

Consistency of surrogate risk minimization methods for multiclass 0-1 classification

Lecturer: Shivani Agarwal

Scribe: Debarghya Ghoshdastidar

1 Introduction

In the previous two lectures, we learned about consistency of surrogate risk minimization methods for binary classification. In this lecture, we will study how the previous results extend to a multiclass setting. As examples, we will study consistency of some multiclass SVM formulations.

We begin by describing how the surrogate risk minimization framework changes for multiclass classification. This is shown in Table 1. For the purpose of this lecture, we will restrict ourselves to the multiclass 0-1 loss.

	Binary classification	Multiclass classification
Label space, \mathcal{Y} and Prediction space, \mathcal{T}	$\{\pm 1\}$	$[n]$
Target 0-1 loss $\ell_{0-1} : \mathcal{Y} \times \mathcal{T} \mapsto \mathbb{R}_+$	$\ell_{0-1} : \{\pm 1\} \times \{\pm 1\} \mapsto \mathbb{R}_+$ $\ell_{0-1}(y, t) = \mathbf{1}(t \neq y)$	$\ell_{0-1} : [n] \times [n] \mapsto \mathbb{R}_+$ $\ell_{0-1}(y, t) = \mathbf{1}(t \neq y)$
Surrogate loss $\psi : \mathcal{Y} \times \mathcal{C} \mapsto \mathbb{R}_+$	$\psi : \{\pm 1\} \times \mathcal{C} \mapsto \mathbb{R}_+$ where $\mathcal{C} \subseteq \mathbb{R}$	$\psi : [n] \times \mathcal{C} \mapsto \mathbb{R}_+$ where $\mathcal{C} \subseteq \mathbb{R}^n$
'pred' function $\text{pred} : \mathcal{C} \mapsto \mathcal{T}$	$\text{pred} : \mathcal{C} \mapsto \{\pm 1\}$ $\text{pred}(\alpha) = \text{sign}(\alpha)$	$\text{pred} : \mathcal{C} \mapsto [n]$ $\text{pred}(\mathbf{u}) = \arg \max_{y \in [n]} u_y$
Conditional ψ -risk	$L_\psi : [0, 1] \times \mathcal{C} \mapsto \mathbb{R}_+$ $L_\psi(\eta, \alpha) = \eta\psi(1, \alpha) + (1 - \eta)\psi(-1, \alpha)$	$L_\psi : \Delta_n \times \mathcal{C} \mapsto \mathbb{R}_+$ $L_\psi(\mathbf{p}, \mathbf{u}) = \sum_{y=1}^n p_y \psi(y, \mathbf{u})$
Conditional Bayes ψ -risk	$H_\psi : [0, 1] \mapsto \mathbb{R}_+$ $H_\psi(\eta) = \inf_{\alpha \in \mathcal{C}} L_\psi(\eta, \alpha)$	$H_\psi : \Delta_n \mapsto \mathbb{R}_+$ $H_\psi(\mathbf{p}) = \inf_{\mathbf{u} \in \mathcal{C}} L_\psi(\mathbf{p}, \mathbf{u})$
Conditional ψ -regret	$R_\psi : [0, 1] \times \mathcal{C} \mapsto \mathbb{R}_+$ $R_\psi(\eta, \alpha) = L_\psi(\eta, \alpha) - H_\psi(\eta)$	$R_\psi : \Delta_n \times \mathcal{C} \mapsto \mathbb{R}_+$ $R_\psi(\mathbf{p}, \mathbf{u}) = L_\psi(\mathbf{p}, \mathbf{u}) - H_\psi(\mathbf{p})$

Table 1: Surrogate risk minimization framework for binary and multiclass classification. We denote the y^{th} coordinate of vectors \mathbf{u} and \mathbf{p} by u_y and p_y , respectively. Here, $p_y = \mathbf{P}(x \in \text{Class } y|x)$.

Example 1 (Multiclass logistic regression). The loss function in multiclass logistic regression is given by the surrogate loss $\psi_{\log} : [n] \times \mathbb{R}^n \mapsto \mathbb{R}_+$ defined as

$$\psi_{\log}(y, \mathbf{u}) = -\ln \left(\frac{e^{-u_y}}{\sum_{y'=1}^n e^{-u_{y'}}} \right) \quad \forall y \in [n] \quad \forall \mathbf{u} = (u_1, \dots, u_n) \in \mathbb{R}^n.$$

One can verify that minimizing ψ_{\log} is equivalent to maximizing the softmax function, which can be used as parametric form for conditional class probability. It is also known that ψ_{\log} is a multiclass proper composite loss [4].

