

On Structured Prediction Theory with Calibrated Convex Surrogate Losses

Rushab Munot

University of Pennsylvania

rushab@seas.upenn.edu

Instructor: Shivani Agarwal

Authors

Anton Osokin

Francis Bach

Simon Lacoste-Julien

October 30, 2018

- Structured Prediction - Obstacles
 - Exponential number of classes
 - Cost-Sensitive Task Loss
- Previous Work
 - Consistent but not efficient (non-convex)
 - Efficient (convex) but not consistent
- What we need
 - Consistency for surrogate losses **and**
 - Efficient algorithm

- Convexity + Consistency $\not\Rightarrow$ Efficient Algorithm
 - Exponential constant in Generalization Error bound (eg. Ciliberto et. al. 2010)
 - To have error ϵ on Task Loss
 - 1 Exponentially small error on surrogate loss
 - 2 Exponential number of iterations
- **Approach**
 - **Calibration function:** Connect actual excess risk and surrogate excess risk
 - **Exponential constants:** Constrain the score vector space
 - **Online SGD**

Recap - Calibrated Surrogates for 0-1 Loss

- Inner Risk: $L_I(\eta, \alpha) = E_{y \sim \eta}[l(y, \alpha)]$
- Bayesian Inner Risk: $H_I(\eta) = \inf_{\alpha} L_I(\eta, \alpha)$
- Inner Regret/Excess Risk: $\mathcal{R}_I(\eta, \alpha) = L_I(\eta, \alpha) - H_I(\eta)$

- **0-1 Calibrated Surrogate Loss** ($\psi : [0, 1] \rightarrow \mathbb{R}$):

$$\inf_{R_{0-1}(\eta, \text{sign}(f)) > 0} R_{\psi}(\eta, f) > 0 \quad \dots \quad \forall \eta \in [0, 1]$$

- Bartlett et al (2006): Calibrated surrogate losses admit a regret transfer bound w.r.t. the 0-1 loss
- Calibrated Multi-class losses

- $\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$ (structured)
 $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$
- $|\mathcal{Y}| = k$ (finite, exponential in size of \mathbf{y})
- Loss function $L \in \mathbb{R}^{k \times k}$ (non-negative)
- Score function: $\mathbf{f} : \mathcal{X} \rightarrow \mathcal{F} \subseteq \mathbb{R}^k$
 $\text{pred}(\mathbf{f}(\mathbf{x})) = \text{argmax}_{\hat{\mathbf{y}} \in \mathcal{Y}} \mathbf{f}_{\hat{\mathbf{y}}}(\mathbf{x})$
- $\mathfrak{F}_{\mathcal{F}} =$ Set of score functions
- Surrogate Loss $\Phi : \mathbb{R}^k \times \mathcal{Y} \rightarrow \mathbb{R}$ (continuous and bounded below)

- Generalization Error:

$$\begin{aligned}er_{\mathcal{D}}^L &= E_{(x,y) \sim \mathcal{D}}[L(\text{pred}(f(x)), y)] \\&= E_{x \sim \mathcal{D}_{\mathcal{X}}}[l(\text{pred}(f(x)), P(\cdot|x))]\end{aligned}$$
$$l(\mathbf{f}, \mathbf{q}) = \sum_{c=1}^k q_c L(\text{pred}(f), c) \dots \text{Inner Risk}$$

- Surrogate Error:

$$\begin{aligned}er_{\mathcal{D}}^{\Phi} &= E_{(x,y) \sim \mathcal{D}}[\Phi(f(x), y)]\end{aligned}$$
$$\phi(\mathbf{f}, \mathbf{q}) = \sum_{c=1}^k q_c \Phi(f, c) \dots \text{Surrogate Inner Risk}$$

- R_l and R_{ϕ} denote the inner regret and surrogate inner regret

Loss functions under consideration

- 01 Loss
- Hamming Loss
- Block loss:
Divide into b blocks of size s
 $k = bs$

- Mixed Loss:
$$L_{mixed,01,b}(\hat{y}, y) = \eta L_{01}(f, q)(\hat{y}, y) + (1 - \eta) L_b(\hat{y}, y)$$

