780–10790, 2019.