2 Classification calibration in multiclass setting

We will now see how the definition of ℓ_{0-1} -calibrated loss function extends to the multiclass framework. We can also state a result on multiclass ℓ_{0-1} -calibrated losses similar to the one in previous lecture. This is given in Theorem 2.1.

Definition. For $\mathcal{C} \subseteq \mathbb{R}^n$, a surrogate loss function $\psi : [n] \times \mathcal{C} \mapsto \mathbb{R}^n$ is said to be ℓ_{0-1} -calibrated if

$$\forall \mathbf{p} \in \Delta^n, \quad \inf_{\mathbf{u} \in \mathcal{C}: \text{pred}(\mathbf{u}) \notin \arg \min_{t \in [n]} L_{0-1}(\mathbf{p}, t)} L_\psi(\mathbf{p}, \mathbf{u}) > \inf_{\mathbf{u} \in \mathcal{C}} L_\psi(\mathbf{p}, \mathbf{u}),$$

or equivalently,

$$\forall \mathbf{p} \in \Delta^n, \quad \inf_{\mathbf{u} \in \mathcal{C}: R_{0-1}(\mathbf{p}, \text{pred}(\mathbf{u})) > 0} R_\psi(\mathbf{p}, \mathbf{u}) > 0.$$

Theorem 2.1. Let $\mathcal{C} \subseteq \mathbb{R}^n$. Let $\psi : [n] \times \mathcal{C} \mapsto \mathbb{R}_+$ be such that $\psi_y(\mathbf{u}) := \psi(y, \mathbf{u})$ is continuous $\forall y \in [n]$. Then the following are equivalent:

1. ψ is ℓ_{0-1} -calibrated.
2. for all distributions D on $\mathcal{X} \times [n]$,

$$\text{regret}_D^\psi[f_S] \xrightarrow{P} 0 \quad \implies \quad \text{regret}_D^{0-1}[\text{pred} \circ f_S] \xrightarrow{P} 0.$$

3. $\exists g : \mathbb{R}_+ \mapsto \mathbb{R}_+$ such that g is continuous at 0, $g(0) = 0$ and $\forall D$ on $\mathcal{X} \times [n]$, $\forall f : \mathcal{X} \mapsto \mathcal{C}$,

$$\text{regret}_D^{0-1}[\text{pred} \circ f] \leq g\left(\text{regret}_D^\psi[f]\right).$$

PROOF. As in the previous lecture, we will skip the proof for (2) \implies (1), whereas (3) \implies (2) is obvious. For the case of (1) \implies (3), we note that the proof is similar to the corresponding proof in previous lecture. The key step in that proof was a result on positivity of the uniform calibration function. We will state and prove a generalized version of that lemma, which can be directly applied to prove that (1) \implies (3) as in previous lecture. \square

We restate the definition of calibration function in the multiclass setting.

Definition. The calibration function at $\mathbf{p} \in \Delta_n$ is the function $\delta_{\mathbf{p}} : [0, 1] \mapsto \mathbb{R}_+$ defined as

$$\delta_{\mathbf{p}}(\epsilon) = \inf_{\mathbf{u} \in \mathcal{C}: R_{0-1}(\mathbf{p}, \text{pred}(\mathbf{u})) \geq \epsilon} R_\psi(\mathbf{p}, \mathbf{u}).$$

The uniform calibration function $\delta : [0, 1] \mapsto \mathbb{R}_+$ is given by $\delta(\epsilon) = \inf_{\mathbf{p} \in \Delta_n} \delta_{\mathbf{p}}(\epsilon)$.

Note that $\delta_{\mathbf{p}}(0) = 0 \forall \mathbf{p} \in \Delta_n$, and hence, $\delta(0) = 0$. Also, from the definition of classification calibration, we have that

$$\psi \text{ is } \ell_{0-1}\text{-calibrated} \quad \iff \quad \delta_{\mathbf{p}}(\epsilon) > 0 \quad \forall \epsilon > 0, \forall \mathbf{p} \in \Delta_n.$$

Hence, we can say that $\delta(\epsilon) \geq 0 \quad \forall \epsilon > 0$. We now state the result by Zhang [6], which guarantees positivity of $\delta(\epsilon) \forall \epsilon > 0$. Before proving this result, recall the following property of conditional Bayes risk.