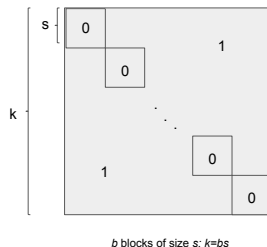


Figure: Block 0-1 Loss

Calibration Function

- Connects actual and surrogate inner regrets

Definition (Calibration Function)

For a task loss L , surrogate loss Φ , and space of allowed score vectors \mathcal{F} ,

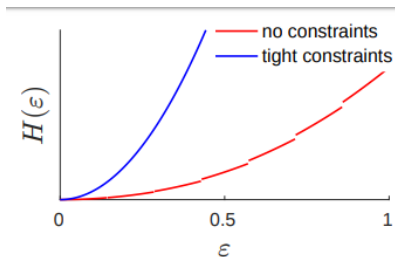
$$\mathcal{H}_{L,\Phi,\mathcal{F}}(\epsilon) = \inf_{\mathbf{f} \in \mathcal{F}, \mathbf{q} \in \Delta_k} R_\Phi(\mathbf{f}, \mathbf{q})$$

s.t. $R_L(\mathbf{f}, \mathbf{q}) > \epsilon$

- Implications:

- 1 $\mathcal{H}_{L,\Phi,\mathcal{F}}(R_L(\mathbf{f}, \mathbf{q})) \leq R_\Phi(\mathbf{f}, \mathbf{q})$
- 2 $\mathcal{H}_{L,\Phi,\mathcal{F}}$ is non-decreasing
- 3 Larger $\mathcal{H}_{L,\Phi,\mathcal{F}}(\epsilon)$ is better

Visualizing the Calibration Function

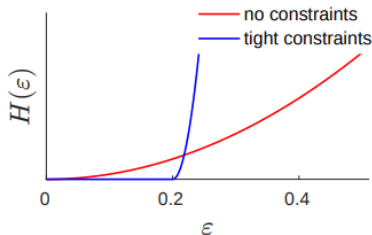


(a)

Surrogate Loss: Quadratic

Task Loss: Hamming

Constraints: $\mathcal{F} = \text{span}(L_{\text{Ham}, T})$



(b)

Surrogate Loss: Quadratic

Task Loss: Mixed $L_{01,b,\eta}$

Constraints: $\mathcal{F} = \text{span}(L_{01,b,\eta})$

Figure: Exponential rise in $\mathcal{H}_{L,\phi,\mathcal{F}}(\epsilon)$ with constraints. $\mathcal{H}_{L,\phi,\mathcal{F}}$ can be inconsistent/ not continuous

Consistency

Note that, $\forall \epsilon, \mathcal{H}_{L,\Phi,\mathcal{F}}(\epsilon) > 0 \implies$ consistency

Theorem

$\check{\mathcal{H}}_{L,\Phi,\mathcal{F}}$ is the lower convex envelope of $\mathcal{H}_{L,\Phi,\mathcal{F}}$. Then

$$\text{regret}_D^\Phi[f] < \check{\mathcal{H}}_{L,\Phi,\mathcal{F}}(\epsilon) \implies \text{regret}_D^L < \epsilon$$

Definition (η consistency)

Φ is η -consistent iff

- 1 $\mathcal{H}_{L,\Phi,\mathcal{F}}(\epsilon) > 0, \forall \epsilon > \eta$
- 2 $\mathcal{H}_{L,\Phi,\mathcal{F}}$ is finite for some $\hat{\epsilon} > 0$

Allows optimization up to a certain accuracy. Can be much faster (Mixed 01 block loss).

- **Approximation Error:**

$$er_{\mathcal{D}}^{L, \mathcal{F}, *} - er_{\mathcal{D}}^{L, *} = \inf_{\mathbf{f} \in \mathcal{F}} er_{\mathcal{D}}^L[\mathbf{f}] - \inf_{\mathbf{f}} er_{\mathcal{D}}^L[\mathbf{f}]$$

