



A!

Aalto University

# A GENERAL ONLINE ALGORITHM FOR OPTIMIZING COMPLEX PERFORMANCE METRICS

Wojciech Kotłowski<sup>1</sup> Marek Wydmuch<sup>1</sup> Erik Schultheis<sup>2</sup> Rohit Babbar<sup>2,3</sup> Krzysztof Dembczyński<sup>1,4</sup>

<sup>1</sup>Poznan University of Technology, Poland

<sup>2</sup>Aalto University, Helsinki, Finland

<sup>3</sup>University of Bath, UK

<sup>4</sup>Yahoo Research, New York, USA



## 1. Online classification

- Classification:  $x \in \mathcal{X} \rightarrow y \in \mathcal{Y}, (x, y) \sim \mathbb{P}$ .
- $\mathbb{P}_{\mathcal{X}}$  – marginal distribution of  $x$
- $\eta(x)$  – label conditional probability,  $\eta(x) = \mathbb{E}_{y|x}[y]$ .
- Multi-class classification:**
  - $y \in \{y' \in \{0, 1\}^m : \|y'\|_1 = 1\}$  – a vector one-hot encoding of one out of  $m$  classes,
  - $\eta(x) = (\eta_1(x), \dots, \eta_m(x))$  – label conditional distribution with  $\eta_j(x) = \mathbb{P}(y_j = 1|x)$
- Multi-label classification:**
  - $y := \{0, 1\}^m$  – a vector of relevant labels
  - $\eta(x) = (\eta_1(x), \dots, \eta_m(x))$  – a vector of **label marginal conditional probabilities**,  $\eta_j(x) = \mathbb{P}(y_j = 1|x)$ .
- Online protocol:**
  - Input:** a sequence  $(x^n, y^n) = ((x_1, y_1), \dots, (x_n, y_n))$
  - for  $t = 1, \dots, n$  do
    - Observe input instance  $x_t$  drawn from  $\mathbb{P}_{\mathcal{X}}$
    - Receive conditional probability estimate  $\hat{\eta}_t(x_t)$
    - Predict label  $\hat{y}_t$  based on  $\hat{\eta}_t(x_t)$
    - Receive true label  $y_t$  drawn from  $\mathbb{P}(\cdot|x_t)$
    - Evaluate based on  $\psi(y^n, \hat{y}^n)$
- Goal:** Maximize performance metric  $\psi(y^n, \hat{y}^n)$

## 2. Confusion matrices

Binary $C^{2 \times 2}$ :		Multi-class $C^{m \times m}$ :	
$\hat{y}$		$\hat{y}$	
0	1	1	2 ... $m$
$y$	0 TN FP	TP <sub>1</sub> E <sub>1,2</sub> ... E <sub>1,m</sub>	
1	FN TP	E <sub>2,1</sub> TP <sub>2</sub> ... E <sub>2,m</sub>	
Multi-label $C^{m \times 2 \times 2}$ :		$\hat{y}_1$	$\hat{y}_2$ ... $\hat{y}_m$
$y_1$	0 TN <sub>1</sub> FP <sub>1</sub>	0 1	0 1
1	FN <sub>1</sub> TP <sub>1</sub>	FN <sub>2</sub> FP <sub>2</sub>	FN <sub>m</sub> FP <sub>m</sub>

## 3. Complex performance metrics

We focus on the online maximization of performance metrics that do **not decompose into a sum over instances** but are general **functions of the confusion matrix**.

Examples of binary and multi-class confusion matrix measures.

Metric	$\psi(C^{2 \times 2})$	$\psi(C^{m \times m})$
Balanced Acc.	$\frac{tp}{2(tp+fn)} + \frac{tn}{2(tn+fp)}$	$\sum_{i=1}^m \frac{C_{ii}}{m \sum_{j=1}^m C_{ij}}$
Recall	$\frac{tp}{tp+fn}$	micro- or macro-average
Precision	$\frac{tp}{tp+fp}$	micro- or macro-average
$F_{\beta}$ -measure	$\frac{(1+\beta^2)tp}{(1+\beta^2)tp+\beta^2fn+fp}$	micro- or macro-average
G-mean	$\sqrt{\frac{tp \cdot tn}{(tp+fn)(tn+fp)}}$	$\left(\prod_{j=1}^m \frac{C_{jj}}{\sum_{i=1}^m C_{ji}}\right)^{1/m}$
H-mean	$2 \left( \frac{tp+fn}{tp} + \frac{tn+fp}{tn} \right)^{-1}$	$m \left( \sum_{j=1}^m \frac{C_{jj}}{\sum_{i=1}^m C_{ji}} \right)^{-1}$
Q-mean	$1 - \sqrt{\frac{1}{2} \left( \left( \frac{fn}{tp+fn} \right)^2 + \left( \frac{fp}{tn+fp} \right)^2 \right)}$	$1 - \sqrt{\frac{1}{m} \sum_{j=1}^m \left( 1 - \frac{C_{jj}}{\sum_{i=1}^m C_{ji}} \right)^2}$