Exercise. The function $H_\psi(\mathbf{p})$ is a continuous function of \mathbf{p} on Δ_n .

Lemma 2.2. Let $\psi : [n] \times \mathcal{C} \mapsto \mathbb{R}_+$ be such that $\psi_y(\mathbf{u})$ is continuous $\forall y \in [n]$. If ψ is ℓ_{0-1} -calibrated, then $\delta(\epsilon) > 0 \forall \epsilon > 0$.

PROOF. Let ψ be ℓ_{0-1} -calibrated. We will prove the claim by contradiction. Let, if possible, $\exists \epsilon > 0$ such that $\delta(\epsilon) = 0$, *i.e.*,

$$\inf_{\mathbf{p} \in \Delta_n} \inf_{\mathbf{u} \in \mathcal{C}: R_{0-1}(\mathbf{p}, \text{pred}(\mathbf{u})) \geq \epsilon} R_\psi(\mathbf{p}, \mathbf{u}) = 0.$$

Note that $R_\psi(\mathbf{p}, \mathbf{u})$ is continuous both in \mathbf{p} and \mathbf{u} . Hence, there exists a sequence $(\mathbf{p}^{(m)}, \mathbf{u}^{(m)})_{m \in \mathbb{N}}$ in $\Delta_n \times \mathcal{C}$ such that

$$R_{0-1}(\mathbf{p}^{(m)}, \text{pred}(\mathbf{u}^{(m)})) \geq \epsilon \quad \forall m \quad (1)$$

$$\text{and } \lim_{m \rightarrow \infty} R_\psi(\mathbf{p}^{(m)}, \mathbf{u}^{(m)}) = 0, \text{ i.e., } \lim_{m \rightarrow \infty} (L_\psi(\mathbf{p}^{(m)}, \mathbf{u}^{(m)}) - H_\psi(\mathbf{p}^{(m)})) = 0. \quad (2)$$

Since Δ_n is compact, \exists a convergent subsequence of $(\mathbf{p}^{(m)})$. Let \mathbf{p} be the limit of that subsequence. We will further work only with this subsequence, and so, for convenience, we denote the subsequence also by the same index m , i.e., $\mathbf{p}^{(m)} \rightarrow \mathbf{p} \in \Delta_n$.

By continuity of H_ψ , we have $H_\psi(\mathbf{p}^{(m)}) \rightarrow H_\psi(\mathbf{p})$ and so,

$$\lim_{m \rightarrow \infty} L_\psi(\mathbf{p}^{(m)}, \mathbf{u}^{(m)}) = H_\psi(\mathbf{p}). \quad (3)$$

Assume without loss of generality that $\mathbf{p} = (p_1, \dots, p_n)$ is such that $p_y = 0$ for $y = 1, \dots, k$ and $p_y > 0$ for $y = k+1, \dots, n$. We have

$$\begin{aligned} \limsup_{m \rightarrow \infty} L_\psi(\mathbf{p}, \mathbf{u}^{(m)}) &= \limsup_{m \rightarrow \infty} \sum_{y=k+1}^n p_y \psi_y(\mathbf{u}^{(m)}) \\ &= \limsup_{m \rightarrow \infty} \sum_{y=k+1}^n p_y^{(m)} \psi_y(\mathbf{u}^{(m)}) && \text{(since } \mathbf{p}^{(m)} \rightarrow \mathbf{p}) \\ &\leq \limsup_{m \rightarrow \infty} \sum_{y=1}^n p_y^{(m)} \psi_y(\mathbf{u}^{(m)}) \\ &= \lim_{m \rightarrow \infty} L_\psi(\mathbf{p}^{(m)}, \mathbf{u}^{(m)}) = H_\psi(\mathbf{p}) && \text{(from Eq. (3)).} \end{aligned}$$

On the other hand, from definition of $H_\psi(\mathbf{p})$, it follows that

$$\liminf_{m \rightarrow \infty} L_\psi(\mathbf{p}, \mathbf{u}^{(m)}) \geq H_\psi(\mathbf{p}).$$

Hence, $\lim_{m \rightarrow \infty} L_\psi(\mathbf{p}, \mathbf{u}^{(m)}) = H_\psi(\mathbf{p})$, i.e., $\lim_{m \rightarrow \infty} R_\psi(\mathbf{p}, \mathbf{u}^{(m)}) = 0$.

