

# Decision-theoretic Machine Learning

Krzysztof Dembczyński and Wojciech Kotłowski

Intelligent Decision Support Systems Laboratory (IDSS)  
Poznań University of Technology, Poland



Poznań University of Technology, Summer 2015

# Agenda

- 1 Introduction to Machine Learning
- 2 Binary Classification
- 3 Bipartite Ranking
- 4 **Multi-Label Classification**

## Outline

- 1 Multi-label classification
- 2 Simple approaches to multi-label classification
- 3 Beyond simple approaches
- 4 Other task losses
- 5 Rank loss minimization
- 6 Summary

# Outline

- 1 Multi-label classification
- 2 Simple approaches to multi-label classification
- 3 Beyond simple approaches
- 4 Other task losses
- 5 Rank loss minimization
- 6 Summary

## Multi-label classification

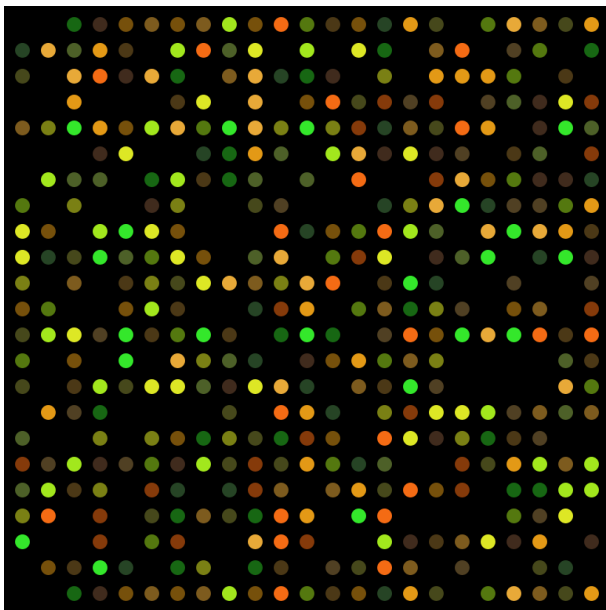
- A classification problem in which we consider **more than one** binary output variable.



Image annotation: cloud? sky? tree?



Ecology: Prediction of the presence or absence of species



Gene function prediction



Stack Overflow is a question and answer site for professional and enthusiast programmers. It's 100% free, no registration required.

[Take the 2-minute tour](#)

Here's how it works:



Anybody can ask  
a question



Anybody can  
answer

Th

## Top Questions

[interesting](#)

[featured](#)
[hot](#)
[week](#)
[month](#)

0

votes

0

answers

1

views

[sending and receiving mails from registered user emailaddresses](#)
[php](#)
[email](#)
[web-applications](#)

asked 34s ago [Angelo A](#) 489

0

votes

0

answers

5

views

[How to create sprites using ConfigParser in Pygame](#)
[python](#)
[pygame](#)

modified 37s ago [Sudoadmin](#) 5

1

votes

1

answers

8

views

[Fortran: possible fibonacci logical error](#)
[fortran](#)
[fibonacci](#)
[fortran95](#)

answered 40s ago [oropendola](#) 326

4

votes

2

answers

1k

views

[Angular - Using one controller for many coherent views across multiple HTTP requests](#)

# Document tagging

## Multi-label classification

- **Multi-label classification:** For a feature vector  $\mathbf{x}$  predict accurately a vector of responses  $\mathbf{y}$  using a function  $\mathbf{h}(\mathbf{x})$ :

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \xrightarrow{\mathbf{h}(\mathbf{x})} \mathbf{y} = (y_1, y_2, \dots, y_m) \in \mathcal{Y} = \{0, 1\}^m$$

## Multi-label classification

- Training data:  $\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ .
- **Predict** a vector  $\mathbf{y} = (y_1, y_2, \dots, y_m)$  for a given  $\mathbf{x}$ .

	$x_1$	$x_2$	$y_1$	$y_2$	$\dots$	$y_m$
$\mathbf{x}_1$	5.0	4.5	1	1		0
$\mathbf{x}_2$	2.0	2.5	0	1		0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	3.0	3.5	0	1		1
$\mathbf{x}$	4.0	2.5	?	?		?

## Multi-label classification

- Training data:  $\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)\}$ .
- **Predict** a vector  $\mathbf{y} = (y_1, y_2, \dots, y_m)$  for a given  $\mathbf{x}$ .

	$x_1$	$x_2$	$y_1$	$y_2$	$\dots$	$y_m$
$\mathbf{x}_1$	5.0	4.5	1	1		0
$\mathbf{x}_2$	2.0	2.5	0	1		0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	3.0	3.5	0	1		1
$\mathbf{x}$	4.0	2.5	1	1		0

## Multi-label classification

- **Example**  $x$  is coming from an unknown input distribution  $P(x)$ .
- **True outcome**  $y$  is generated from  $P(y | x)$ .
- **Predicted outcome** is given by  $\hat{y} = h(x)$ .
- The (**task**) **loss** of a single prediction is  $\ell(y, \hat{y})$ .

## Multi-label classification

- The overall goal is to minimize the **risk**:

$$L_\ell(\mathbf{h}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y})}(\ell(\mathbf{y}, \mathbf{h}(\mathbf{x})))$$

- The optimal prediction function, the so-called **Bayes classifier**, is:

$$\mathbf{h}_\ell^* = \arg \min_{\mathbf{h}} L_\ell(\mathbf{h})$$

- The **regret** of a classifier  $\mathbf{h}$  with respect to  $\ell$  is defined as:

$$\text{Reg}_\ell(\mathbf{h}) = L_\ell(\mathbf{h}) - L_\ell(\mathbf{h}_\ell^*) = L_\ell(\mathbf{h}) - L_\ell^*$$

## Multi-label classification

- We use training examples  $\{\mathbf{x}_i, \mathbf{y}_i\}_1^n$  to find either:
  - ▶ a good approximation of  $\mathbf{h}^*$ , or
  - ▶ a good estimation of  $P(\mathbf{y} | \mathbf{x})$  (or a function of it).
- In the second case, we need to apply an inference procedure to approximate  $\mathbf{h}^*$ .

## Main challenges

- Appropriate modeling of dependencies between labels

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

- A multitude of multivariate loss functions defined over the output vector

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



## Label interdependence

- **Marginal** and **conditional dependence**:

$$P(\mathbf{y}) \neq \prod_{i=1}^m P(y_i) \quad P(\mathbf{y} | \mathbf{x}) \neq \prod_{i=1}^m P(y_i | \mathbf{x})$$

marginal (in)dependence  $\nRightarrow$  conditional (in)dependence

## Label interdependence

- **Marginal dependence  $\nrightarrow$  conditional dependence**

## Label interdependence

- **Marginal dependence  $\nrightarrow$  conditional dependence**

- ▶ Consider two independent labels  $y_1$  and  $y_2$  generated by the same logistic model:

$$P(y_i = 1|\mathbf{x}) = (1 + \exp(-\phi f(\mathbf{x})))^{-1},$$

where  $\phi$  controls the Bayes error rate.

## Label interdependence

- **Marginal dependence  $\nrightarrow$  conditional dependence**

- ▶ Consider two independent labels  $y_1$  and  $y_2$  generated by the same logistic model:

$$P(y_i = 1|\mathbf{x}) = (1 + \exp(-\phi f(\mathbf{x})))^{-1},$$

where  $\phi$  controls the Bayes error rate.

- ▶ Thus, the two labels are conditionally independent, having the conditional distribution:

$$P(\mathbf{y}|\mathbf{x}) = P(y_1|\mathbf{x})P(y_2|\mathbf{x}).$$

# Label interdependence

- **Marginal dependence  $\nrightarrow$  conditional dependence**

- ▶ Consider two independent labels  $y_1$  and  $y_2$  generated by the same logistic model:

$$P(y_i = 1|\mathbf{x}) = (1 + \exp(-\phi f(\mathbf{x})))^{-1},$$

where  $\phi$  controls the Bayes error rate.

- ▶ Thus, the two labels are conditionally independent, having the conditional distribution:

$$P(\mathbf{y}|\mathbf{x}) = P(y_1|\mathbf{x})P(y_2|\mathbf{x}).$$

- ▶ Depending on the value of  $\phi$ , however, they will be stronger or weaker marginally dependent.

# Label interdependence

- **Marginal dependence  $\nrightarrow$  conditional dependence**

- ▶ Consider two independent labels  $y_1$  and  $y_2$  generated by the same logistic model:

$$P(y_i = 1|\mathbf{x}) = (1 + \exp(-\phi f(\mathbf{x})))^{-1},$$

where  $\phi$  controls the Bayes error rate.

- ▶ Thus, the two labels are conditionally independent, having the conditional distribution:

$$P(\mathbf{y}|\mathbf{x}) = P(y_1|\mathbf{x})P(y_2|\mathbf{x}).$$

- ▶ Depending on the value of  $\phi$ , however, they will be stronger or weaker marginally dependent.
- ▶ For  $\phi \rightarrow \infty$  (Bayes error rate tends to 0), the marginal dependence increases towards the deterministic one ( $y_1 = y_2$ ).

## Label interdependence

- Marginal dependence  $\nrightarrow$  conditional dependence

## Label interdependence

- **Marginal dependence  $\not\sim$  conditional dependence**

- ▶ Consider two labels  $y_1$  and  $y_2$  and a single binary feature  $x_1$  with the joint distribution  $P(x_1, y_1, y_2)$  given as:

$x_1$	$y_1$	$y_2$	$P$	$x_1$	$y_1$	$y_2$	$P$
0	0	0	0.25	1	0	0	0
0	0	1	0	1	0	1	0.25
0	1	0	0	1	1	0	0.25
0	1	1	0.25	1	1	1	0



## Label interdependence

- **Marginal dependence  $\not\leftarrow$  conditional dependence**

- ▶ Consider two labels  $y_1$  and  $y_2$  and a single binary feature  $x_1$  with the joint distribution  $P(x_1, y_1, y_2)$  given as:

$x_1$	$y_1$	$y_2$	$P$	$x_1$	$y_1$	$y_2$	$P$
0	0	0	0.25	1	0	0	0
0	0	1	0	1	0	1	0.25
0	1	0	0	1	1	0	0.25
0	1	1	0.25	1	1	1	0

- ▶ The labels are conditionally dependent, since:

$$P(y_1 = 0|x_1 = 1)P(y_2 = 0|x_1 = 1) = 0.5 \times 0.5 = 0.25,$$

but the joint probability is

$$P(y_1 = 0, y_2 = 0|x_1 = 1) = 0.$$

In fact  $y_1 = y_2$  for  $x_1 = 0$  and  $y_2 = 1 - y_1$  for  $x_1 = 1$ .

## Label interdependence

- **Marginal dependence  $\not\leftarrow$  conditional dependence**

- ▶ Consider two labels  $y_1$  and  $y_2$  and a single binary feature  $x_1$  with the joint distribution  $P(x_1, y_1, y_2)$  given as:

$x_1$	$y_1$	$y_2$	$P$	$x_1$	$y_1$	$y_2$	$P$
0	0	0	0.25	1	0	0	0
0	0	1	0	1	0	1	0.25
0	1	0	0	1	1	0	0.25
0	1	1	0.25	1	1	1	0

- ▶ The labels are conditionally dependent, since:

$$P(y_1 = 0|x_1 = 1)P(y_2 = 0|x_1 = 1) = 0.5 \times 0.5 = 0.25,$$

but the joint probability is

$$P(y_1 = 0, y_2 = 0|x_1 = 1) = 0.$$

In fact  $y_1 = y_2$  for  $x_1=0$  and  $y_2 = 1 - y_1$  for  $x_1 = 1$ .

- ▶ However, the labels are marginally independent, since

$$P(y_1) = P(y_2) = 0.5, \text{ and } P(y_1, y_2) = P(y_1)P(y_2).$$

## Label interdependence

- **Deterministic dependencies:**

- ▶ Consider labels  $y_1$  and  $y_2$  with the following conditional distribution for a given  $\mathbf{x}$ :

$$P(y_1 = 1, y_2 = 1 \mid \mathbf{x}) = 1,$$

$$P(y_1 = 0, y_2 = 1 \mid \mathbf{x}) = 0,$$

$$P(y_1 = 1, y_2 = 0 \mid \mathbf{x}) = 0,$$

$$P(y_1 = 0, y_2 = 0 \mid \mathbf{x}) = 0.$$

- ▶ Are  $y_1$  and  $y_1$  conditionally dependent?

