

# OPTIMIZING THE F-MEASURE IN MULTI-LABEL CLASSIFICATION: PLUG-IN RULE APPROACH VERSUS STRUCTURED LOSS MINIMIZATION

Krzysztof Dembczyński<sup>1</sup>, Arkadiusz Jachnik<sup>1</sup>, Wojciech Kotłowski<sup>1</sup>,  
Willem Waegeman<sup>2</sup>, and Eyke Hüllermeier<sup>3</sup>

<sup>1</sup> Intelligent Decision Support Systems Laboratory, Poznań University of Technology, Poland  
<sup>2</sup> NGDATA-Europe, Belgium <sup>3</sup> Mathematics and Computer Science, Marburg University, Germany



NGDATA

## Multi-Label Classification (MLC)

- For a feature vector  $\mathbf{x}$  predict a binary vector of responses  $\mathbf{y}$  using a prediction function  $\mathbf{h}(\mathbf{x})$ :

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \xrightarrow{\mathbf{h}(\mathbf{x})} \mathbf{y} = (y_1, y_2, \dots, y_m)$$

- Main challenges in multi-label classification:**

- Appropriate modeling of label dependencies between labels

$$y_1, y_2, \dots, y_m$$

- A multitude of multivariate loss functions defined over the binary vectors

$$\ell(\mathbf{y}, \mathbf{h}(\mathbf{x}))$$

## $F_\beta$ -measure

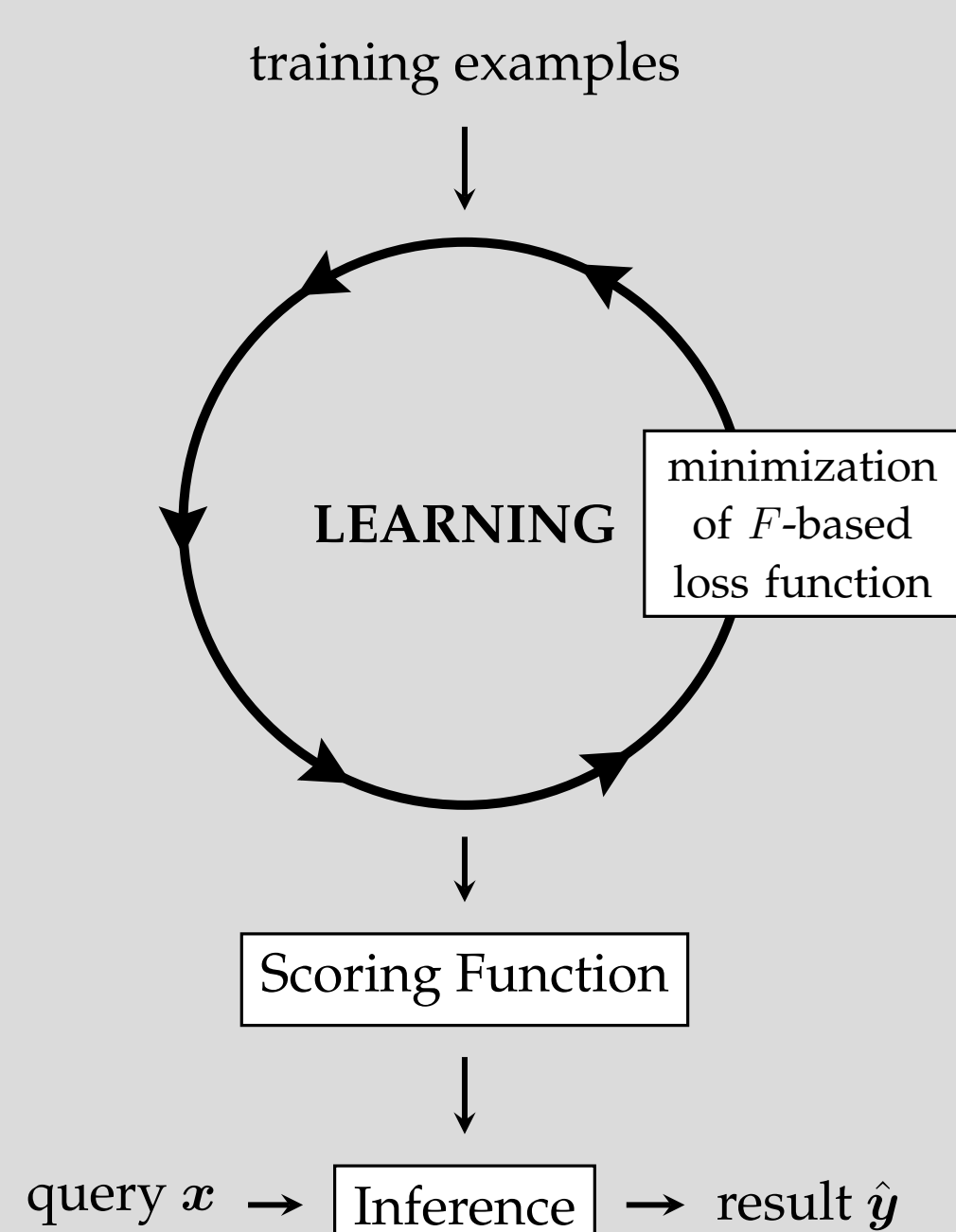
- We focus on the  $F_\beta$ -measure-based loss function ( $F_\beta$ -loss):