But from Eq. (1) and the fact $\mathbf{p}^{(m)} \rightarrow \mathbf{p}$, we can derive that

$$R_{0-1}(\mathbf{p}, \text{pred}(\mathbf{u}^{(m)})) \geq \epsilon \quad \forall m.$$

Thus, we have $\delta_{\mathbf{p}}(\epsilon) = \inf_{\mathbf{u} \in \mathcal{C}: R_{0-1}(\mathbf{p}, \text{pred}(\mathbf{u})) \geq \epsilon} R_\psi(\mathbf{p}, \mathbf{u}) = 0$, which contradicts the fact that $\epsilon > 0$.

Hence, the claim is proved. \square

In the above results (Theorem 2.1 and Lemma 2.2), the assumption on continuity of ψ_y will not be required if the loss is a margin based binary loss. This follows from the exact expression that can be obtained for margin based losses (as discussed in the previous lecture).

Exercise. We have proved Theorem 2.1 and Lemma 2.2 based on the assumption that $\mathcal{Y} = \mathcal{T} = [n]$ is finite. Identify which of the statements in the proofs make use of this assumption. Can we modify the proofs such that the results hold for $\mathcal{Y} = \mathcal{T} = [0, 1]$?

3 Consistency of some multiclass SVM formulations

Theorem 2.1 can be used to study consistency of different surrogate risk minimization methods for multiclass 0-1 classification. In particular, we may study consistency of various multiclass formulations for SVM that have been proposed in the literature [5, 1, 2]. This has been studied in [6]. In all the multiclass SVM formulations, a multiclass extension of the hinge loss is minimized. However, one may consider any margin function ϕ in general as used below. In particular, $\phi(z) = (1 - z)_+$ would correspond to the algorithms in [5, 1, 2].

3.1 Formulation by Weston & Watkins [5]

$\mathcal{C} = \mathbb{R}^n$ and surrogate loss $\psi_{WW} : [n] \times \mathbb{R}^n \mapsto \mathbb{R}_+$ is

$$\psi_{WW}(y, \mathbf{u}) = \sum_{y' \neq y} \phi(u_y - u_{y'}) \quad \text{for some } \phi : \mathbb{R} \mapsto \mathbb{R}_+$$

Theorem 3.1. Let ϕ be a decreasing function that is differentiable everywhere and $\phi'(0) < 0$. Then ψ_{WW} is ℓ_{0-1} -calibrated.

Though the above includes exponential and logistic losses, it does not include the hinge loss. In fact Zhang [6] showed a counter-example in case of hinge loss. Thus ψ_{WW} with the hinge loss is not ℓ_{0-1} -calibrated.

3.2 Formulation by Crammer & Singer [1]

$\mathcal{C} = \mathbb{R}^n$ and surrogate loss $\psi_{CS} : [n] \times \mathbb{R}^n \mapsto \mathbb{R}_+$ is

$$\psi_{CS}(y, \mathbf{u}) = \phi \left(u_y - \max_{y' \neq y} u_{y'} \right) \quad \text{for some } \phi : \mathbb{R} \mapsto \mathbb{R}_+$$

This formulation is popularly used in structured prediction. However, for convex ϕ , in general ψ_{CS} is not ℓ_{0-1} -calibrated.

3.3 Formulation by Lee, Lin & Wahba [2]

$\mathcal{C} = \left\{ \mathbf{u} \in \mathbb{R}^n : \sum_{y=1}^n u_y = 0 \right\}$ and surrogate loss $\psi_{LLW} : [n] \times \mathcal{C} \mapsto \mathbb{R}_+$ is

$$\psi_{LLW}(y, \mathbf{u}) = \sum_{y' \neq y} \phi(-u_{y'}) \quad \text{for some } \phi : \mathbb{R} \mapsto \mathbb{R}_+$$

Theorem 3.2. Let ϕ be a convex function that is differentiable on $(-\infty, 0]$ and $\phi'(0) < 0$. Then ψ_{LLW} is ℓ_{0-1} -calibrated.

This result includes the hinge loss, and therefore this yields a universally consistent formulation of multiclass SVMs.

