

Alex Nowak-Vila, Francis Bach and Alessandro Rudi
INRIA - Ecole Normale Supérieure - PSL Research University

Challenges in Structured Prediction

In structured prediction:

- the **number of possible labels is exponentially large** w.r.t to the natural dimension of the data.
- we generally have **more possible outputs than training data**.

Examples of typical structured prediction problems:

- **Multilabel prediction**: predict a subset of labels.
- **Sequence prediction**: predict a sequence over a fixed dictionary.
- **Ranking**: predict a permutation.

Q: When is learning statistically and computationally feasible in structured prediction?

Supervised Learning Setting

- **Spaces**: input space \mathcal{X} , **discrete** label space \mathcal{Y} and **discrete** output space \mathcal{Z} .
- **Data**: n i.i.d observations $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ from a distribution P .
- **Structured loss**: a loss between outputs and labels

$$L : \mathcal{Z} \times \mathcal{Y} \rightarrow \mathbb{R}.$$

- **Expected risk and Bayes risk**: minimize the expected risk $\mathcal{E}(f)$:

$$\mathcal{E}(f) = \mathbb{E}_{(X,Y) \sim P} L(f(X), Y) = \int_{\mathcal{X}} \ell(f(x), x) dP(x),$$

where $\ell(z, x) = \int_{\mathcal{Y}} L(z, y) dP(y|x)$ is the Bayes risk.

- **Bayes classifier**: the Bayes classifier $f^* : \mathcal{X} \rightarrow \mathcal{Z}$ has the form:

$$f^*(x) = \arg \min_{z \in \mathcal{Z}} \ell(z, x), \quad f^* = \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Z}} \mathcal{E}(f).$$

Quadratic Surrogate (QS) Estimator

Given a kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ defined on the input space and $\lambda > 0$, the QS estimator introduced in [1] has the form

$$\hat{f}_n(x) = \arg \min_{z \in \mathcal{Z}} \sum_{i=1}^n \alpha_i(x) L(z, y_i), \quad (1)$$

where $\alpha(x) = (K + n\lambda I)^{-1} K_x \in \mathbb{R}^n$ with $K_x = (k(x, x_1), \dots, k(x, x_n)) \in \mathbb{R}^n$ and $K \in \mathbb{R}^{n \times n}$ is defined by $K_{ij} = k(x_i, x_j)$.

General Analysis of the Estimator

It is known that the resulting estimator is consistent. Moreover, we have the following generalization bound ([1]):

$$\mathcal{E}(\hat{f}_n) - \mathcal{E}(f^*) \leq C n^{-1/4},$$

where C is a constant. A priori, there is no control over the magnitude of the constant C ([2]).

- If $C \sim |\mathcal{Y}|, |\mathcal{Z}|$, then the bound is generally non-informative.

Goal: Characterize C for discrete losses.

Affine Decomposition of the Loss

We consider an **affine decomposition** of the loss matrix $L \in \mathbb{R}^{|\mathcal{Z}| \times |\mathcal{Y}|}$:

$$L = FU^T + c\mathbf{1}. \quad (2)$$

where $F \in \mathbb{R}^{|\mathcal{Z}| \times r}$, $U \in \mathbb{R}^{|\mathcal{Y}| \times r}$, $c \in \mathbb{R}$, $\mathbf{1} \in \mathbb{R}^{|\mathcal{Z}| \times |\mathcal{Y}|}$ is the matrix of ones and $r \in \mathbb{N}$.

- We will use it to characterize the statistical and computational complexity of learning with loss L .

Statistical Complexity Analysis

Let $n \in \mathbb{N}, \tau > 0$ and $\lambda_n = n^{-1/2}$. Assume that the loss L decomposes as Eq. (2). If $g^* \in \mathcal{G}$, we have that with probability $1 - 8e^{-\tau}$,

$$\mathcal{E}(\hat{f}_n) - \mathcal{E}(f^*) \leq A Q c \tau^2 n^{-1/4},$$

where $A = \sqrt{r} \|F\|_{\infty} U_{\max}$, $Q = \max_j \|g_j^*/U_{\max}\|$ and $U_{\max} = \max_{j,k} |U_{kj}|$.

The statistical complexity is characterized by the constant

$$A = \sqrt{r} \|F\|_{\infty} U_{\max}.$$

- Provide **generalization of low-noise conditions for general losses**: improved rates with conditions of the form $P_{\mathcal{X}}(\gamma(X) \leq \varepsilon) = o(\varepsilon^p)$ for $p > 0$, where $\gamma(x) = \min_{z' \neq f^*(x)} \ell(z', x) - \ell(f^*(x), x)$.
- **Tightness**: Use calibration dimension introduced by [3] to analyze tightness of statistical bounds.

Computational Complexity Analysis

Computational complexity of solving Eq. (1) is the same as the one of solving the following minimization problem:

$$\min_{z \in \mathcal{Z}} F_z \theta, \quad (3)$$

where $F_z \in \mathbb{R}^r$ is the z -th row of F and $\theta \in \mathbb{R}^r$.

The computational complexity is characterized by the cost of solving problem (3).

Analysis of Multilabel and Ranking Losses

- The label space of the following losses / scores is $\mathcal{Y} = \{0, 1\}^m$.