- Given $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$, $\mathbf{q} = P(\cdot | \mathbf{x}) \in \mathbb{R}^k$ and $L \in \mathbb{R}^{k \times k}$
 - Expected loss if you predict $y = c$ is $(L\mathbf{q})_c$
 - Observe that $\mathbf{f} = -L\mathbf{q}$ is optimal
- $\text{span}(L) \subseteq \mathcal{F} \implies 0$ approximation error
- Can restrict $\mathcal{F} = \text{span}(L)$
- Low Rank L may get rid of exponential constants

Quadratic Surrogate Loss

Definition

$\mathbf{f} \in \mathcal{F}, \mathbf{y} \in \mathcal{Y}$

$$\Phi_{quad}(\mathbf{f}, \mathbf{y}) = \frac{1}{2k} \|\mathbf{f} + L(:, \mathbf{y})\|^2 = \frac{1}{2k} \sum_{c=1}^k \left(\mathbf{f}_c + L(c, \mathbf{y}) \right)^2$$

- Let $F \in \mathbb{R}^{k \times r}$ (e.g. F is the basis of $\text{span}(L)$)
- Constrain $\mathcal{F} = \text{span}(F) = \{F\boldsymbol{\theta} \mid \boldsymbol{\theta} \in \mathbb{R}^r\}$
- $\dim(\mathcal{F}) = \text{rank}(F) \leq r \ll k$
- Surrogate inner regret: $R_\phi(\mathbf{f} = F\boldsymbol{\theta}) = \frac{1}{2k} \|F\boldsymbol{\theta} + L\mathbf{q}\|^2$ (Convex in \mathbf{f} and \mathbf{q})

Lower Bound on $\mathcal{H}_{\Phi_{quad}, L, \mathcal{F}}$

- i = Label with lowest expected loss
 j = Label with highest predicted score

$$\mathcal{H}_{ij}(\epsilon) = \inf_{i, j \in \text{pred}(\mathcal{F})} R_{\phi}(f, q)$$

s.t.

$$(Lq)_i \leq (Lq)_j - \epsilon$$

$$(Lq)_i \leq (Lq)_c \quad \forall c \in \text{pred}(\mathcal{F})$$

$$f_j \geq f_c \quad \forall c \in \text{pred}(\mathcal{F})$$

$$\mathbf{f} \in \mathcal{F}$$

$$\mathbf{q} \in \Delta_k$$

- The union of feasibility sets is $\mathcal{F} \times \Delta_k$
- $\mathcal{H}_{L, \Phi, \mathcal{F}}(\epsilon) = \min_{i, j \in \text{pred}(\mathcal{F}); i \neq j} \mathcal{H}_{ij}(\epsilon)$

Lower Bound on $\mathcal{H}_{\Phi_{quad}, L, \mathcal{F}}$

- For the quadratic surrogate, this translates to

$$\begin{aligned}\mathcal{H}_{ij}(\epsilon) &= \inf_{\theta, \mathbf{q}} \|F\theta + L\mathbf{q}\|^2 \\ (\mathbf{e}_i - \mathbf{e}_j)^T L\mathbf{q} &= -\epsilon \\ (\mathbf{e}_i - \mathbf{e}_j)^T \theta &\leq 0\end{aligned}$$

Can prove that

$$\mathcal{H}_{\Phi_{quad}, L, \mathcal{F}} \geq \frac{\epsilon^2}{2k} \left(\frac{1}{\max_{i \neq j} \|P_{\mathcal{F}}(\mathbf{e}_i - \mathbf{e}_j)\|_2^2} \right) \geq \frac{\epsilon^2}{4k}$$

$P_{\mathcal{F}}$: Operator for projection onto \mathcal{F}

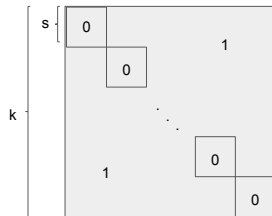
\mathbf{e}_c : c^{th} standard basis vector in \mathbb{R}^k

Task Loss: 0-1

- $\mathcal{H}_{L_{01}, \Phi_{quad}, \mathcal{F}}(\epsilon) \geq \frac{\epsilon^2}{4k}$
- This bound is tight
- Note that $\mathcal{F} = \text{span}(L_{01}) = \mathbb{R}^k$
Hence the projection is always $e_i - e_j$
- k is exponential in the dimension of \mathbf{y}
- 0-1 Loss has poor guarantees as expected