## Label interdependence

- **Deterministic dependencies:**

- ▶ Consider labels  $y_1$  and  $y_2$  with the following conditional distribution for a given  $\mathbf{x}$ :

$$P(y_1 = 1, y_2 = 1 \mid \mathbf{x}) = 1,$$

$$P(y_1 = 0, y_2 = 1 \mid \mathbf{x}) = 0,$$

$$P(y_1 = 1, y_2 = 0 \mid \mathbf{x}) = 0,$$

$$P(y_1 = 0, y_2 = 0 \mid \mathbf{x}) = 0.$$

- ▶ Are  $y_1$  and  $y_2$  conditionally dependent?
- ▶ **No**, since it holds:

$$P(y_1, y_2 \mid \mathbf{x}) = P(y_1 \mid \mathbf{x})P(y_2 \mid \mathbf{x})$$

# Label interdependence

- **Model similarities:**

- ▶ Similarities in the structural parts  $g_i(\mathbf{x})$  of the models:

$$f_i(\mathbf{x}) = g_i(\mathbf{x}) + \epsilon_i, \text{ for } i = 1, \dots, m$$

- **Structure** imposed (domain knowledge) on targets

- ▶ Chains,
- ▶ Hierarchies,
- ▶ General graphs,
- ▶ ...

# Label interdependence

- **Interdependence** vs. **hypothesis** and **feature space**:
  - ▶ Regularization constraints the hypothesis space.
  - ▶ Modeling dependencies may increase the expressiveness of the model.
  - ▶ Using a more complex model on individual labels may also help.
  - ▶ Comparison of models is difficult in general, as they differ in many respects (e.g., complexity)!

## Multivariate loss functions

- **Decomposable** and **non-decomposable** losses over **examples**

$$L = \sum_{i=1}^n \ell(\mathbf{y}_i, \mathbf{h}(\mathbf{x}_i)) \quad L \neq \sum_{i=1}^n \ell(\mathbf{y}_i, \mathbf{h}(\mathbf{x}_i))$$

- **Decomposable** and **non-decomposable** losses over **labels**

$$\ell(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \sum_{i=1}^m \ell(y_i, h_i(\mathbf{x})) \quad \ell(\mathbf{y}, \mathbf{h}(\mathbf{x})) \neq \sum_{i=1}^m \ell(y_i, h_i(\mathbf{x}))$$

- Different formulations of loss functions possible:
  - ▶ Set-based losses.
  - ▶ Ranking-based losses.

## Multi-label loss functions

- **Subset 0/1 loss:**  $\ell_{0/1}(\mathbf{y}, \mathbf{h}) = \llbracket \mathbf{y} \neq \mathbf{h} \rrbracket$
- **Hamming loss:**  $\ell_H(\mathbf{y}, \mathbf{h}) = \frac{1}{m} \sum_{i=1}^m \llbracket y_i \neq h_i \rrbracket$
- **F-measure-based loss:**  $\ell_F(\mathbf{y}, \mathbf{h}) = 1 - \frac{2 \sum_{i=1}^m y_i h_i}{\sum_{i=1}^m y_i + \sum_{i=1}^m h_i}$
- **Rank loss:**  $\ell_{\text{rk}}(\mathbf{y}, \mathbf{h}) = w(\mathbf{y}) \sum_{y_i > y_j} \left( \llbracket h_i < h_j \rrbracket + \frac{1}{2} \llbracket h_i = h_j \rrbracket \right)$
- ...



## Relations between losses

- The set-based loss function  $\ell(\mathbf{y}, \mathbf{h})$  should fulfill some basic conditions:
  - ▶  $\ell(\mathbf{y}, \mathbf{h}) = 0$  if and only if  $\mathbf{y} = \mathbf{h}$ .
  - ▶  $\ell(\mathbf{y}, \mathbf{h})$  is maximal when  $y_i \neq h_i$  for every  $i = 1, \dots, m$ .
  - ▶ Should be monotonically non-decreasing with respect to the number of  $y_i \neq h_i$ .

## Relations between losses

- The set-based loss function  $\ell(\mathbf{y}, \mathbf{h})$  should fulfill some basic conditions:
  - ▶  $\ell(\mathbf{y}, \mathbf{h}) = 0$  if and only if  $\mathbf{y} = \mathbf{h}$ .
  - ▶  $\ell(\mathbf{y}, \mathbf{h})$  is maximal when  $y_i \neq h_i$  for every  $i = 1, \dots, m$ .
  - ▶ Should be monotonically non-decreasing with respect to the number of  $y_i \neq h_i$ .
- **In case of deterministic data (no-noise):** the optimal prediction should have the same form for all loss functions and the risk for this prediction should be 0.

## Relations between losses

- The set-based loss function  $\ell(\mathbf{y}, \mathbf{h})$  should fulfill some basic conditions:
  - ▶  $\ell(\mathbf{y}, \mathbf{h}) = 0$  if and only if  $\mathbf{y} = \mathbf{h}$ .
  - ▶  $\ell(\mathbf{y}, \mathbf{h})$  is maximal when  $y_i \neq h_i$  for every  $i = 1, \dots, m$ .
  - ▶ Should be monotonically non-decreasing with respect to the number of  $y_i \neq h_i$ .
- **In case of deterministic data (no-noise):** the optimal prediction should have the same form for all loss functions and the risk for this prediction should be 0.
- **In case of non-deterministic data (noise):** the optimal prediction and its risk can be different for different losses.

## Learning and inference with multi-label losses

- The loss functions, like Hamming loss or subset 0/1 loss, often referred to as **task losses**, are usually neither convex nor differentiable.

## Learning and inference with multi-label losses

- The loss functions, like Hamming loss or subset 0/1 loss, often referred to as **task losses**, are usually neither convex nor differentiable.
- Therefore learning is a hard optimization problem.

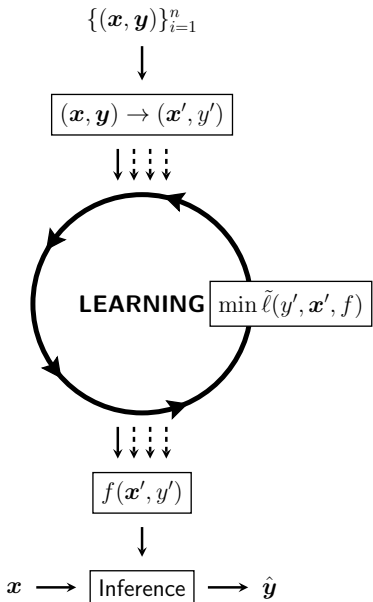
## Learning and inference with multi-label losses

- The loss functions, like Hamming loss or subset 0/1 loss, often referred to as **task losses**, are usually neither convex nor differentiable.
- Therefore learning is a hard optimization problem.
- Two approaches try to make this task easier
  - ▶ Reduction.
  - ▶ Surrogate loss minimization.

## Learning and inference with multi-label losses

- The loss functions, like Hamming loss or subset 0/1 loss, often referred to as **task losses**, are usually neither convex nor differentiable.
- Therefore learning is a hard optimization problem.
- Two approaches try to make this task easier
  - ▶ Reduction.
  - ▶ Surrogate loss minimization.
- Two phases in solving multi-label problems:
  - ▶ Learning: Estimate parameters of a scoring function  $f(\mathbf{x}, \mathbf{y})$ .
  - ▶ Inference: Use the scoring function  $f(\mathbf{x}, \mathbf{y})$  to classify new instances by finding the best  $\mathbf{y}$  for a given  $\mathbf{x}$ .

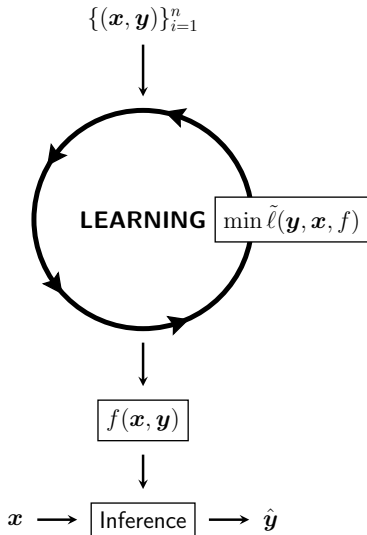
## Reduction



- **Reduce** the original problem into simple problems, for which efficient algorithmic solutions are available.
- Reduction to one or a sequence of problems.
- Plug-in rule classifiers.



## Structured loss minimization



- Replace the task loss by a **surrogate loss** that is easier to cope with.
- Surrogate loss is typically a differentiable approximation of the task loss or a convex upper bound of it.

## Statistical consistency

- Analysis of algorithms in terms of their infinite sample performance.<sup>1</sup>
- We say that a proxy loss  $\tilde{\ell}$  is **consistent** (**calibrated**) with the task loss  $\ell$  when the following holds:

$$\text{Reg}_{\tilde{\ell}}(\mathbf{h}) \rightarrow 0 \Rightarrow \text{Reg}_{\ell}(\mathbf{h}) \rightarrow 0.$$

---

<sup>1</sup> A. Tewari and P.L. Bartlett. On the consistency of multiclass classification methods. *JMLR*, 8:1007–1025, 2007

D. McAllester and J. Keshet. Generalization bounds and consistency for latent structural probit and ramp loss. In *NIPS*, pages 2205–2212, 2011

W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199-200:22–44, 2013

## Statistical consistency

- Analysis of algorithms in terms of their infinite sample performance.<sup>1</sup>
- We say that a proxy loss  $\tilde{\ell}$  is **consistent** (**calibrated**) with the task loss  $\ell$  when the following holds:

$$\text{Reg}_{\tilde{\ell}}(\mathbf{h}) \rightarrow 0 \Rightarrow \text{Reg}_{\ell}(\mathbf{h}) \rightarrow 0.$$

- The definition concerns both surrogate loss minimization and reduction:
  - ▶ Surrogate loss minimization:  $\tilde{\ell}$  = surrogate loss.
  - ▶ Reduction:  $\tilde{\ell}$  = loss used in the reduced problem.

---

<sup>1</sup> A. Tewari and P.L. Bartlett. On the consistency of multiclass classification methods. *JMLR*, 8:1007–1025, 2007

D. McAllester and J. Keshet. Generalization bounds and consistency for latent structural probit and ramp loss. In *NIPS*, pages 2205–2212, 2011

W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199-200:22–44, 2013

# Outline

- 1 Multi-label classification
- 2 Simple approaches to multi-label classification**
- 3 Beyond simple approaches
- 4 Other task losses
- 5 Rank loss minimization
- 6 Summary

## Basic reductions: Binary relevance

- **Binary relevance:** Decomposes the problem to  $m$  binary classification problems:

$$(\mathbf{x}, \mathbf{y}) \longrightarrow (\mathbf{x}, y = y_i), \quad i = 1, \dots, m$$

	$X_1$	$X_2$	$Y_1$	$Y_2$	$\dots$	$Y_m$
$\mathbf{x}_1$	5.0	4.5	1	1		0
$\mathbf{x}_2$	2.0	2.5	0	1		0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	3.0	3.5	0	1		1

## Basic reductions: Binary relevance

- **Binary relevance:** Decomposes the problem to  $m$  binary classification problems:

$$(\mathbf{x}, \mathbf{y}) \longrightarrow (\mathbf{x}, y = y_i), \quad i = 1, \dots, m$$

	$X_1$	$X_2$	$Y_1$	$Y_2$	$\dots$	$Y_m$
$\mathbf{x}_1$	5.0	4.5	1	1		0
$\mathbf{x}_2$	2.0	2.5	0	1		0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	3.0	3.5	0	1		1

- Seems to be very simplistic.
- Ignores any dependencies.
- Is it good for any loss function?