$$\begin{aligned} \ell_{F_\beta}(\mathbf{y}, \mathbf{h}(\mathbf{x})) &= 1 - F_\beta(\mathbf{y}, \mathbf{h}(\mathbf{x})) \\ &= 1 - \frac{(1 + \beta^2) \sum_{i=1}^m y_i h_i(\mathbf{x})}{\beta^2 \sum_{i=1}^m y_i + \sum_{i=1}^m h_i(\mathbf{x})} \in [0, 1]. \end{aligned}$$

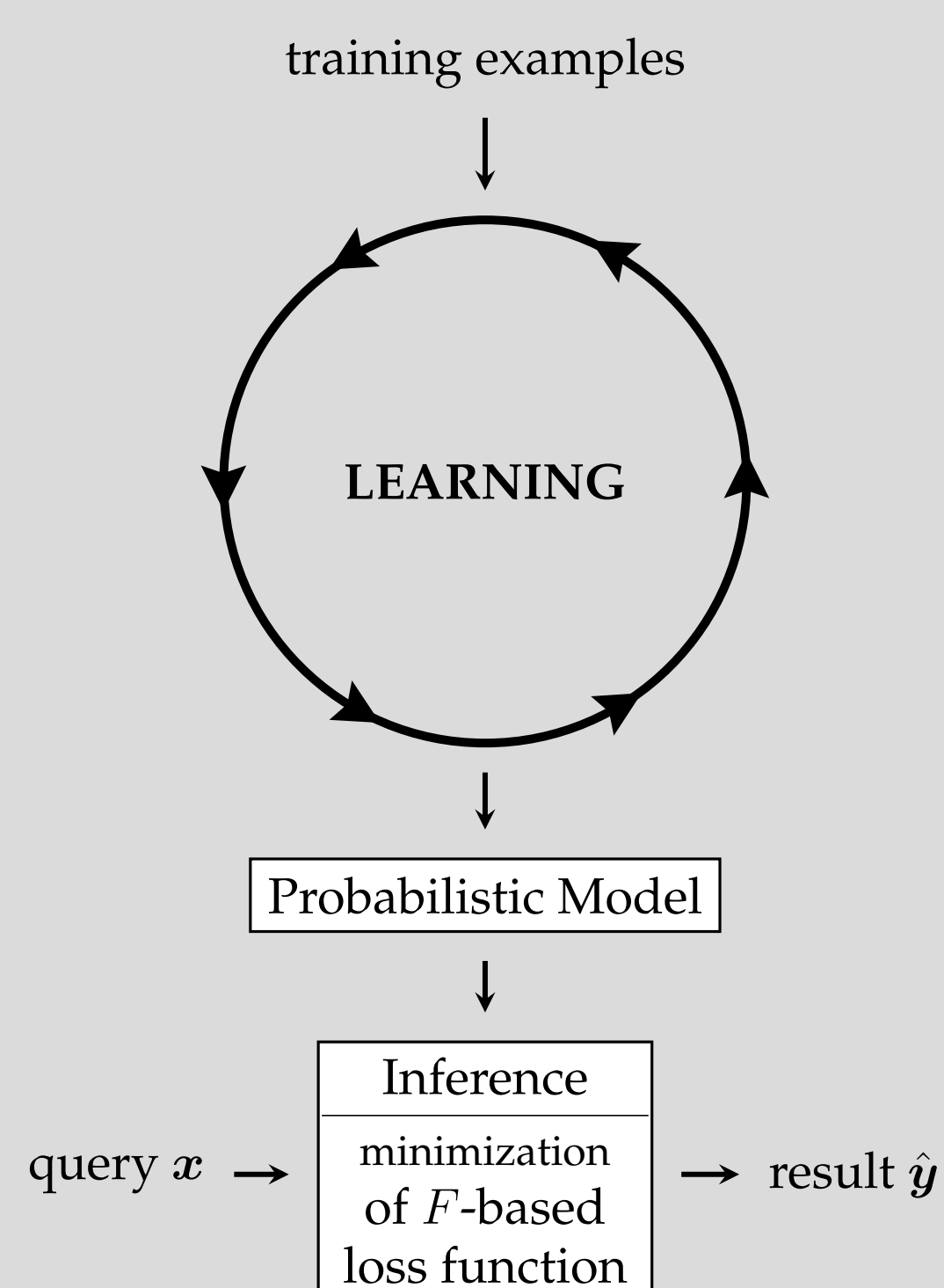
- Provides a **better balance** between relevant and irrelevant labels.
- However, it is **not easy** to optimize.

