

Similarity Learning for Pointwise ROC Optimization

Robin Vogel^{1,3}, **Stéphan Clémençon**¹ **and Aurélien Bellet**² ¹ **Télécom ParisTech**, ² **Inria**, ³ **IDEMIA**

MOTIVATION

Biometric identification = checks correspondance of two measurements (x, x').

Given a similarity S and a threshold t,

(x, x') is a match $\Leftrightarrow S(x, x') > t.$ (1)

The ROC curve of S gives the true positive rate (TPR) given the false positive rate (FPR) for eq. (1).

Biometric systems are deployed to function at fixed FPR, see [1], hence we study *pointwise* ROC *opti-*

CONTRIBUTIONS

Proposition of an appropriate probabilistic framework for a novel perspective on similarity learning.

Statistical guarantees for the constrained optimization problem corresponding to the empirical version of our theoretical objective, i.e. pointwise ROC optimization in the context of pairwise ranking.

Faster rates under a weak low-noise assumption.

Empirical illustration of the faster rates through numerical simulations.

SCALABILITY

When n large and K large, calculating $R_n^-(G)$ is computationally costly. A sensible approach is to drastically subsample the negative pairs, while keeping all positive pairs.

We studied the equivalent of eq. (3) when replacing $\hat{R}_n^-(G)$ by the following approximation:

$$\bar{R}_B^-(G) := \frac{1}{B} \sum_{(i,j)\in\mathcal{P}_B} G(X_i, X_j),$$

where \mathcal{P}_B is a set of cardinality *B* built by sampling with replacement in the set of negative training pairs Λ_P , with:

mization.



PRELIMINARIES

Classification setting. Assume $(X, Y) \sim P$, with:

- $Y \in \{1, \ldots, K\}$ the output label,
- $X \in \mathcal{X} \subset \mathbb{R}^d$ input random variable.

Similarity learning. Select a similarity S s.t.

the larger S(X, X') the higher $\mathbb{P}\{Y = Y' \mid X, X'\}$, with $(X, Y) \perp (X', Y') \sim P$.

Optimal similarity rules S^* are increasing transforms of the posterior probability η :

Study of sampling strategies for scalability issues that arise from the high number of negative pairs.

GENERALIZATION

Generalization guarantees. Our theorem describes the generalization capacities of the solution of **eq. (3)**, under some conditions on \mathcal{G}_0 and a suitable choice of Φ .

Theorem 1. Suppose that:

- G_0 is a VC-major class of VC-dimension V,
- $\forall G \in \mathcal{G}_0, \|G\|_{\infty} \leq 1$,
- $\exists \kappa \in (0,1)$ such that $\kappa \leq \mathbb{P}\{Y = Y'\} \leq 1 \kappa$,

For all $\delta \in (0, 1)$ and n > 1, :

• set $\Phi_{n,\delta} = C_{V,\delta,\kappa} \cdot n^{-1/2}$, where $C_{V,\delta,\kappa}$ is known and depends on V, δ, κ , $\Lambda_P = \{(i,j) \mid i,j \in \{1,\ldots,n\}; Y_i \neq Y_j\}.$

We show that the results of **theorem 1** still hold, with a different Φ of the order $O(\sqrt{\log(n)/B}))$, using results from [2].

It implies that it is sufficient to sample B = O(n)pairs to get a learning rate of order $O(\sqrt{\log(n)/n})$.

EXPERIMENT ON FAST RATES

We illustrate the results presented in **theorem 2**.

How?

Introduce distributions that satisfy (NA) with different *a*'s, and show the difference in the generalization speed.

Defining $K = 2, X \sim \mathcal{U}[0, 1], \mathbb{P}\{Y = 1\} = 1/2$ and the density μ_1 of $X \mid Y = 1$ fully caracterizes P. Hence we define a family of μ_1 's parameterized by a.

 $\eta(x, x') = \mathbb{P}\{Y = Y' \mid (X, X') = (x, x')\}.$

Pointwise ROC optimization. Given $\alpha \in (0, 1)$,

$$\max_{G \in \mathcal{G}_0} R^+(G) \quad \text{subject to} \quad R^-(G) \le \alpha, \quad (2)$$
where $R^+(G) = \mathbb{E}[G(X, X') \mid Y = Y']$,
 $R^-(G) = \mathbb{E}[G(X, X') \mid Y \ne Y']$ and \mathcal{G}_0 a class of functions.