## Basic reductions: Label powerset

- **Label powerset:** Treats each label combination as a new meta-class in multi-class classification:

$$(\mathbf{x}, \mathbf{y}) \longrightarrow (\mathbf{x}, y = \text{metaclass}(\mathbf{y}))$$

	$X_1$	$X_2$	$Y_1$	$Y_2$	$\dots$	$Y_m$
$\mathbf{x}_1$	5.0	4.5	1	1		0
$\mathbf{x}_2$	2.0	2.5	0	1		0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	3.0	3.5	0	1		1

## Basic reductions: Label powerset

- **Label powerset:** Treats each label combination as a new meta-class in multi-class classification:

$$(\mathbf{x}, \mathbf{y}) \longrightarrow (\mathbf{x}, y = \text{metaclass}(\mathbf{y}))$$

	$X_1$	$X_2$	$Y_1$	$Y_2$	$\dots$	$Y_m$
$\mathbf{x}_1$	5.0	4.5	1	1		0
$\mathbf{x}_2$	2.0	2.5	0	1		0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	3.0	3.5	0	1		1

- Any multi-class classification algorithm can be used, but the number of classes is huge.
- Takes other labels into account, but ignores internal structure of classes (label vectors).



**What about task losses minimized by BR and LP?**

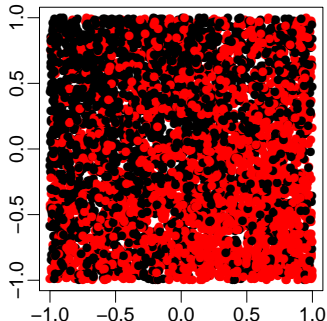
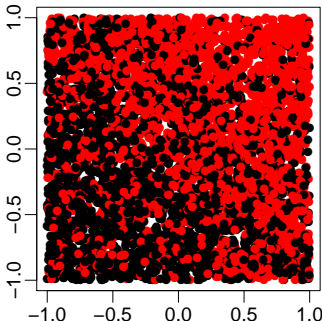
## Synthetic data

- Two independent models:

$$f_1(\mathbf{x}) = \frac{1}{2}x_1 + \frac{1}{2}x_2, \quad f_2(\mathbf{x}) = \frac{1}{2}x_1 - \frac{1}{2}x_2$$

- Logistic model to get labels:

$$P(y_i = 1) = \frac{1}{1 + \exp(-2f_i)}$$



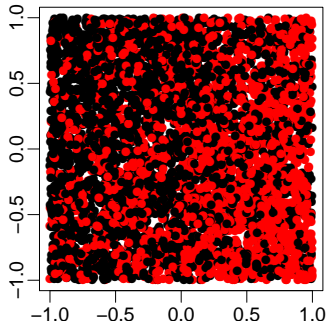
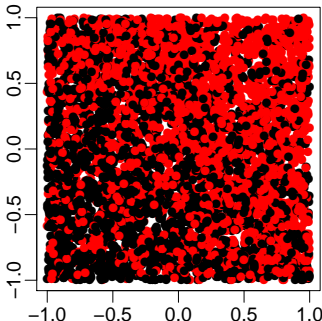
## Synthetic data

- Two dependent models:

$$f_1(\mathbf{x}) = \frac{1}{2}x_1 + \frac{1}{2}x_2 \quad f_2(y_1, \mathbf{x}) = y_1 + \frac{1}{2}x_1 - \frac{1}{2}x_2 - \frac{2}{3}$$

- Logistic model to get labels:

$$P(y_i = 1) = \frac{1}{1 + \exp(-2f_i)}$$



## Results for two performance measures

- Hamming loss:  $\ell_H(\mathbf{y}, \mathbf{h}) = \frac{1}{m} \sum_{i=1}^m \llbracket y_i \neq h_i \rrbracket$ ,
- Subset 0/1 loss:  $\ell_{0/1}(\mathbf{y}, \mathbf{h}) = \llbracket \mathbf{y} \neq \mathbf{h} \rrbracket$ .

CONDITIONAL INDEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR		
LP LR		
CONDITIONAL DEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR		
LP LR		

## Results for two performance measures

- Hamming loss:  $\ell_H(\mathbf{y}, \mathbf{h}) = \frac{1}{m} \sum_{i=1}^m \llbracket y_i \neq h_i \rrbracket$ ,
- Subset 0/1 loss:  $\ell_{0/1}(\mathbf{y}, \mathbf{h}) = \llbracket \mathbf{y} \neq \mathbf{h} \rrbracket$ .

CONDITIONAL INDEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR	0.4232	0.6723
LP LR	0.4232	0.6725

CONDITIONAL DEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR		
LP LR		

## Results for two performance measures

- Hamming loss:  $\ell_H(\mathbf{y}, \mathbf{h}) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}[y_i \neq h_i]$ ,
- Subset 0/1 loss:  $\ell_{0/1}(\mathbf{y}, \mathbf{h}) = \mathbb{I}[\mathbf{y} \neq \mathbf{h}]$ .

CONDITIONAL INDEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR	0.4232	0.6723
LP LR	0.4232	0.6725

CONDITIONAL DEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR	0.3470	0.5499
LP LR	0.3610	0.5146

## Linear + XOR synthetic data

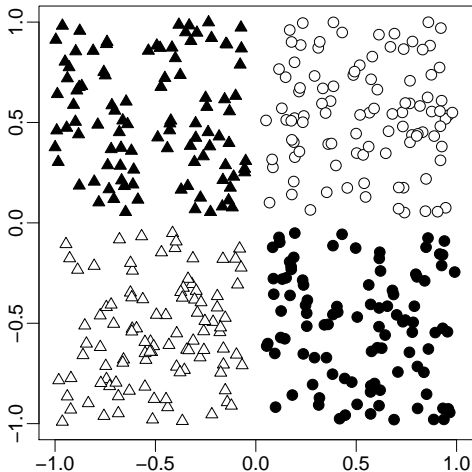


Figure : Problem with two targets: shapes ( $\triangle$  vs.  $\circ$ ) and colors ( $\square$  vs.  $\blacksquare$ ).

## Linear + XOR synthetic data

CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR	0.2399( $\pm$ .0097)	0.4751( $\pm$ .0196)
LP LR	0.0143( $\pm$ .0020)	0.0195( $\pm$ .0011)
BAYES OPTIMAL	0	0

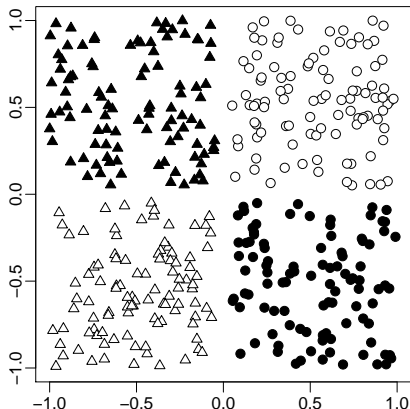


## Linear + XOR synthetic data

CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR	0.2399( $\pm$ .0097)	0.4751( $\pm$ .0196)
LP LR	0.0143( $\pm$ .0020)	0.0195( $\pm$ .0011)
<b>BR MLRules</b>	<b>0.0011(<math>\pm</math>.0002)</b>	<b>0.0020(<math>\pm</math>.0003)</b>
BAYES OPTIMAL	0	0

## Linear + XOR synthetic data

- BR LR uses two linear classifiers: cannot handle the label color ( $\square$  vs.  $\blacksquare$ ) – the XOR problem.
- LP LR uses four linear classifiers to solve 4-class problem ( $\triangle$ ,  $\blacktriangle$ ,  $\circ$ ,  $\bullet$ ): extends the hypothesis space.
- BR MLRules uses two non-linear classifiers (based on decision rules): XOR problem is not a problem.
- There is no noise in the data.
- Easy to perform unfair comparison.



## Multi-label loss functions

- The conditional risk in multi-label classification of  $\mathbf{h}$  at  $\mathbf{x}$ :

$$L_\ell(\mathbf{h} \mid \mathbf{x}) = \mathbb{E}_{\mathbf{y}} [\ell(\mathbf{y}, \mathbf{h}(\mathbf{x}))] = \sum_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} \mid \mathbf{x}) \ell(\mathbf{y}, \mathbf{h}(\mathbf{x}))$$

- The risk-minimizing classifier for a given  $\mathbf{x}$ :

$$\mathbf{h}^*(\mathbf{x}) = \arg \min_{\mathbf{h}} L_\ell(\mathbf{h} \mid \mathbf{x})$$

- Let us start with Hamming loss and subset 0/1 loss ...<sup>2</sup>

---

<sup>2</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. On loss minimization and label dependence in multi-label classification. *Machine Learning*, 88:5–45, 2012

## Hamming loss vs. subset 0/1 loss

- The risk minimizer for the Hamming loss is

## Hamming loss vs. subset 0/1 loss

- The risk minimizer for the Hamming loss is the **marginal mode**:

$$h_i^*(\boldsymbol{x}) = \arg \max_{y_i \in \{0,1\}} P(y_i | \boldsymbol{x}), \quad i = 1, \dots, m,$$

## Hamming loss vs. subset 0/1 loss

- The risk minimizer for the Hamming loss is the **marginal mode**:

$$h_i^*(\mathbf{x}) = \arg \max_{y_i \in \{0,1\}} P(y_i | \mathbf{x}), \quad i = 1, \dots, m,$$

while for the subset 0/1 loss is

## Hamming loss vs. subset 0/1 loss

- The risk minimizer for the Hamming loss is the **marginal mode**:

$$h_i^*(\mathbf{x}) = \arg \max_{y_i \in \{0,1\}} P(y_i | \mathbf{x}), \quad i = 1, \dots, m,$$

while for the subset 0/1 loss is the **joint mode**:

$$\mathbf{h}^*(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x}).$$

## Hamming loss vs. subset 0/1 loss

- The risk minimizer for the Hamming loss is the **marginal mode**:

$$h_i^*(\mathbf{x}) = \arg \max_{y_i \in \{0,1\}} P(y_i | \mathbf{x}), \quad i = 1, \dots, m,$$

while for the subset 0/1 loss is the **joint mode**:

$$\mathbf{h}^*(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x}).$$

- Marginal mode vs. joint mode.

$\mathbf{y}$	$P(\mathbf{y})$	
0 0 0 0	0.30	
0 1 1 1	0.17	Marginal mode: 1 1 1 1
1 0 1 1	0.18	Joint mode: 0 0 0 0
1 1 0 1	0.17	
1 1 1 0	0.18	



## Equivalence of risk minimizers and mutual risk bounds

- The risk minimizers for  $\ell_H$  and  $\ell_{0/1}$  are **equivalent**,

$$\mathbf{h}_H^*(\mathbf{x}) = \mathbf{h}_{0/1}^*(\mathbf{x}),$$

under specific conditions, for example, when:

## Equivalence of risk minimizers and mutual risk bounds

- The risk minimizers for  $\ell_H$  and  $\ell_{0/1}$  are **equivalent**,

$$\mathbf{h}_H^*(\mathbf{x}) = \mathbf{h}_{0/1}^*(\mathbf{x}),$$

under specific conditions, for example, when:

- ▶ Targets  $y_1, \dots, y_m$  are conditionally independent, i.e,

$$P(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^m P(y_i|\mathbf{x}).$$

## Equivalence of risk minimizers and mutual risk bounds

- The risk minimizers for  $\ell_H$  and  $\ell_{0/1}$  are **equivalent**,

$$\mathbf{h}_H^*(\mathbf{x}) = \mathbf{h}_{0/1}^*(\mathbf{x}),$$

under specific conditions, for example, when:

- ▶ Targets  $y_1, \dots, y_m$  are conditionally independent, i.e,

$$P(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^m P(y_i|\mathbf{x}).$$

- ▶ The probability of the joint mode satisfies

$$P(\mathbf{h}_{0/1}^*(\mathbf{x})|\mathbf{x}) > 0.5.$$

## Equivalence of risk minimizers and mutual risk bounds

- The risk minimizers for  $\ell_H$  and  $\ell_{0/1}$  are **equivalent**,

$$\mathbf{h}_H^*(\mathbf{x}) = \mathbf{h}_{0/1}^*(\mathbf{x}),$$

under specific conditions, for example, when:

- Targets  $y_1, \dots, y_m$  are conditionally independent, i.e.,

$$P(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^m P(y_i|\mathbf{x}).$$

- The probability of the joint mode satisfies

$$P(\mathbf{h}_{0/1}^*(\mathbf{x})|\mathbf{x}) > 0.5.$$

- The following bounds hold for any  $P(\mathbf{y}|\mathbf{x})$  and  $\mathbf{h}$ :

$$\frac{1}{m}L_{0/1}(\mathbf{h}|\mathbf{x}) \leq L_H(\mathbf{h}|\mathbf{x}) \leq L_{0/1}(\mathbf{h}|\mathbf{x})$$

## Regret analysis

- The previous results may suggest that one of the loss functions can be used as a proxy (surrogate) for the other:
  - ▶ For some situations both risk minimizers coincide.
  - ▶ One can provide mutual bounds for both loss functions.