3.4 Formulation using one vs. all approach

$\mathcal{C} = \mathbb{R}^n$ and surrogate loss $\psi_{OvA} : [n] \times \mathbb{R}^n \mapsto \mathbb{R}_+$ is

$$\psi_{OvA}(y, \mathbf{u}) = \phi(u_y) + \sum_{y' \neq y} \phi(-u_{y'}) \quad \text{for some } \phi : \mathbb{R} \mapsto \mathbb{R}_+$$

Exercise. Show that one vs. all approach effectively minimizes the above loss on the training sample.

Theorem 3.3. Let ϕ be a convex function that is differentiable everywhere and $\phi(z) < \phi(-z) \forall z > 0$. Then ψ_{OvA} is ℓ_{0-1} -calibrated.

This result also does not hold for the hinge loss.

4 Consistency over a class of distributions

The above discussion shows that most of the multiclass SVM formulations (in Sections 3.1,3.2 and 3.4) are not consistent. However, in practice, they are found to perform quite well. This is due to the fact that these losses are ℓ_{0-1} -calibrated under certain conditions on the underlying distributions. This can be formalized through the following definition.

Definition. For $\mathcal{P} \subseteq \Delta_n$ and $\mathcal{C} \subseteq \mathbb{R}^n$, a surrogate loss function $\psi : [n] \times \mathcal{C} \mapsto \mathbb{R}^n$ is said to be ℓ_{0-1} -calibrated over \mathcal{P} if

$$\forall \mathbf{p} \in \mathcal{P}, \quad \inf_{\mathbf{u} \in \mathcal{C}: R_{0-1}(\mathbf{p}, \text{pred}(\mathbf{u})) > 0} R_{\psi}(\mathbf{p}, \mathbf{u}) > 0.$$

Based on this definition, we can redefine the calibration and uniform calibration functions. The statement of Theorem 2.1 holds for all distributions D on $\mathcal{X} \times [n]$ such that $D_{y|x} \in \mathcal{P}$, *i.e.*, we do not get universal consistency, but consistency on \mathcal{P} . In Sections 3.1,3.2 and 3.4, the formulation with hinge loss is consistent over $\mathcal{P} = \left\{ \mathbf{p} \in \Delta_n : \max_{y \in [n]} p_y > \frac{1}{2} \right\}$, *i.e.*, when one of the classes strictly dominates over others.

5 Additional pointers

Most of the results discussed in this lecture can be found in [6]. An alternative approach was studied in [3] using geometric interpretations. The authors of this paper also study if consistency is possible when the ‘pred’ function is not of ‘argmax’ type. One of their observations is that if $\psi : [n] \times \mathcal{C} \mapsto \mathbb{R}_+$ is a symmetric surrogate loss, then it suffices to show only the 0-1 calibration for argmax pred function. Here, the term ‘symmetric’ loss means that for all $\mathbf{u} \in \mathcal{C}$ and permutations $\sigma : [n] \rightarrow [n]$, $\exists \mathbf{u}' \in \mathcal{C}$ such that $\psi(\sigma(y), \mathbf{u}) = \psi(y, \mathbf{u}') \forall y \in [n]$.

We had seen in last class that better bounds are achieved using classification calibration than using proper losses. This apparently seems counter-intuitive as stronger assumptions are used in the latter case. However, in case of proper losses, we tried to estimate distribution (given by η) which implied that we tried to learn more than just the classes of data. Hence, the bounds were looser than the case of classification calibration, where we try to achieve consistency only with respect to problem of learning the classes. In fact, it is known that the classification calibration based bounds are the tightest in the distribution-free setting.

References

- [1] K. Crammer and Y. Singer. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research*, 2:265–292, 2001.
- [2] Y. Lee, Y. Lin, and G. Wahba. Multicategory support vector machines: Theory and application to the classification of microarray data and satellite radiance data. *Journal of the American Statistical Association*, 99(465):67–81, 2004.
- [3] A. Tewari and P. L. Bartlett. On the consistency of multiclass classification methods. *Journal of Machine Learning Research*, 8:1007–1025, 2007.
- [4] E. Vernet, R. C. Williamson, and M. Reid. Composite multiclass losses. In *Advances in Neural Information Processing Systems*, 2011.
- [5] J. Weston and C. Watkins. Multi-class support vector machines. Technical report, CSD-TR-98-04, Department of Computer Science, Royal Holloway, University of London, 1998.
- [6] T. Zhang. Statistical analysis of some multi-category large margin classification methods. *Journal of Machine Learning Research*, 5:1225–1251, 2004.