By the Neyman-Pearson lemma:

 $G_{\alpha}^* := \mathbb{I}_{\mathcal{R}_{\alpha}^*}$ optimal solution of **eq. (2)**,

where $\mathcal{R}^*_{\alpha} = \{(x, x') \in \mathcal{X}^2 : \eta(x, x') \ge Q^*_{\alpha}\}$, with Q^*_{α} quantile of $\eta(X, X') \mid Y \neq Y'$ at level $1 - \alpha$.

Empirical problem. Using a training sample:

 $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\},\$

composed of *n* i.i.d. copies of (X, Y), we form esti-

• Let \hat{G}_n solution of eq. (3) with $\Phi = \Phi_{n,\delta/2}$,

We have w.p. $\geq 1 - \delta$, $\forall n \geq 1 + 4\kappa^{-2} \log(3/\delta)$,

 $R^{+}(\hat{G}_{n}) \geq \sup_{G \in \mathcal{G}_{0}: R^{-}(G) \leq \alpha} R^{+}(G) - \Phi_{n,\delta/2},$ and $R^{-}(\hat{G}_{n}) \leq \alpha + \Phi_{n,\delta/2}.$

FAST RATES

In some situations, rates faster than $O(1/\sqrt{n})$ can be achieved by solutions of eq. (3). These rates hold when the following noise assumption is verified:

Noise assumption (NA). $\exists c \in \mathbb{R}^*_+, a \in [0, 1] \text{ s.t.},$

 $\mathbb{E}_{X'}\left[|\eta(X,X') - Q^*_{\alpha}|^{-a}\right] \le c \quad a.s.$

Our next theorem establishes fast rate bounds under the (NA) condition on the data distribution. It For some a, the 90-quantile of $\log(R^+(G_{\alpha}^*) - R^+(\hat{G}_n))$ for different n's is fitted to $C_a \times \log(n) + D_a$ to get the empirical generalization speed C_a .





mates of $R^+(G)$ and $R^-(G)$:

 $\hat{R}_{n}^{+}(G) = \frac{1}{n_{+}} \sum_{1 \le i < j \le n} G(X_{i}, X_{j}) \cdot \mathbb{I}\{Y_{i} = Y_{j}\},\$ $\hat{R}_{n}^{-}(G) = \frac{1}{n_{-}} \sum_{1 \le i < j \le n} G(X_{i}, X_{j}) \cdot \mathbb{I}\{Y_{i} \ne Y_{j}\},\$

with $n_+ = \sum_{1 \le i < j \le n} \mathbb{I}\{Y_i = Y_j\} = n(n-1)/2 - n_-.$

One can then derive the empirical version of **eq. (2)** :

 $\max_{G \in \mathcal{G}_0} R_n^+(G) \quad \text{subject to} \quad R_n^-(G) \le \alpha + \Phi, \quad (3)$

where $\Phi > 0$ is some tolerance parameter.

relies on a variant of the Bernstein inequality for U-statistics.

Theorem 2. *Suppose that:*

- the assumptions of *theorem 1* are satisfied,
- NA holds true,

• $G^*_{\alpha} \in \mathcal{G}_0$,

Fix $\delta > 0$, there exists C', that depends of δ , κ , Q_{α}^* , a, c and V such that, w.p. $\geq 1 - \delta$,

 $R^+(\hat{G}_n) \ge R^+(G^*_{\alpha}) - C' \cdot n^{-(2+a)/4},$ and $R^-(\hat{G}_n) \le \alpha + \Phi_{n,\delta/2}.$

REFERENCES

[1] Anil K. Jain, Arun A. Ross, and Karthik Nandakumar. *Introduction to Biometrics*. Springer, 2011.

[2] Stephan Clémençon, Igor Colin, and Aurélien Bellet.
 Scaling-up Empirical Risk Minimization: Optimization of Incomplete U-statistics. Journal of Machine Learning Research, 17(76):1–36, 2016.

[3] Stephan Clémençon and Nicolas Vayatis. Overlaying classifiers: A practical approach to optimal scoring. *Constructive Approximation*, 32(3):619–648, 2010.

[4] Stephan Clémençon, Gàbor Lugosi, and Nicolas Vayatis. Ranking and Empirical Minimization of U-Statistics. *The Annals of Statistics*, 36(2):844–874, 2008.