## Two Approaches

### Structured loss minimization



### Plug-in rule approach



## Structured Loss Minimization with SSVM

- Use a scoring function  $f(\mathbf{y}, \mathbf{x})$
- Minimize the **structured hinge loss** [5]:

$$\tilde{\ell}_h(\mathbf{y}, \mathbf{x}, f) = \max_{\mathbf{y}' \in \mathcal{Y}} \{ \ell_F(\mathbf{y}, \mathbf{y}') + f(\mathbf{y}', \mathbf{x}) \} - f(\mathbf{y}, \mathbf{x}),$$

- With  $\ell_F(\mathbf{y}, \mathbf{y}')$  used for margin rescaling.
- Predict according to:

$$\mathbf{h}(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{y}, \mathbf{x}).$$

- Requires solving the **arg max** and **constraint generation** problem.
- Two algorithms:

### RML [3]

No label interactions:

$$f(\mathbf{y}, \mathbf{x}) = \sum_{i=1}^m f_i(y_i, \mathbf{x})$$

Quadratic learning and linear prediction

### SML [4]

Submodular interactions:

$$f(\mathbf{y}, \mathbf{x}) = \sum_{i=1}^m f_i(y_i, \mathbf{x}) + \sum_{y_k, y_l} f_{k,l}(y_k, y_l)$$

More complex (graph-cut and approximate algorithms)

## Plug-in Rule Approaches with LR

- Plug estimates of required parameters into the **Bayes classifier**.
- The brute-force algorithm is **intractable**:

$$\mathbf{h}^* = \arg \min_{\mathbf{h} \in \mathcal{Y}} \mathbb{E} [\ell_{F_\beta}(\mathbf{Y}, \mathbf{h})] = \arg \max_{\mathbf{h} \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} \Pr(\mathbf{y}) \frac{(\beta + 1) \sum_{i=1}^m y_i h_i}{\beta^2 \sum_{i=1}^m y_i + \sum_{i=1}^m h_i}$$

- Approximation needed?** Not really. The exact solution is **tractable**!

### LFP

- Assumes label independence
- Linear number of parameters:  $\Pr(y_i = 1)$
- Inference based on dynamic programming [6]
- Reduction to LR for each label

### EFP

- No assumptions
- Quadratic number of parameters:  $\Pr(y_i = 1, s = \sum_i y_i)$
- Inference based on matrix multiplication and top  $k$  selection [1]
- Reduction to multinomial LR for each label

## Theoretical Analysis

- Computational complexity** (with respect to the number of labels):

	RML	SML	LFP	EFP
learning	$\mathcal{O}(m^2)$	$\mathcal{O}(m^4)$	$\mathcal{O}(m)$	$\mathcal{O}(m^2)$
prediction	$\mathcal{O}(m)$	$\mathcal{O}(m^3)$	$\mathcal{O}(m^2)$	$\mathcal{O}(m^3)$

- Statistical consistency** of multi-label classifiers [2]:

- RML and SML are **not consistent**
- EFP is **consistent**.