## Regret analysis

- The previous results may suggest that one of the loss functions can be used as a proxy (surrogate) for the other:
  - ▶ For some situations both risk minimizers coincide.
  - ▶ One can provide mutual bounds for both loss functions.
- **However, the regret analysis of the worst case shows that minimization of the subset 0/1 loss may result in a large error for the Hamming loss and vice versa.**

## Regret analysis

- The **regret** of a classifier with respect to  $\ell$  is defined as:

$$\text{Reg}_\ell(\mathbf{h}) = L_\ell(\mathbf{h}) - L_\ell(\mathbf{h}_\ell^*),$$

where  $\mathbf{h}_\ell^*$  is the Bayes classifier for a given loss  $\ell$ .

- Regret measures how worse is  $\mathbf{h}$  by comparison with the optimal classifier for a given loss.
- To simplify the analysis we will consider the conditional regret:

$$\text{Reg}_\ell(\mathbf{h} \mid \mathbf{x}) = L_\ell(\mathbf{h} \mid \mathbf{x}) - L_\ell(\mathbf{h}_\ell^* \mid \mathbf{x}).$$

- We will analyze the regret between:
  - ▶ the Bayes classifier for Hamming loss  $\mathbf{h}_H^*$
  - ▶ the Bayes classifier for subset 0/1 loss  $\mathbf{h}_{0/1}^*$

with respect to both functions.

- It is a bit an unusual analysis.

## Regret analysis

- The following **upper bound** holds:

$$\text{Reg}_{0/1}(\mathbf{h}_H^* | \mathbf{x}) = L_{0/1}(\mathbf{h}_H^* | \mathbf{x}) - L_{0/1}(\mathbf{h}_{0/1}^* | \mathbf{x}) < 0.5$$

- Moreover, this **bound is tight**.
- Example:**

$\mathbf{y}$	$P(\mathbf{y})$
0 0 0 0	0.02
0 0 1 1	0.49
1 1 0 0	0.49

Marginal mode:                      0 0 0 0  
Joint mode:                      0 0 1 1 or 1 1 0 0



## Regret analysis

- The following **upper bound** holds  $m > 3$ :

$$\text{Reg}_H(\mathbf{h}_{0/1}^* | \mathbf{x}) = L_H(\mathbf{h}_{0/1}^* | \mathbf{x}) - L_H(\mathbf{h}_H^* | \mathbf{x}) < \frac{m-2}{m+2}$$

- Moreover, this **bound is tight**.
- Example:**

$y$	$P(y)$
0 0 0 0	0.170
0 1 1 1	0.166
1 0 1 1	0.166
1 1 0 1	0.166
1 1 1 0	0.166
1 1 1 1	0.166

Marginal mode: 1 1 1 1  
Joint mode: 0 0 0 0

## Relations between losses

- The risk minimizers of Hamming and subset 0/1 loss are different: marginal mode vs. joint mode.

## Relations between losses

- The risk minimizers of Hamming and subset 0/1 loss are different: marginal mode vs. joint mode.
- Under specific conditions, like label independence or high probability of the joint mode ( $> 0.5$ ), these two risk minimizers are equivalent.

## Relations between losses

- The risk minimizers of Hamming and subset 0/1 loss are different: marginal mode vs. joint mode.
- Under specific conditions, like label independence or high probability of the joint mode ( $> 0.5$ ), these two risk minimizers are equivalent.
- The risks of these loss functions are mutually upper bounded.

## Relations between losses

- The risk minimizers of Hamming and subset 0/1 loss are different: marginal mode vs. joint mode.
- Under specific conditions, like label independence or high probability of the joint mode ( $> 0.5$ ), these two risk minimizers are equivalent.
- The risks of these loss functions are mutually upper bounded.
- Minimization of the subset 0/1 loss may cause a high regret for the Hamming loss and vice versa.

## BR vs. LP

## BR vs. LP

- Binary relevance (BR)

## BR vs. LP

- Binary relevance (BR)
  - ▶ BR is **consistent** for Hamming loss **without** any additional assumptions on **label (in)dependence**.



## BR vs. LP

- Binary relevance (BR)
  - ▶ BR is **consistent** for Hamming loss **without** any additional assumptions on **label (in)dependence**.
  - ▶ If this would not be true, then **we could not optimally solve binary classification problems!!!**

## BR vs. LP

- Binary relevance (BR)
  - ▶ BR is **consistent** for Hamming loss **without** any additional assumptions on **label (in)dependence**.
  - ▶ If this would not be true, then **we could not optimally solve binary classification problems!!!**
  - ▶ For other losses, one should take **additional assumptions**:
    - For subset 0/1 loss: label independence, high probability of the joint mode ( $> 0.5$ ), ...

## BR vs. LP

- Binary relevance (BR)
  - ▶ BR is **consistent** for Hamming loss **without** any additional assumptions on **label (in)dependence**.
  - ▶ If this would not be true, then **we could not optimally solve binary classification problems!!!**
  - ▶ For other losses, one should take **additional assumptions**:
    - For subset 0/1 loss: label independence, high probability of the joint mode ( $> 0.5$ ), ...
  - ▶ Learning and inference is **linear** in  $m$  (however, faster algorithms exist).

## BR vs. LP

- Label powerset (LP)

## BR vs. LP

- Label powerset (LP)
  - ▶ LP is **consistent** for the subset 0/1 loss.

## BR vs. LP

- Label powerset (LP)
  - ▶ LP is **consistent** for the subset 0/1 loss.
  - ▶ In its basic formulation it is **not consistent** for Hamming loss.

## BR vs. LP

- Label powerset (LP)
  - ▶ LP is **consistent** for the subset 0/1 loss.
  - ▶ In its basic formulation it is **not consistent** for Hamming loss.
  - ▶ However, if used with a probabilistic multi-class classifier, it estimates the joint conditional distribution for a given  $x$ : inference for **any loss** would be then possible.

## BR vs. LP

- Label powerset (LP)
  - ▶ LP is **consistent** for the subset 0/1 loss.
  - ▶ In its basic formulation it is **not consistent** for Hamming loss.
  - ▶ However, if used with a probabilistic multi-class classifier, it estimates the joint conditional distribution for a given  $x$ : inference for **any loss** would be then possible.
  - ▶ Similarly, by reducing to cost-sensitive multi-class classification LP can be used with **almost any loss function**.



## BR vs. LP

- Label powerset (LP)
  - ▶ LP is **consistent** for the subset 0/1 loss.
  - ▶ In its basic formulation it is **not consistent** for Hamming loss.
  - ▶ However, if used with a probabilistic multi-class classifier, it estimates the joint conditional distribution for a given  $x$ : inference for **any loss** would be then possible.
  - ▶ Similarly, by reducing to cost-sensitive multi-class classification LP can be used with **almost any loss function**.
  - ▶ LP may gain from the implicit expansion of the **feature** or **hypothesis space**.

## BR vs. LP

- Label powerset (LP)
  - ▶ LP is **consistent** for the subset 0/1 loss.
  - ▶ In its basic formulation it is **not consistent** for Hamming loss.
  - ▶ However, if used with a probabilistic multi-class classifier, it estimates the joint conditional distribution for a given  $x$ : inference for **any loss** would be then possible.
  - ▶ Similarly, by reducing to cost-sensitive multi-class classification LP can be used with **almost any loss function**.
  - ▶ LP may gain from the implicit expansion of the **feature** or **hypothesis space**.
  - ▶ Unfortunately, learning and inference is basically **exponential** in  $m$  (however, this complexity is constrained by the number of training examples).

## Relations between losses

- Both are commonly used.
- Hamming loss:
  - ▶ Not too many labels.
  - ▶ Well-balanced labels.
  - ▶ **Application**: Gene function prediction.
- Subset 0/1 loss:
  - ▶ Very restrictive.
  - ▶ Small number of labels.
  - ▶ Low noise problems.
  - ▶ **Application**: Prediction of diseases of a patient.



# Outline

- 1 Multi-label classification
- 2 Simple approaches to multi-label classification
- 3 Beyond simple approaches**
- 4 Other task losses
- 5 Rank loss minimization
- 6 Summary

## Beyond LP

- **Classical multi-class classification algorithms:**
  - ▶  $k$ -nearest neighbors,
  - ▶ Decision trees,
  - ▶ Logistic regression,
  - ▶ Multi-class SVMs,
  - ▶ ...
- **Reduction algorithms:**
  - ▶ 1 vs All,
  - ▶ 1 vs 1 and Weighted All-Pairs (WAP),
  - ▶ Directed acyclic graphs (DAG),
  - ▶ ECOC, PECOC, SECOC,
  - ▶ Filter Trees,
  - ▶ Conditional Probability Trees,
  - ▶ ...
- **Can we adapt these algorithms to multi-label classification and different task losses in a more direct way?**

## Beyond LP

- Naive reduction to 1 vs. All:

$$(x, y) \longrightarrow (x, y = \text{metaclass}(y))$$

## Beyond LP

- Naive reduction to 1 vs. All:

$$(\mathbf{x}, \mathbf{y}) \longrightarrow (\mathbf{x}, y = \text{metaclass}(\mathbf{y}))$$

- Reduction of **multi-class classification** to **binary classification**:

$$(\mathbf{x}, y = \text{metaclass}(\mathbf{y})) \longrightarrow \{(\mathbf{x}, y, 1)\} \cup \{(\mathbf{x}, y', 0) : \forall y' \neq y\}$$

## Beyond LP

- Naive reduction to 1 vs. All:

$$(x, y) \longrightarrow (x, y = \text{metaclass}(y))$$

- Reduction of **multi-class classification** to **binary classification**:

$$(x, y = \text{metaclass}(y)) \longrightarrow \{(x, y, 1)\} \cup \{(x, y', 0) : \forall y' \neq y\}$$

- But we can reduce directly **multi-label classification** to **binary classification**:

$$(x, y) \longrightarrow \{(x, y, 1)\} \cup \{(x, y', 0) : \forall y' \neq y\}$$



## Beyond LP

- Naive reduction to 1 vs. All:

$$(x, y) \longrightarrow (x, y = \text{metaclass}(y))$$

- Reduction of **multi-class classification** to **binary classification**:

$$(x, y = \text{metaclass}(y)) \longrightarrow \{(x, y, 1)\} \cup \{(x, y', 0) : \forall y' \neq y\}$$

- But we can reduce directly **multi-label classification** to **binary classification**:

$$(x, y) \longrightarrow \{(x, y, 1)\} \cup \{(x, y', 0) : \forall y' \neq y\}$$

- **We can exploit now the internal structure of label vectors!!!**

## Internal structure of classes

- The model can be given by a **scoring function**  $f(x, y)$ .

## Internal structure of classes

- The model can be given by a **scoring function**  $f(\mathbf{x}, \mathbf{y})$ .
- Different forms of  $f(\mathbf{x}, \mathbf{y})$  are possible, for example:

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

## Internal structure of classes

- The model can be given by a **scoring function**  $f(\mathbf{x}, \mathbf{y})$ .
- Different forms of  $f(\mathbf{x}, \mathbf{y})$  are possible, for example:

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

where the second term models pairwise interactions.

## Internal structure of classes

- The model can be given by a **scoring function**  $f(\mathbf{x}, \mathbf{y})$ .
- Different forms of  $f(\mathbf{x}, \mathbf{y})$  are possible, for example:

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

where the second term models pairwise interactions.