Micro- and macro-averages of metrics for  $C^{m \times 2 \times 2}$ :

$$\text{Micro-}\psi(\mathbf{C}) := \psi\left(\frac{1}{m} \sum_{j=1}^m \mathbf{C}_j\right), \quad \text{Macro-}\psi(\mathbf{C}) := \frac{1}{m} \sum_{j=1}^m \psi(\mathbf{C}_j).$$

## 5. Theoretical results

If  $\psi(\mathbf{C})$  is **concave, Lipschitz, and smooth**:

- With  $\eta(x)$  available, the regret of OMMA is of order  $\mathcal{O}(\frac{\ln n}{n})$ .
- Otherwise, the regret is bounded by  $\mathbb{E}[\|\eta(x) - \hat{\eta}(x)\|]$ .
- If the estimation error converges to zero with  $n \rightarrow \infty$ , OMMA becomes a **no-regret** learning algorithm.

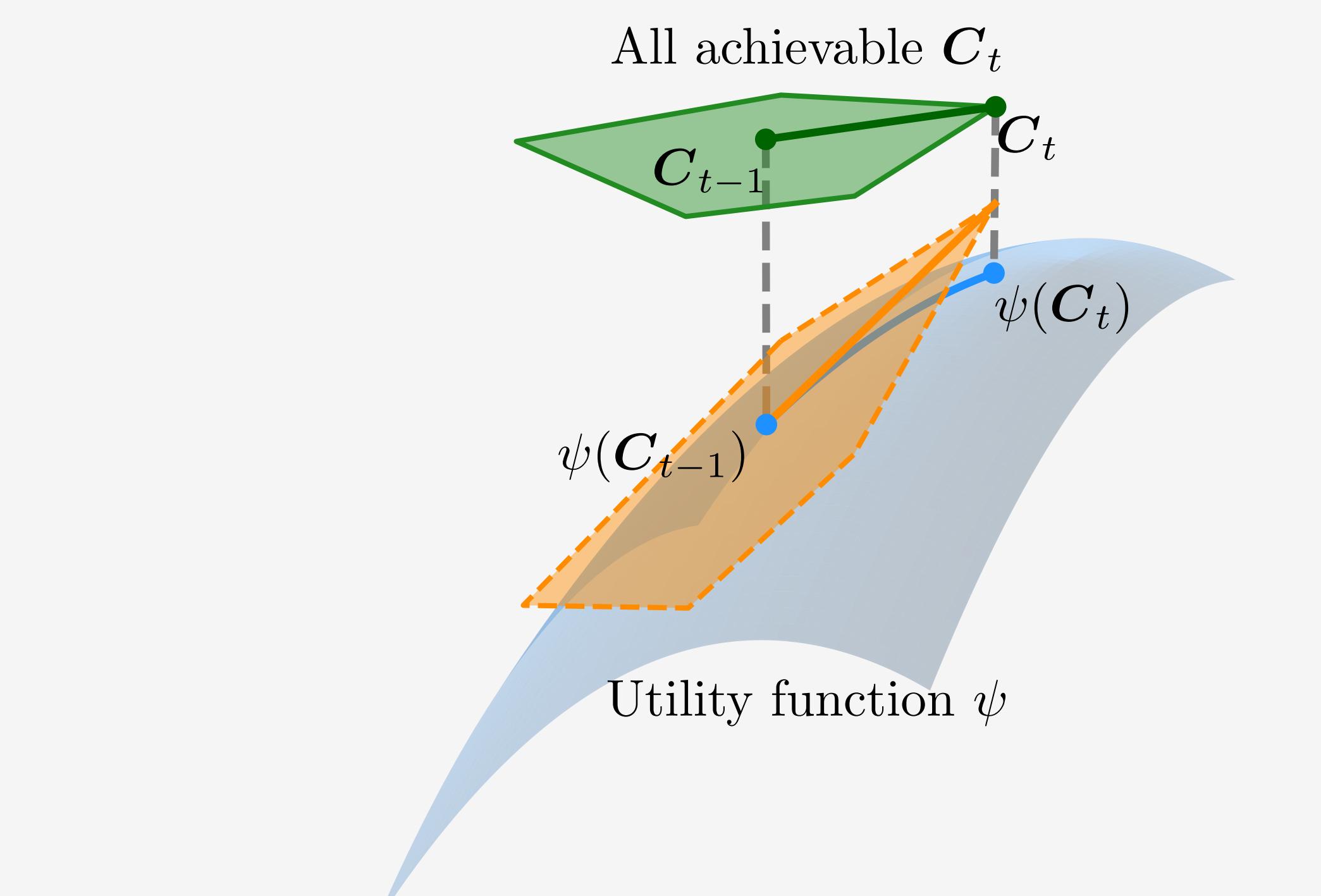
For some  $\psi$ , Lipschitzness and smoothness might **not hold globally**:

- These properties are invoked only along the parameter path of the algorithm,  $\{\mathbf{C}_t\}_{t=1}^n$ .
- Many  $\psi$  are Lipschitz and smooth when label frequencies are bounded away from zero
- $\Rightarrow$  adding a small value  $\lambda$  to  $\mathbf{C}$  stabilizes the algorithm.
- $\Rightarrow$  OMMA is a **robust** approach for online learning scenarios.

## 4. Online Metric Maximi. Algo. (OMMA)

Under assumption that the utility  $\psi(\mathbf{C})$  is **differentiable** in  $\mathbf{C}$ :

- Initialization:** confusion matrix  $\mathbf{C}_0$
- for**  $t = 1, \dots, n$  **do**
- 3: Receive input  $x_t$  and probability estimate  $\hat{\eta}_t(x_t)$
- 4: Predict  $\hat{y}_t = \text{argmax}_{\hat{y}} \nabla \psi(\mathbf{C}_{t-1}) \cdot (\mathbb{E}_{y_t \sim \hat{\eta}_t(x_t)} [\mathbf{C}(y_t, \hat{y}_t)])$
- 5: Receive label  $y_t$  and update  $\mathbf{C}_t = \frac{t-1}{t} \mathbf{C}_{t-1} + \frac{1}{t} \mathbf{C}(y_t, \hat{y}_t)$



## 6. Example: binary classification

$-y_t \in \{0, 1\}$ , and  $\eta(x_t) = P(y_t = 1|x_t)$ .

-Because:

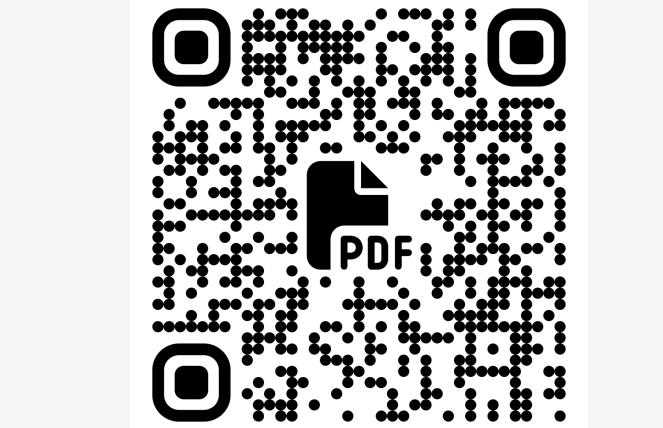
$$\mathbb{E}_{y_t \sim \hat{\eta}_t(x_t)} [\mathbf{C}(y_t, \hat{y}_t)] = \begin{bmatrix} (1 - \hat{\eta}_t(x_t))(1 - \hat{y}_t) & (1 - \hat{\eta}_t(x_t))\hat{y}_t \\ \hat{\eta}_t(x_t)(1 - \hat{y}_t) & \hat{\eta}_t(x_t)\hat{y}_t \end{bmatrix},$$

line 3 of the algorithm boils down to maximizing a cost-sensitive classification accuracy  $\hat{y}(\alpha\hat{\eta}(x_t) - \beta)$ , with:

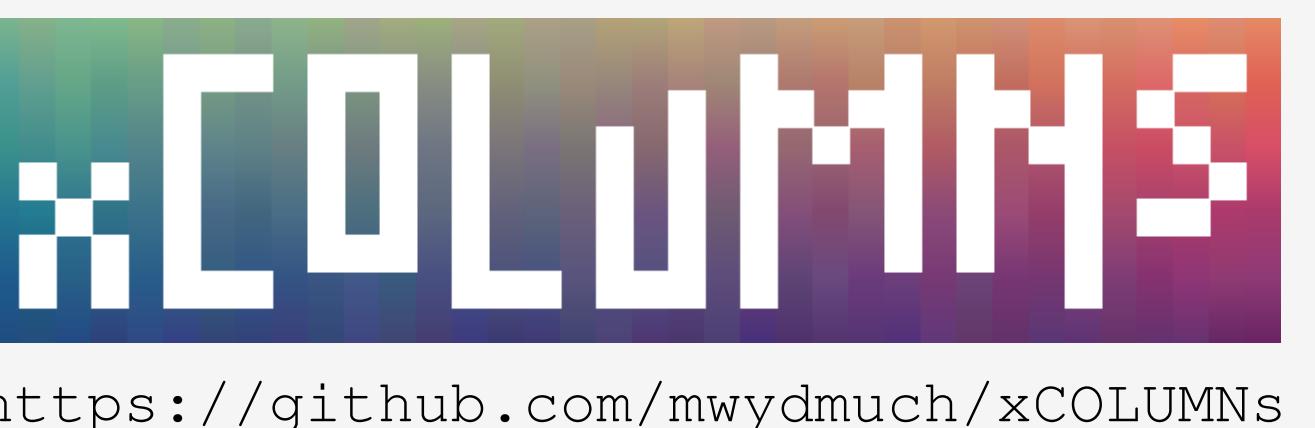
$$\alpha = \nabla_{11} + \nabla_{00} - \nabla_{01} - \nabla_{10}, \quad \beta = \nabla_{00} - \nabla_{01}.$$

- All practical utilities are non-decreasing with TP and TN ( $\nabla_{11}, \nabla_{00} \geq 0$ ), and non-increasing with FN and FP ( $\nabla_{01}, \nabla_{10} \leq 0$ ).
- $a \geq 0 \rightarrow$  maximizing  $\hat{y}(\alpha\hat{\eta}(x_t) - \beta)$  boils down to choosing  $\hat{y}_t = 1$  when  $\hat{\eta}(x_t) > \beta/\alpha$  – well-known rule for maximizing the population-level version of  $\psi$  (OMMA **mimics** this optimal rule).

Read the paper



Try our library



<https://github.com/mwydmuch/xCOLUMNS>

## 7. Experimental results

Results of different online algorithms on multi-class and multi-label problems compared to OMMA (all use the same  $\hat{\eta}(x)$ ):

Method	LEDGAR-LEXGLUE ( $m = 100, n = 10000$ )				FLICKR ( $m = 195, n = 24154$ )			
	Macro F1	F1 @ 3 Rec. @ 3	Pr. @ 3	G-mean	Macro F1	F1	F1 @ 3 Rec. @ 3	G-mean
Top-k / $\hat{\eta} > 0.5$	79.06	51.89	92.04	38.88	0.00	0.00	69.21	29.46
OFO	x	x	x	x	x	x	x	x
Greedy	79.30	78.08	93.26	91.17	x	x	30.90	30.42
Online-FW*	79.22	70.94	93.30	54.52	62.31	75.81	74.59	46.41
Online-FW( $\hat{\eta}$ )*	79.22	73.88	93.38	49.63	78.02	77.58	47.02	83.37
OMMA (ours)	79.28	78.10	93.26	91.41	77.48	74.62	41.01	30.90
OMMA( $\hat{\eta}$ ) (ours)	79.34	78.22	93.39	90.10	78.03	76.08	74.53	46.33

Method	RCV1X ( $m = 2456, n = 155062$ )				AMAZONCAT ( $m = 306784$ )			
	Macro F1	F1	F1 @ 3 Rec. @ 3	Pr. @ 3	Macro F1	F1	F1 @ 3 Rec. @ 3	G-mean H-mean
Top-k / $\hat{\eta} > 0.5$	68.57	11.29	5.34	4.59	13.24	16.01	12.36	67.77
OFO	69.83	20.26	x	x	x	x	x	x
Greedy	x	20.80	16.01	21.20	30.91	69.07	67.04	44.20
Online-FW*	69.83	19.82	15.33	21.09	19.88	69.07	67.04	42.42
Online-FW( $\hat{\eta}$ )*	69.79	20.40	15.55	22.21	22.17	69.06	67.04	47.32
OMMA (ours)	69.77	20.57	15.87	21.08	30.39	69.07	67.04	43.00
OMMA( $\hat{\eta}$ ) (ours)	69.71	20.71	16.07	22.23	30.38	69.06	67.04	47.67