## Empirical Evaluation On Benchmark Datasets

	HL[%]↓	$F_1$ [%]↑	$t_{cv}$	$t_{train}$	$t_{inf}$	HL[%]↓	$F_1$ [%]↑	$t_{cv}$	$t_{train}$	$t_{inf}$
IMAGE			SCENE							
BR	<b>19.90</b>	43.63	<b>9</b>	<b>0.392</b>	0.087	10.51	55.73	<b>29</b>	<b>0.733</b>	0.241
LFP	27.55	58.86	<b>9</b>	<b>0.392</b>	0.119	12.18	74.38	<b>29</b>	<b>0.733</b>	0.270
EFP	26.07	<b>59.77</b>	24	0.606	0.183	12.22	<b>74.44</b>	72	0.995	0.399
RML	25.07	57.49	94	1.104	<b>0.051</b>	<b>9.70</b>	73.92	73	1.001	<b>0.118</b>
SML	28.82	56.99	156	7.116	0.052	15.65	68.50	52	1.129	0.123
YEAST			MEDICAL							
BR	<b>20.03</b>	60.59	<b>12</b>	<b>0.429</b>	0.128	<b>1.17</b>	70.19	<b>9</b>	<b>1</b>	0.952
LFP	22.24	65.02	<b>12</b>	<b>0.429</b>	0.146	1.18	<b>81.27</b>	<b>9</b>	<b>1</b>	1.513
EFP	22.82	<b>65.47</b>	101	2.004	0.367	1.23	80.39	16	1	1.883
RML	22.82	64.78	206	5.194	<b>0.056</b>	1.20	80.63	1253	30	<b>0.144</b>
SML	24.52	63.96	319	4.385	0.070	2.50	67.90	715	23	0.773
ENRON			MEDIAMILL							
BR	<b>4.54</b>	55.49	<b>52</b>	<b>4</b>	1.016	<b>3.19</b>	51.21	<b>3238</b>	<b>118</b>	13
LFP	6.09	56.86	<b>52</b>	<b>4</b>	1.519	3.67	55.15	<b>3238</b>	<b>118</b>	20
EFP	5.34	<b>61.04</b>	214	6	2.628	3.63	<b>55.16</b>	24620	440	30
RML	6.35	57.69	3897	41	<b>0.143</b>	4.12	49.35	–	1125	<b>7</b>
SML	7.82	54.61	18780	62	0.887	4.18	50.02	–	10365	131

Experimental results for Hamming loss (HL),  $F_1$ , and running times (in seconds) of cross-validation ( $t_{cv}$ ), training ( $t_{train}$ ) (for the best set of parameters) and inference ( $t_{inf}$ ). The best results are marked by a '\*'. BR denotes Binary Relevance which is a baseline in the comparison.

## References

- K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An exact algorithm for F-measure maximization. In *NIPS*, volume 25, 2011.
- W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Journal of Machine Learning Research - Proceedings Track*, 19:341–358, 2011.
- J. Petterson and T. S. Caetano. Reverse multi-label learning. In *Advances in Neural Information Processing Systems 24*, pages 1912–1920, 2010.
- J. Petterson and T. S. Caetano. Submodular multi-label learning. In *Advances in Neural Information Processing Systems 24*, pages 1512–1520, 2011.
- Y. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun. Large margin methods for structured and interdependent output variables. *JMLR*, 6:1453–1484, 2005.
- N. Ye, K. Chai, W. Lee, and H. Chieu. Optimizing F-measures: a tale of two approaches. In *ICML*, 2012.

This project is partially supported by the Foundation of Polish Science under the Homing Plus programme, co-financed by the European Regional Development Fund.



INNOVATIVE ECONOMY  
NATIONAL COHESION STRATEGY



Foundation for Polish Science



EUROPEAN UNION  
EUROPEAN REGIONAL  
DEVELOPMENT FUND