- Prediction is given by:

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

## Internal structure of classes

- Generalization of logistic regression and SVMs for  $f(\mathbf{x}, \mathbf{y})$ :
  - ▶ Conditional random fields (CRFs),<sup>3</sup>
  - ▶ Structured support vector machines (SSVMs).<sup>4</sup>

---

<sup>3</sup> John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *ICML*, pages 282–289, 2001

<sup>4</sup> Y. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun. Large margin methods for structured and interdependent output variables. *JMLR*, 6:1453–1484, 2005

## CRFs and SSVMs

- CRFs use **logistic loss** as a surrogate:

$$\tilde{\ell}_{\log}(\mathbf{y}, \mathbf{x}, f) = -\log P(\mathbf{y}|\mathbf{x}) = \log \left( \sum_{\mathbf{y} \in \mathcal{Y}} \exp(f(\mathbf{x}, \mathbf{y})) \right) - f(\mathbf{x}, \mathbf{y}).$$

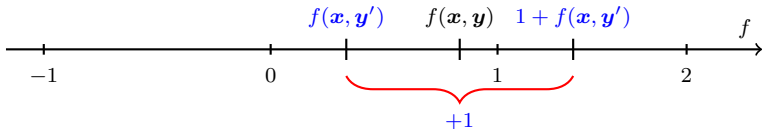
## CRFs and SSVMs

- CRFs use **logistic loss** as a surrogate:

$$\tilde{\ell}_{\log}(\mathbf{y}, \mathbf{x}, f) = -\log P(\mathbf{y}|\mathbf{x}) = \log \left( \sum_{\mathbf{y} \in \mathcal{Y}} \exp(f(\mathbf{x}, \mathbf{y})) \right) - f(\mathbf{x}, \mathbf{y}).$$

- SSVMs minimize the **structured hinge loss**:

$$\tilde{\ell}_h(\mathbf{y}, \mathbf{x}, f) = \max_{\mathbf{y}' \in \mathcal{Y}} \{ \mathbb{I}[\mathbf{y}' \neq \mathbf{y}] + f(\mathbf{x}, \mathbf{y}') \} - f(\mathbf{x}, \mathbf{y}).$$





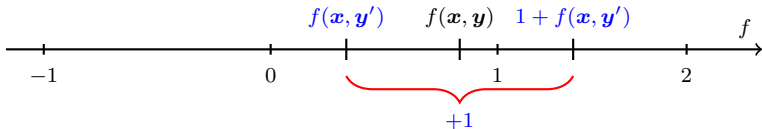
## CRFs and SSVMs

- CRFs use **logistic loss** as a surrogate:

$$\tilde{\ell}_{\log}(\mathbf{y}, \mathbf{x}, f) = -\log P(\mathbf{y}|\mathbf{x}) = \log \left( \sum_{\mathbf{y} \in \mathcal{Y}} \exp(f(\mathbf{x}, \mathbf{y})) \right) - f(\mathbf{x}, \mathbf{y}).$$

- SSVMs minimize the **structured hinge loss**:

$$\tilde{\ell}_h(\mathbf{y}, \mathbf{x}, f) = \max_{\mathbf{y}' \in \mathcal{Y}} \{ \mathbb{I}[\mathbf{y}' \neq \mathbf{y}] + f(\mathbf{x}, \mathbf{y}') \} - f(\mathbf{x}, \mathbf{y}).$$



- SSVMs and CRFs are quite similar to each other:
  - ▶ max vs. soft-max
  - ▶ margin vs. no-margin

## CRFs and SSVMs

- Follow the general LP strategy, but can exploit the internal structure of classes within scoring function  $f(\mathbf{x}, \mathbf{y})$ .
- Convex optimization problem, but its hardness depends on the structure of  $f(\mathbf{x}, \mathbf{y})$ .
- Similarly, the inference (also known as decoding problem) is hard in the general case.
- For sequence and tree structures, the problem can be solved in polynomial time.

## CRFs and SSVMs for different task losses

- In SSVMs, task loss  $\ell(\mathbf{y}, \mathbf{y}')$  can be used for **margin rescaling**:

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

---

<sup>5</sup> W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199-200:22–44, 2013

A. Tewari and P.L. Bartlett. On the consistency of multiclass classification methods. *JMLR*, 8:1007–1025, 2007

D. McAllester. *Generalization Bounds and Consistency for Structured Labeling in Predicting Structured Data*. MIT Press, 2007

## CRFs and SSVMs for different task losses

- In SSVMs, task loss  $\ell(\mathbf{y}, \mathbf{y}')$  can be used for **margin rescaling**:

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

- SSVMs with Hamming loss and

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

decompose to BR with SVMs.

---

<sup>5</sup> W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199-200:22–44, 2013

A. Tewari and P.L. Bartlett. On the consistency of multiclass classification methods. *JMLR*, 8:1007–1025, 2007

D. McAllester. *Generalization Bounds and Consistency for Structured Labeling in Predicting Structured Data*. MIT Press, 2007

## CRFs and SSVMs for different task losses

- In SSVMs, task loss  $\ell(\mathbf{y}, \mathbf{y}')$  can be used for **margin rescaling**:

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

- SSVMs with Hamming loss and

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

decompose to BR with SVMs.

### Question

**Prove** why this is true.

---

<sup>5</sup> W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199-200:22–44, 2013

A. Tewari and P.L. Bartlett. On the consistency of multiclass classification methods. *JMLR*, 8:1007–1025, 2007

D. McAllester. *Generalization Bounds and Consistency for Structured Labeling in Predicting Structured Data*. MIT Press, 2007

## CRFs and SSVMs for different task losses

- In SSVMs, task loss  $\ell(\mathbf{y}, \mathbf{y}')$  can be used for **margin rescaling**:

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

- SSVMs with Hamming loss and

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

decompose to BR with SVMs.

### Question

**Prove** why this is true.

- In general SSVMs are inconsistent.<sup>5</sup>

---

<sup>5</sup> W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199-200:22–44, 2013

A. Tewari and P.L. Bartlett. On the consistency of multiclass classification methods. *JMLR*, 8:1007–1025, 2007

D. McAllester. *Generalization Bounds and Consistency for Structured Labeling in Predicting Structured Data*. MIT Press, 2007

## CRFs and SSVMs for different task losses

- CRFs are tailored for the subset 0/1 loss and cannot directly take other task losses into account.

---

<sup>6</sup> P. Pletscher, C.S. Ong, and J.M. Buhmann. Entropy and margin maximization for structured output learning. In *ECML/PKDD*. Springer, 2010

Q. Shi, M. Reid, and T. Caetano. Hybrid model of conditional random field and support vector machine. In *Workshop at NIPS*, 2009

K. Gimpel and N. Smith. Softmax-margin crfs: Training log-linear models with cost functions. In *HLT*, pages 733–736, 2010

## CRFs and SSVMs for different task losses

- CRFs are tailored for the subset 0/1 loss and cannot directly take other task losses into account.
- CRFs with the scoring function of the form

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

minimize Hamming loss ( $\rightarrow$  BR with logistic regression).

---

6

P. Pletscher, C.S. Ong, and J.M. Buhmann. Entropy and margin maximization for structured output learning. In *ECML/PKDD*. Springer, 2010

Q. Shi, M. Reid, and T. Caetano. Hybrid model of conditional random field and support vector machine. In *Workshop at NIPS*, 2009

K. Gimpel and N. Smith. Softmax-margin crfs: Training log-linear models with cost functions. In *HLT*, pages 733–736, 2010



## CRFs and SSVMs for different task losses

- CRFs are tailored for the subset 0/1 loss and cannot directly take other task losses into account.
- CRFs with the scoring function of the form

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

minimize Hamming loss ( $\rightarrow$  BR with logistic regression).

### Question

**Prove** why this is true.

- 
- <sup>6</sup> P. Pletscher, C.S. Ong, and J.M. Buhmann. Entropy and margin maximization for structured output learning. In *ECML/PKDD*. Springer, 2010
- Q. Shi, M. Reid, and T. Caetano. Hybrid model of conditional random field and support vector machine. In *Workshop at NIPS*, 2009
- K. Gimpel and N. Smith. Softmax-margin crfs: Training log-linear models with cost functions. In *HLT*, pages 733–736, 2010

## CRFs and SSVMs for different task losses

- CRFs are tailored for the subset 0/1 loss and cannot directly take other task losses into account.
- CRFs with the scoring function of the form

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

minimize Hamming loss ( $\rightarrow$  BR with logistic regression).

### Question

**Prove** why this is true.

- Some works on incorporating margin into CRFs.<sup>6</sup>

<sup>6</sup> P. Pletscher, C.S. Ong, and J.M. Buhmann. Entropy and margin maximization for structured output learning. In *ECML/PKDD*. Springer, 2010

Q. Shi, M. Reid, and T. Caetano. Hybrid model of conditional random field and support vector machine. In *Workshop at NIPS*, 2009

K. Gimpel and N. Smith. Softmax-margin crfs: Training log-linear models with cost functions. In *HLT*, pages 733–736, 2010

## SSVMs vs. BR

Table : SSVMs with pairwise term<sup>7</sup> vs. BR with LR<sup>8</sup>.

DATASET	SSVM BEST	BR LR
SCENE	0.101±.003	0.102±.003
YEAST	0.202±.005	0.199±.005
SYNTH1	0.069±.001	0.067±.002
SYNTH2	0.058±.001	0.084±.001

- There is almost no difference between both algorithms.

<sup>7</sup> Thomas Finley and Thorsten Joachims. Training structural SVMs when exact inference is intractable. In *ICML*. Omnipress, 2008

<sup>8</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An analysis of chaining in multi-label classification. In *ECAI*, 2012

## Probabilistic classifier chains

- Probabilistic classifier chains (PCCs)<sup>9</sup> are an efficient reduction method similar to conditional probability trees.<sup>10</sup>
- They estimate the joint conditional distribution  $P(\mathbf{y} | \mathbf{x})$  as CRFs.
- Their idea is to repeatedly apply the **product rule of probability**:

$$P(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^m P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}).$$

---

<sup>9</sup> J. Read, B. Pfahringer, G. Holmes, and E. Frank. Classifier chains for multi-label classification. *Machine Learning Journal*, 85:333–359, 2011

K. Dembczyński, W. Cheng, and E. Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In *ICML*, pages 279–286. Omnipress, 2010

<sup>10</sup> A. Beygelzimer, J. Langford, Y. Lifshits, G. B. Sorkin, and A. L. Strehl. Conditional probability tree estimation analysis and algorithms. In *UAI*, pages 51–58, 2009

## Probabilistic classifier chains

- Probabilistic classifier chains (PCCs)<sup>9</sup> are an efficient reduction method similar to conditional probability trees.<sup>10</sup>
- They estimate the joint conditional distribution  $P(\mathbf{y} | \mathbf{x})$  as CRFs.
- Their idea is to repeatedly apply the **product rule of probability**:

$$P(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^m P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}).$$

- **Example:**

$$P(y_1, y_2 | \mathbf{x}) = \frac{P(y_1, \mathbf{x})}{P(\mathbf{x})} \frac{P(y_1, y_2, \mathbf{x})}{P(y_1, \mathbf{x})} = P(y_1 | \mathbf{x}) P(y_2 | y_1, \mathbf{x}).$$

---

<sup>9</sup> J. Read, B. Pfahringer, G. Holmes, and E. Frank. Classifier chains for multi-label classification. *Machine Learning Journal*, 85:333–359, 2011

K. Dembczyński, W. Cheng, and E. Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In *ICML*, pages 279–286. Omnipress, 2010

<sup>10</sup> A. Beygelzimer, J. Langford, Y. Lifshits, G. B. Sorkin, and A. L. Strehl. Conditional probability tree estimation analysis and algorithms. In *UAI*, pages 51–58, 2009

## Probabilistic classifier chains

- PCCs follow a reduction to a sequence of subproblems:

$$(\mathbf{x}, \mathbf{y}) \longrightarrow (\mathbf{x}' = (\mathbf{x}, y_1, \dots, y_{i-1}), y = y_i), \quad i = 1, \dots, m$$

- Learning of PCCs relies on constructing probabilistic classifiers for estimating

$$P(y_i | \mathbf{x}, y_1, \dots, y_{i-1}),$$

independently for each  $i = 1, \dots, m$ .

- Let us denote these estimates by

$$Q(y_i | \mathbf{x}, y_1, \dots, y_{i-1}).$$

- The final model is:

$$Q(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^m Q(y_i | \mathbf{x}, y_1, \dots, y_{i-1}).$$

## Probabilistic classifier chains

- We can use scoring functions of the form  $f_i(\mathbf{x}', y_i)$  and train logistic regression (or any probabilistic classifier) to get  $Q(y_i|\mathbf{x}')$ .
- By using the linear models, the overall scoring function takes the form:

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

- Theoretically the order of labels does not matter, but practically it may.