Multilabel and Ranking measures					
Measure	\mathcal{Z}	Definition	r	A	$\text{INF}_F(m)$
0-1 (\downarrow)	\mathcal{P}_m	$1(z \neq y)$	2^m	$2^{m/2}$	$\mathcal{O}(n \wedge 2^m)$
Block 0-1 (\downarrow)	\mathcal{P}_m	$1(z \in B_j, y \notin B_j, j \in [b])$	b	\sqrt{b}	$\mathcal{O}(b)$
Hamming (\downarrow)	\mathcal{P}_m	$\frac{1}{m} \sum_{j=1}^m 1([z]_j \neq [y]_j)$	m	$\frac{1}{2}$	$\mathcal{O}(m)$
F-score (\uparrow)	\mathcal{P}_m	$2 \frac{ z \cap y }{ z + y }$	$m^2 + 1$	$\sqrt{2}m$	$\mathcal{O}(m^2)$
Prec@k (\uparrow)	$\mathcal{P}_{m,k}$	$\frac{ z \cap y }{k}$	m	$\sqrt{\frac{m}{k}}$	$\mathcal{O}(m \log k)$
NDCG (\uparrow)	\mathfrak{S}_m	$\frac{1}{N(r)} \sum_{j=1}^m G([r]_j) D_{\sigma(j)}$	m	$\sqrt{m} (\sum_j D_j^2)^{1/2} G_{\max}$	$\mathcal{O}(m \log m)$
PD (\downarrow)	\mathfrak{S}_m	$\frac{1}{N(y)} \sum_{j,\ell=1}^m 1([y]_j < [y]_{\ell}) 1(\sigma(j) > \sigma(\ell))$	$\frac{m(m-1)}{2}$	$\frac{m}{4}$	$\text{MWFAS}(m)$
MAP (\uparrow)	\mathfrak{S}_m	$\frac{1}{ y } \sum_{j=1}^m \frac{ y_j }{\sigma(j)} \sum_{\ell=1}^{\sigma(j)} y_{\sigma^{-1}(\ell)}$	$\frac{m(m+1)}{2}$	$\frac{1}{2} m \sqrt{\log(m+1)}$	$\text{QAP}(m)$

Experiments

- Perform experiments on multilabel and ranking datasets: compare with SSVM and threshold-based method.

Multilabel	bibtext									Ranking		Observed
	n	d	m	THBM	SSVM	QS	0-1	Ham	F-score	NDCG@3	SSVM	
	7395	1836	159	0.82	0.91	0.78	0.82	0.57	0.25	0.47	0.47	106
	645	260	19	1.0	0.53	0.52	1.0	0.16	0.28	0.28	0.51	25
	502	68	174	0.99	1.0	1.0	0.99	0.33	0.47	0.47	0.95	150
	5000	499	374	0.92	0.99	0.95	0.92	0.11	0.26	0.26	0.86	0.47
	1702	1001	53	0.93	0.90	0.86	0.90	0.49	0.52	0.52	0.86	0.47
	43907	120	101	0.31	1.0	0.29	0.31	0.80	0.83	0.83	0.34	0.47
	978	1449	45	0.49	0.35	0.51	0.49	0.63	0.68	0.68	0.47	0.47
	2407	294	6	0.93	0.95	0.76	0.93	0.48	0.48	0.48	0.47	0.47
0-1 (\downarrow)				0.82	0.91	0.78	0.82	0.57	0.25	0.47	0.47	0.47
				1.3e-2	6.4e-2	1.3e-2	1.3e-2	7.9e-2	4.9e-2	0.14	0.14	0.45
				1.1e-2	1.0e-2	8.6e-2	1.1e-2	5.9e-2	9.4e-3	0.14	0.14	0.48
Ham (\downarrow)				0.44	0.19	0.47	0.44	0.25	0.28	0.47	0.47	0.47
				0.25	0.16	0.28	0.25	0.46	0.47	0.47	0.47	0.47
				0.25	0.16	0.28	0.25	0.46	0.47	0.47	0.47	0.47
				0.51	0.49	0.56	0.51	0.49	0.56	0.56	0.48	0.47
				0.80	0.74	0.83	0.80	0.63	0.68	0.68	0.48	0.47
F-score (\uparrow)				0.44	0.19	0.47	0.44	0.25	0.28	0.47	0.47	0.47
				0.25	0.16	0.28	0.25	0.46	0.47	0.47	0.47	0.47
				0.25	0.16	0.28	0.25	0.46	0.47	0.47	0.47	0.47
				0.51	0.49	0.56	0.51	0.49	0.56	0.56	0.48	0.47
				0.80	0.74	0.83	0.80	0.63	0.68	0.68	0.48	0.47
				0.44	0.19	0.47	0.44	0.25	0.28	0.47	0.47	0.47
				0.25	0.16	0.28	0.25	0.46	0.47	0.47	0.47	0.47
				0.25	0.16	0.28	0.25	0.46	0.47	0.47	0.47	0.47
				0.51	0.49	0.56	0.51	0.49	0.56	0.56	0.48	0.47
				0.80	0.74	0.83	0.80	0.63	0.68	0.68	0.48	0.47

- **Take-away message**: importance of being consistent and calibrated to the measure of interest.

Main References

- [1] Carlo Ciliberto, Lorenzo Rosasco and Alessandro Rudi. A consistent regularization approach for structured prediction. In *Advances in Neural Information Processing Systems*, 2016.
- [2] Anton Osokin, Francis Bach and Simon Lacoste-Julien. On structured prediction theory with calibrated convex surrogate losses. In *Advances in Neural Information Processing Systems*, 2017.
- [3] Harish G. Ramaswamy and Shivani Agarwal. Convex calibration dimension for multiclass loss matrices. *The Journal of Machine Learning Research*, 17(1):397–441, 2016.