Block 0-1 Loss

- $\mathcal{F} = \text{span}(L_{\text{block},b})$
- $\text{rank}(L) = b$
- $\|P_{\mathcal{F}}(\mathbf{e}_i - \mathbf{e}_j)\|_2^2 = \frac{2}{s} \dots i \neq j, \text{ else } 0$
- $\mathcal{H}_{L_{\text{block},T}, \Phi_{\text{quad}}, \mathcal{F}}(\epsilon) \geq \frac{\epsilon^2}{4b}$ (not tight)



b blocks of size $s: k=bs$

Figure: Block 0-1 Loss

Hamming Loss

Hamming Loss: Fraction of wrong labels

- T = Sequence Length
- $\mathcal{F} = \text{span}(L_{Ham, T})$
- $\text{rank}(L) = T + 1$
- $\mathcal{H}_{L_{Ham, T}, \Phi_{quad}, \mathcal{F}}(\epsilon) \geq \frac{\epsilon^2}{8T}$ (tight)

Upper Bound on $\mathcal{H}_{\Phi_{quad}, L, \mathcal{F}}$

- Turns out we cannot give good guarantees if \mathcal{F} is not constrained

Theorem

For a Loss Matrix L which is a pseudometric, and unconstrained \mathcal{F} (i.e. $\mathcal{F} = \mathbb{R}^k$),

$$\mathcal{H}_{\Phi_{quad}, L, \mathcal{F}}(\epsilon) < \frac{\epsilon^2}{2k}$$

- This is why we need to constrain the score vector space.
- Also this bound is the reason why previous algorithms were not efficient.

- Online (kernel) Projected Averaged SGD
- $F \in \mathbb{R}^{k \times r}$ $\psi : \mathcal{X} \rightarrow \mathbb{R}^d$ Feature Map; $W \in \mathbb{R}^{r \times d}$
- $\mathbf{f}(x) = FW\psi(x)$
- SGD Update for $(\mathbf{x}^n, \mathbf{y}^n) \sim \mathcal{D}$:

$$W^{(n)} = P_D[W^{(n-1)} - \gamma^{(n)} F^T \nabla \Phi \psi(\mathbf{x}^n)^T]$$

P_D is the projection on the ball of radius D wrt the Hilbert Schmidt norm. γ is the step size. (No need to calculate the feature map - kernel trick)

- **Assumption:** The surrogate regret $er_{\mathcal{D}}^{\Phi}$ has a global minimum \mathbf{f}^* over the function class $\mathfrak{F}_{\mathcal{F}} = \{\mathbf{f} \mid \mathbf{f}(\mathbf{x}) = FW\psi(\mathbf{x})\}$.
Not required - but analysis becomes complicated!

- If we have

- 1 $\Phi(f, y)$ is bounded below and convex wrt f for all y
- 2 $\|F^T \nabla \Phi \psi(\mathbf{x}^n)^T\|_{HS}^2 \leq M^2$
- 3 $\|W^*\|_{HS} < D$

then with $\gamma = \frac{2D}{M\sqrt{N}}$, we have

$$\mathbf{E} \left[\text{er}_{\Phi} \left[\bar{\mathbf{f}}^{(N)} \right] \right] - \text{er}^{\Phi, \mathcal{F}, *} \leq \frac{2DM}{\sqrt{N}}$$

$$\bar{\mathbf{f}}^{(N)} = \frac{1}{N} \sum_{i=1}^N \mathbf{f}^{(i)}; \quad \mathbf{f}^{(n)} = FW^{(n)}\psi(\mathbf{x})$$

- Under the same assumptions, for any $\epsilon > 0$, we have $\mathbf{E} \left[\text{er}_L \left[\bar{\mathbf{f}}^{(N)} \right] \right] - \text{er}^{L, \mathcal{F}, *}< \epsilon$ if

$$N > \frac{4D^2M^2}{\check{\mathcal{H}}_{L, \Phi, \mathcal{F}}^2(\epsilon)}$$

- 01 Loss $DM = O(k)$
- Hamming Loss $DM = O(\log_2 k^3)$
- Block 01 Loss $DM = O(b)$

Thank You