## Probabilistic classifier chains

- PCCs enable estimation of probability of any label vector  $\mathbf{y}$ .
- To get such an estimate it is enough to compute:

$$Q(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^m Q(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$$

- There is, however, a problem how to compute the optimal decision  $\mathbf{h}(\mathbf{x})$  (with respect to  $Q$ ) for a given loss function.



## Probabilistic classifier chains

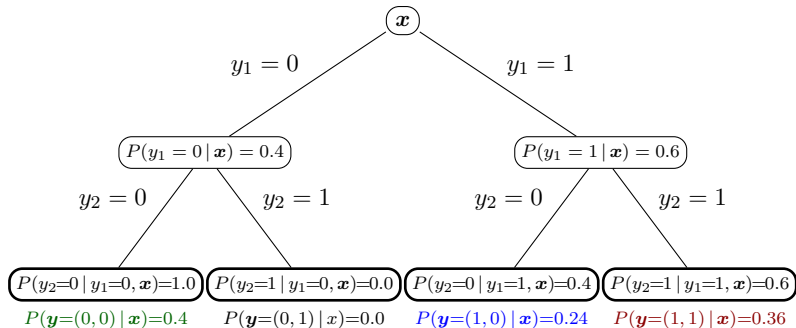
- Inference in PCCs:
  - ▶ Greedy search,
  - ▶ Advanced search techniques: beam search, uniform-cost search,
  - ▶ Exhaustive search,
  - ▶ Sampling + inference.

## Greedy search

- Greedy search follows the chain by using predictions from previous steps as inputs in the consecutive steps:
  - ▶  $f_1 : \mathbf{x} \mapsto \hat{y}_1$
  - ▶  $f_2 : \mathbf{x}, \hat{y}_1 \mapsto \hat{y}_2$
  - ▶  $f_3 : \mathbf{x}, \hat{y}_1, y_2 \mapsto \hat{y}_3$
  - ▶ ...
  - ▶  $f_m : \mathbf{x}, \hat{y}_1, \hat{y}_2, \dots, \hat{y}_{m-1} \mapsto \hat{y}_m$
- Greedy search is fast ( $O(m)$ ).
- Does not require probabilistic classifiers.
- The resulting  $\hat{\mathbf{y}}$  is neither the joint nor the marginal mode.
- Optimal if labels are independent or the probability of the joint mode  $> 0.5$ .

## Greedy search

- Greedy search fails for the joint mode and the marginal mode:



## Advanced search techniques

- Advanced search techniques: beam search,<sup>11</sup> a variant of uniform-cost search.<sup>12</sup>
- Finding the joint mode relies on finding the most probable path in the tree.
- The use of a priority queue and a cut point gives a fast algorithm with provable guarantees.

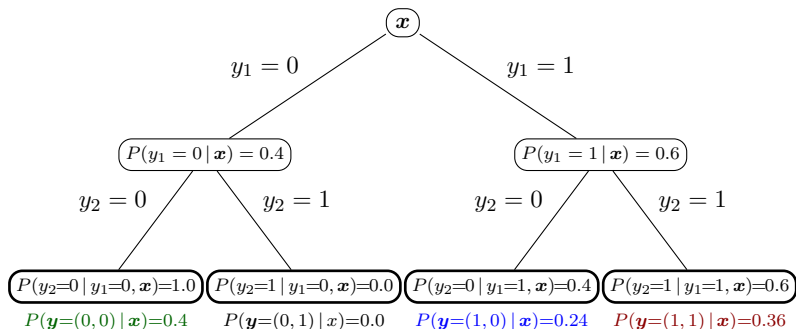
---

<sup>11</sup> A. Kumar, S. Vembu, A.K. Menon, and C. Elkan. Beam search algorithms for multilabel learning. In *Machine Learning*, 2013

<sup>12</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An analysis of chaining in multi-label classification. In *ECAI*, 2012

## Advanced search techniques

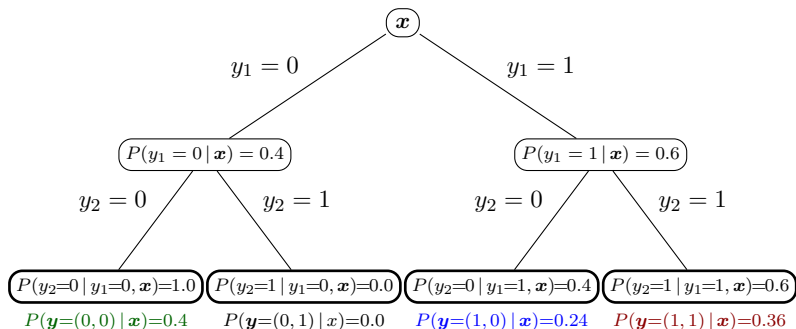
- Uniform-cost search



- Priority list  $Q$ :

## Advanced search techniques

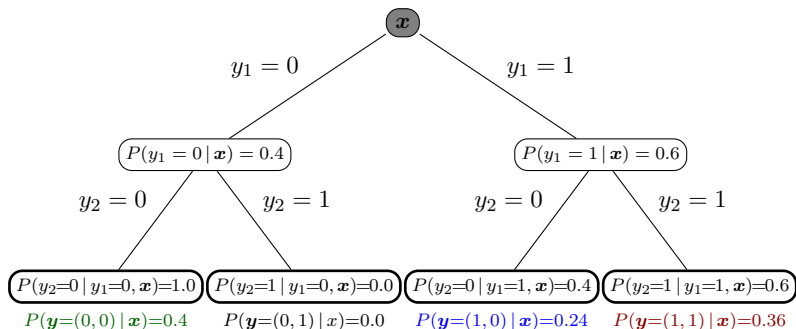
- Uniform-cost search



- Priority list  $Q$ : root

## Advanced search techniques

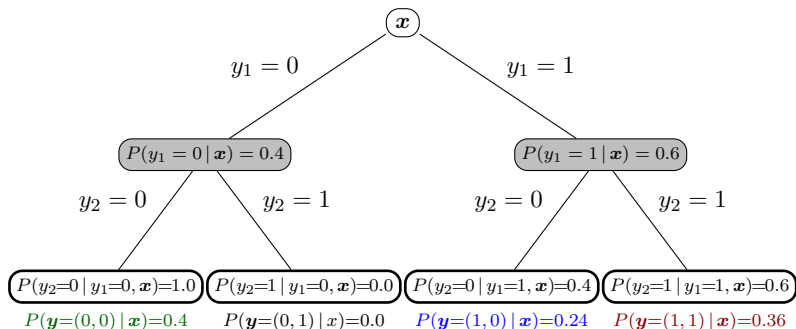
- Uniform-cost search



- Priority list  $Q$ :

## Advanced search techniques

- Uniform-cost search

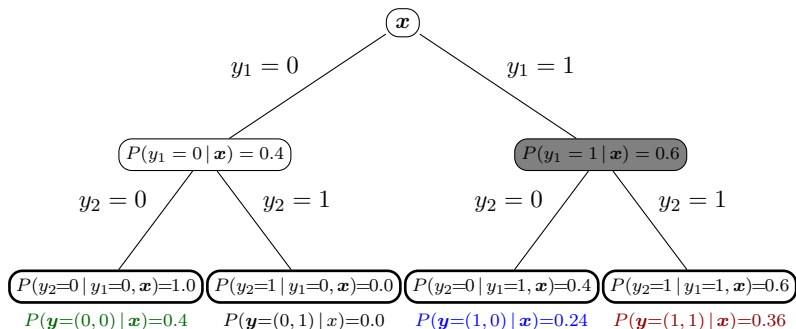


- Priority list  $Q$ :  $[(1), 0.6], [(0), 0.4]$



## Advanced search techniques

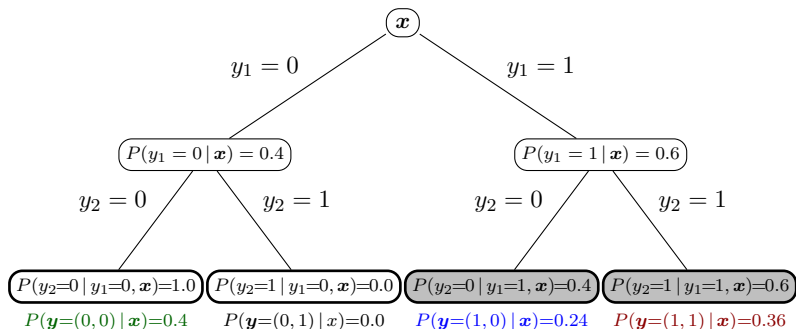
- Uniform-cost search



- Priority list  $Q$ :  $[(0), 0.4]$

## Advanced search techniques

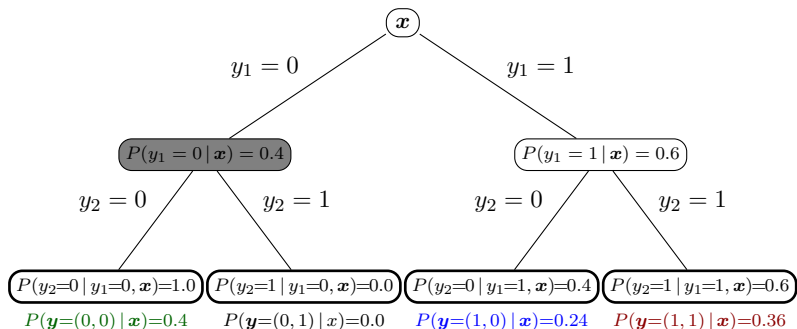
- Uniform-cost search



- Priority list  $\mathcal{Q}$ :  $[(0), 0.4], [(1, 1), 0.36], [(1, 0), 0.24]$

## Advanced search techniques

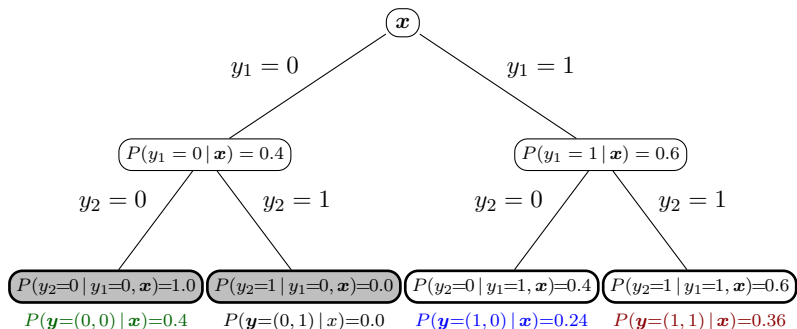
- Uniform-cost search



- Priority list  $\mathcal{Q}$ :  $[(1,1),0.36], [(1,0),0.24]$

## Advanced search techniques

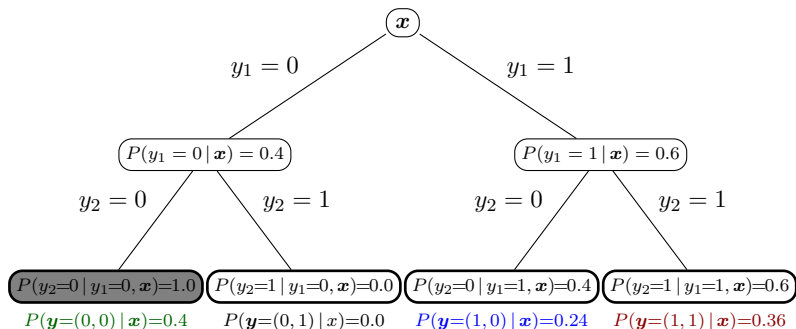
- Uniform-cost search



- Priority list  $\mathcal{Q}$ :  $[(0,0),0.4], [(1,1),0.36], [(1,0),0.24], [(0,1),0.0]$

## Advanced search techniques

- Uniform-cost search



- Priority list  $Q$ : Solution is found

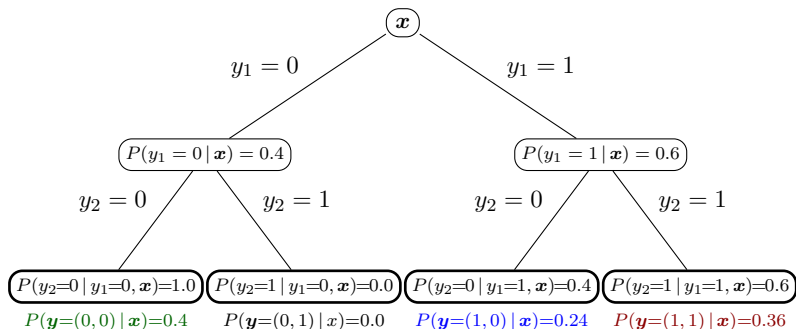
- $\epsilon$ -approximation inference:<sup>13</sup>
  - ▶ Insert items to priority queue  $Q$  with partial probabilities  $> \epsilon$ .
  - ▶ If solution has not been found, then perform greedy search from nodes without survived children.

---

<sup>13</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An analysis of chaining in multi-label classification. In *ECAI*, 2012

## $\epsilon$ -approximation inference

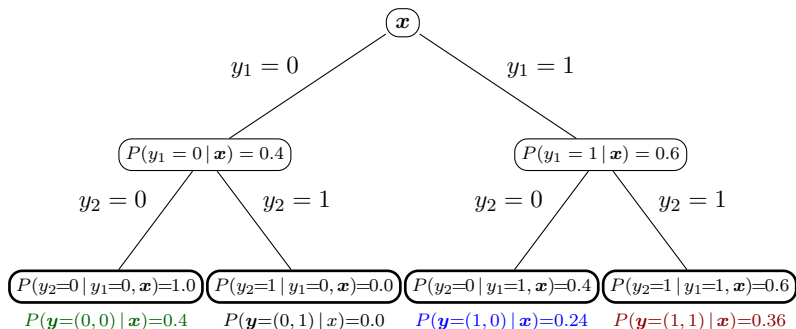
- $\epsilon = 0.5$



- Priority list  $Q$ :

## $\epsilon$ -approximation inference

- $\epsilon = 0.5$

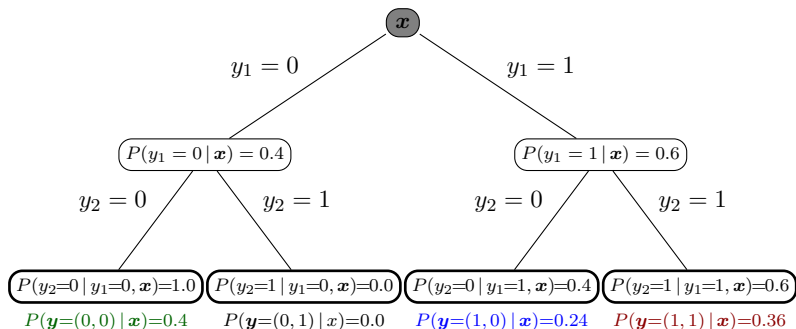


- Priority list  $Q$ : root



## $\epsilon$ -approximation inference

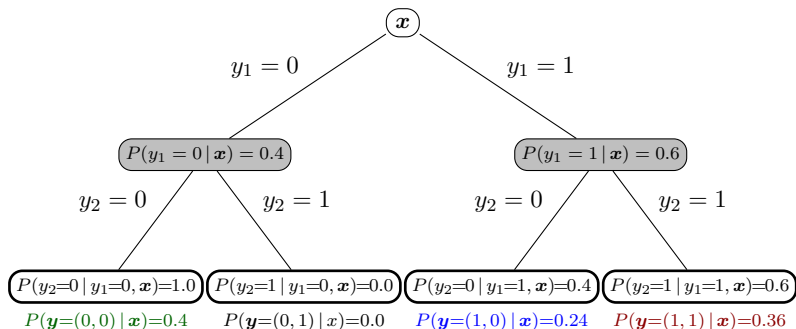
- $\epsilon = 0.5$



- Priority list  $\mathcal{Q}$ :  $\epsilon = 0.5$

## $\epsilon$ -approximation inference

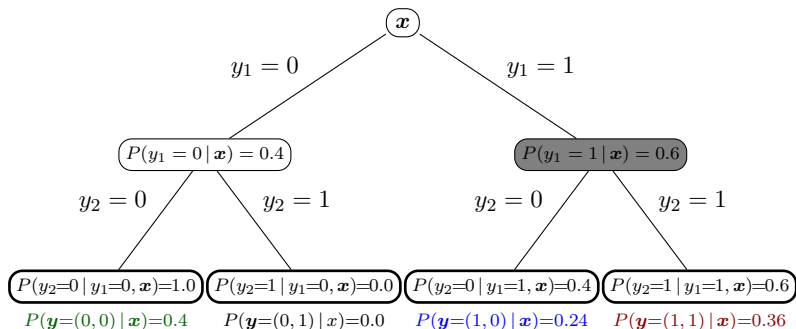
- $\epsilon = 0.5$



- Priority list  $\mathcal{Q}$ :  $[(1),0.6], \epsilon = 0.5, [(0),0.4]$

## $\epsilon$ -approximation inference

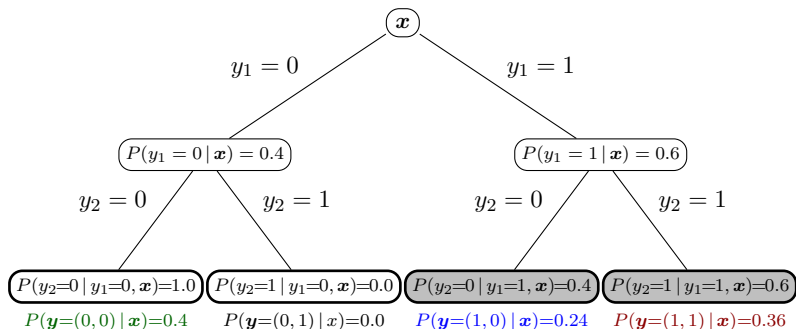
- $\epsilon = 0.5$



- Priority list  $\mathcal{Q}$ :  $\epsilon = 0.5$ ,  $[(0), 0.4]$

## $\epsilon$ -approximation inference

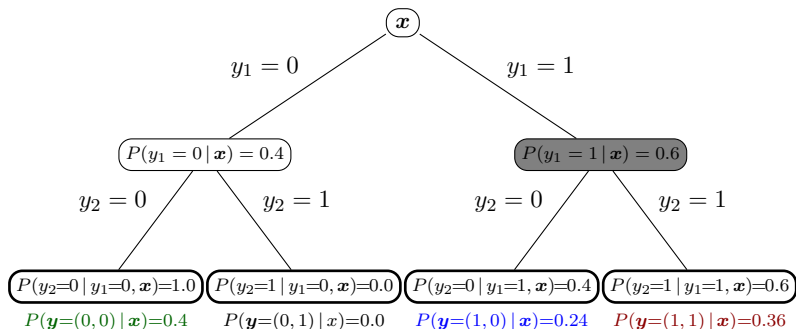
- $\epsilon = 0.5$



- Priority list  $\mathcal{Q}$ :  $\epsilon = 0.5$ ,  $[(0, 0), 0.4]$ ,  $[(1, 1), 0.36]$ ,  $[(1, 0), 0.24]$

## $\epsilon$ -approximation inference

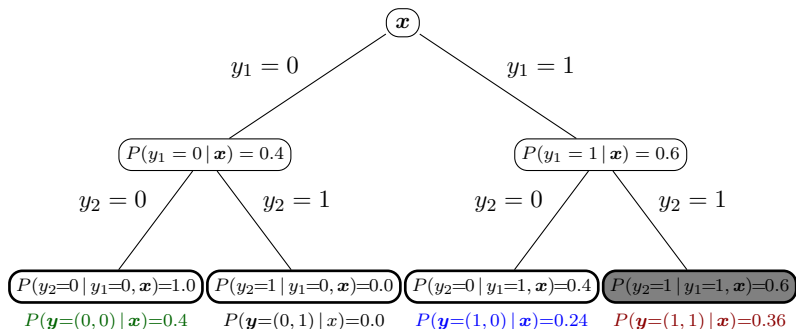
- $\epsilon = 0.5$



- Priority list  $\mathcal{Q}$ : Start the greedy search from (1).

## $\epsilon$ -approximation inference

- $\epsilon = 0.5$



- Priority list  $\mathcal{Q}$ : Suboptimal solution (1,1) is found.

## $\epsilon$ -approximation inference

- For  $\epsilon = 0.5$ , it is equivalent to greedy search.
- For  $\epsilon = 0.0$ , it is equivalent to uniform-cost search.
- For a given  $\epsilon$ , the following guarantees can be given:<sup>14</sup>

**Theorem:** Let  $\epsilon = 2^{-c}$ , where  $1 \leq c \leq m$ . The label vector  $\hat{\mathbf{y}}$  will be returned in time  $\mathcal{O}(m2^c)$  with a guarantee that:

$$Q(\mathbf{y}^* | \mathbf{x}) - Q(\hat{\mathbf{y}} | \mathbf{x}) \leq \epsilon - 2^{-m}$$

---

<sup>14</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An analysis of chaining in multi-label classification. In *ECAI*, 2012

## $\epsilon$ -approximation inference

- The  $\epsilon$ -approximate inference will always find the joint mode if its probability mass  $\geq \epsilon$ .
- In other words, the algorithm finds the solution in a linear time of  $1/p_{\max}$ , where  $p_{\max}$  is the probability mass of the joint mode.
- For problems with low noise (high values of  $p_{\max}$ ), this method should work very fast.
- Greedy search has very bad guarantees:

$$Q(\mathbf{y}^* | \mathbf{x}) - Q(\hat{\mathbf{y}} | \mathbf{x}) \leq 0.5 - 2^{-m}.$$



## Regret bound for PCC

- The typical approach for estimating probabilities of  $\mathbf{y}$  is minimization of the logistic loss:

$$\ell_{\log}(\mathbf{y}, \mathbf{x}, f) = -\log Q(\mathbf{y} | \mathbf{x}),$$

where  $f$  is a model that delivers estimate  $Q(\mathbf{y} | \mathbf{x})$  of  $P(\mathbf{y} | \mathbf{x})$ .

- By using the chain rule of probability, we get:

$$\begin{aligned}\ell_{\log}(\mathbf{y}, \mathbf{x}, f) &= -\log \prod_{i=1}^m Q(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) \\ &= -\sum_{i=1}^m \log Q(y_i | \mathbf{x}, y_1, \dots, y_{i-1}) = -\sum_{i=1}^m \log Q_i(\mathbf{y}),\end{aligned}$$

where we use the notation  $Q_i(\mathbf{y}) = Q(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$ .

- This is a sum of univariate log losses on a path from the root to the leaf corresponding to  $\mathbf{y}$ .

## Regret bound for PCC

- **Theorem:** For all distributions and all PCCs trained with logistic regression  $f$  and used with the  $\epsilon$ -approximate inference algorithm,

$$\text{Reg}_{0/1}(\text{PCC}_\epsilon(f)) \leq \sqrt{2m \overline{\text{Reg}}_{\log}(f)} + \epsilon$$

where  $\overline{\text{Reg}}_{\log}(f)$  is the average logistic regret over the paths from the root to the leafs.

## PCC for other losses

- **Exhaustive search:**

- ▶ Compute the entire distribution  $Q(\mathbf{y} | \mathbf{x})$  by traversing the probability tree.
- ▶ Use an appropriate inference for a given loss  $\ell$  on the estimated joint distribution:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{h} \in \mathcal{Y}} \sum_{\mathbf{y} \in \mathcal{Y}} Q(\mathbf{y} | \mathbf{x}) \ell(\mathbf{y}, \mathbf{h}(\mathbf{x}))$$

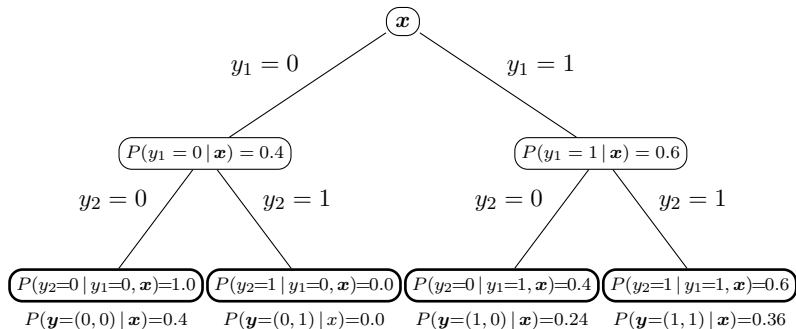
- ▶ This approach is extremely costly.

- **Ancestral sampling:**

- ▶ Sampling can be easily performed by using the probability tree.
- ▶ Make inference based on the empirical distribution.
- ▶ Hamming loss: estimate marginal probabilities.

## Probabilistic classifier chains

- Exhaustive search and ancestral sampling:



- Sample:  $(1,1), (1,0), (0,0), (0,0), (1,1), (0,0), (1,0), (1,1), (0,0) \dots$

## Probabilistic classifier chains

Table : PCC vs. SSVMs on Hamming loss and PCC vs. BR on subset 0/1 loss.

DATASET	PCC	SSVM BEST	PCC	BR
	HAMMING LOSS		SUBSET 0/1 LOSS	
SCENE	0.104±.004	0.101±.003	0.385±.014	0.509±.014
YEAST	0.203±.005	0.202±.005	0.761±.014	0.842±.012
SYNTH1	0.067±.001	0.069±.001	0.239±.006	0.240±.006
SYNTH2	0.000±.000	0.058±.001	0.000±.000	0.832±.004
REUTERS	0.060±.002	0.045±.001	0.598±.009	0.689±.008
MEDIAMILL	0.172±.001	0.182±.001	0.885±.003	0.902±.003

## Recurrent classifiers

- PCCs are similar to Maximum Entropy Markov Models (MEMMs)<sup>15</sup> introduced for sequence learning:
  - ▶ One logistic classifier that takes dependences up to the  $k$ -th degree.
  - ▶ Inference by dynamic programming.
- Searn<sup>16</sup> is another approach that is based on recurrent classifiers:
  - ▶ Linear inference.
  - ▶ Learning is performed in the iterative way to solve the egg and the chicken problem: output of the classifier is also used as input to the classifier.

---

<sup>15</sup> A. K. McCallum, D. Freitag, and F. (2000) Pereira. Maximum entropy markov models for information extraction and segmentation. In *ICML*, 2000

<sup>16</sup> H. Daumé III, J. Langford, and D. Marcu. Search-based structured prediction. *Machine Learning*, 75:297–325, 2009

## Output search space

- More advanced search techniques.
- Popular topic in structured output prediction.
- Search techniques for different task losses.<sup>17</sup>

---

<sup>17</sup> J.R. Dopper, A. Fern, and P. Tadepalli. Structured prediction via output space search. *JMLR*, 15:1317–1350, 2014

## PCC for multi-class classification

- PCC can be used for multi-class classification:
  - ▶ Map each class label to a label vector: binary coding, hierarchical clustering, ...
  - ▶ The same idea as in conditional probability trees (CPT).<sup>18</sup>
  - ▶ Label tree classifiers for efficient multi-class classification.<sup>19</sup>

---

<sup>18</sup> A. Beygelzimer, J. Langford, Y. Lifshits, G. B. Sorkin, and A. L. Strehl. Conditional probability tree estimation analysis and algorithms. In *UAI*, pages 51–58, 2009

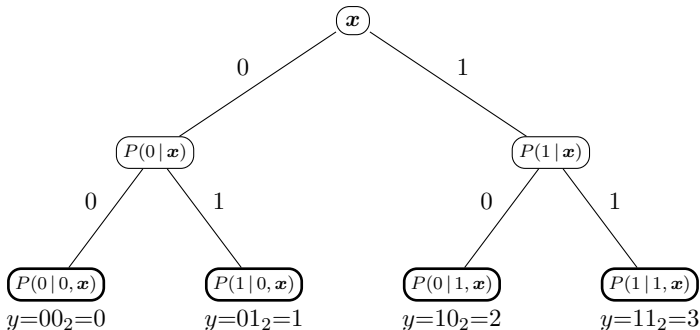
<sup>19</sup> S. Bengio, J. Weston, and D. Grangier. Label embedding trees for large multi-class tasks. In *NIPS*, pages 163–171. Curran Associates, Inc., 2010

J. Deng, S. Satheesh, A. C. Berg, and Fei Fei F. Li. Fast and balanced: Efficient label tree learning for large scale object recognition. In *NIPS*, pages 567–575. 2011



## PCC for multi-class classification

- We assign each class an integer from 0 to  $k - 1$  and code it by its binary representation on  $m$  bits.
- Example:  $k = 4$ ,  $\mathcal{Y} = \{0, 1, 2, 3\}$ .
- $k$  leaves, one for each class.



## Consistent and efficient label tree classifiers

- PCC: fast learning but inference can be costly.
- Greedy search is the most efficient, but is not consistent.
- How to ensure a linear inference in  $m$  for any loss?

## Filter trees

- Filter trees (FT)<sup>20</sup> have been originally introduced for cost-sensitive multi-class classification, but can be easily adapted to multi-label classification.

---

<sup>20</sup> A. Beygelzimer, J. Langford, and P. D. Ravikumar. Error-correcting tournaments. In *ALT*, pages 247–262, 2009

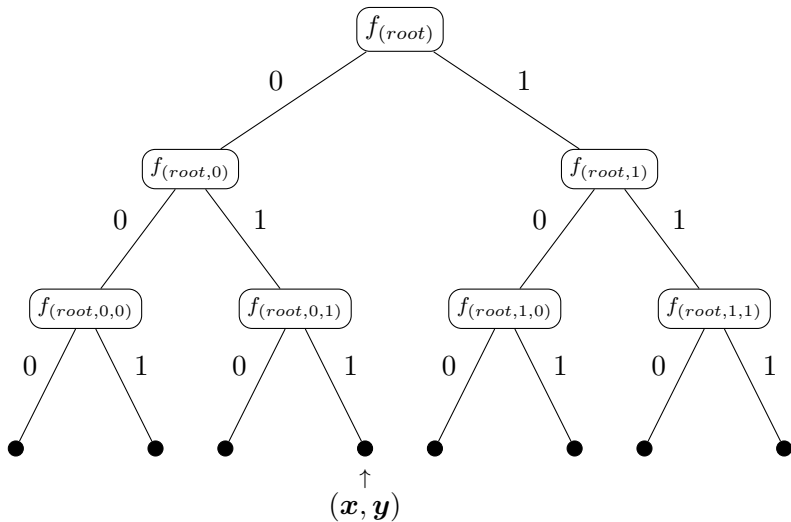
## Filter trees

- Filter trees (FT)<sup>20</sup> have been originally introduced for cost-sensitive multi-class classification, but can be easily adapted to multi-label classification.
- They use a **bottom-up** learning algorithm to train the label tree.
- Based on a single **elimination tournament** on the set of classes/label combinations.

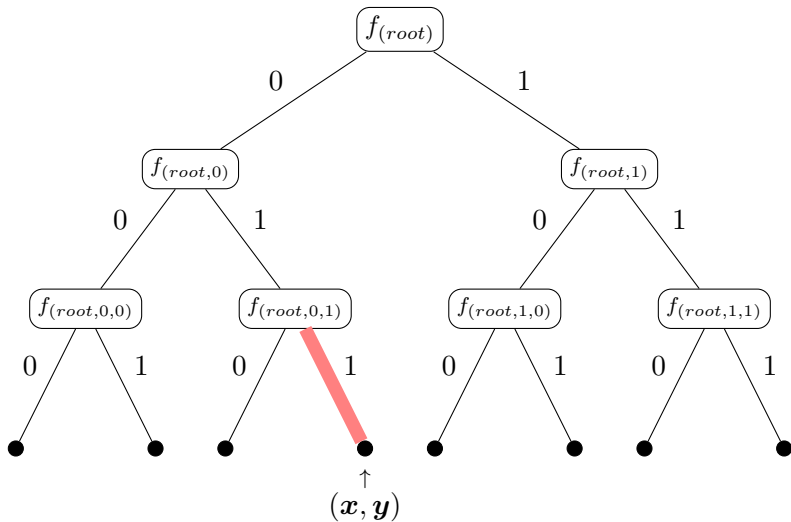
---

<sup>20</sup> A. Beygelzimer, J. Langford, and P. D. Ravikumar. Error-correcting tournaments. In *ALT*, pages 247–262, 2009

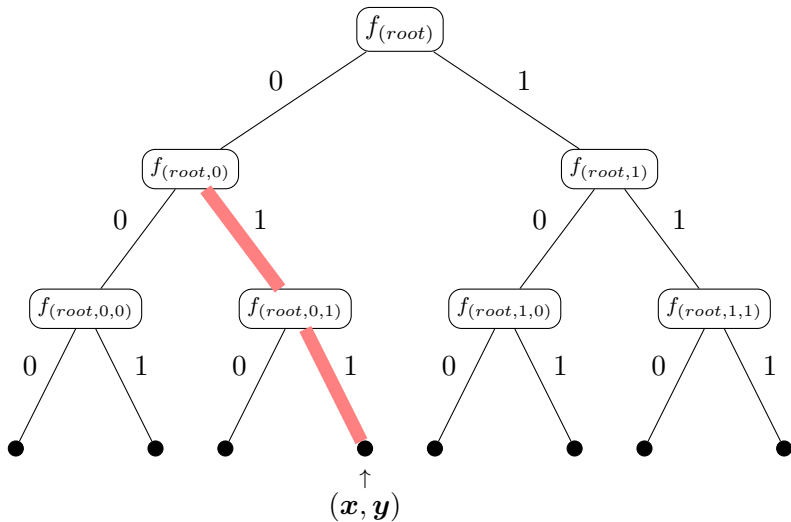
## Filter trees : Example



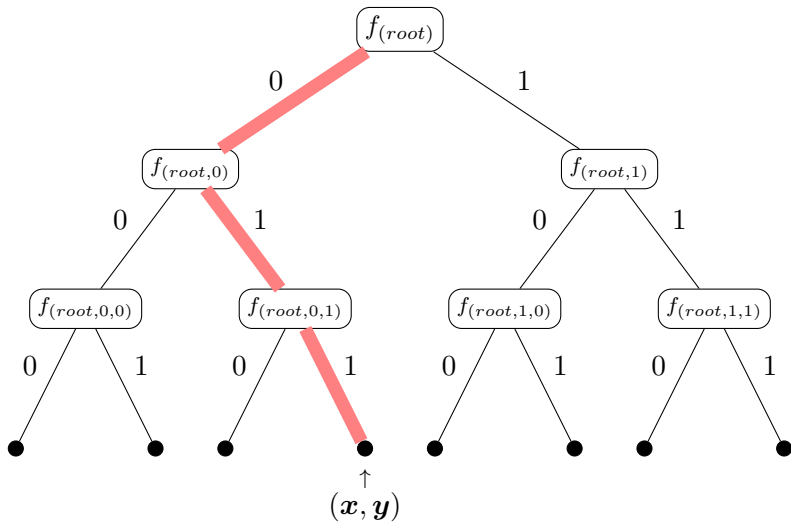
## Filter trees : Example



## Filter trees : Example



## Filter trees : Example





## Filter trees

- FT are trained to predict  $y_{i+1}$  based on previous labels.
- FT implicitly transforms the underlying distribution  $P$  over multi-class/multi-label examples into a **specific distribution**  $P^{\text{FT}}$  over weighted binary examples.
- The inference procedure of FT is straight-forward and uses the **greedy search**.
- FT are **consistent** for any cost function.

## Filter trees

- Filter tree training:

- 1: **Input:** training set  $\{(x_i, y_i)\}_{i=1}^n$ , importance-weighted binary learner *Learn*
- 2: **for** each non-leaf node  $v = (\text{root}, y_1, \dots, y_{i-1})$  in the order from leaves to root **do**
- 3:    $S_v = \emptyset$
- 4:   **for** each training example  $(x, y)$  **do**
- 5:     Let  $y_l$  and  $y_r$  be the two label vectors input to  $v$
- 6:      $y_i \leftarrow \arg \min_{l,r} \{\ell(y, y_l), \ell(y, y_r)\}$
- 7:      $w = |\ell(y, y_l) - \ell(y, y_r)|$
- 8:      $S_v \leftarrow S_v \cup (x, y_i, w)$
- 9:   **end for**
- 10:    $f_v = \text{Learn}(S_v)$
- 11: **end for**
- 12: **return**  $f = \{f_v\}$

## Filter trees

- Different training schemes possible:
  - ▶ Train a classifier in each node,
  - ▶ Train a classifier on each level,
  - ▶ Train one global binary classifier (in several loops).
- The tree in multi-label classification is given naturally, but the order of labels may influence the performance.
- In general case, training can be costly ( $O(2^m)$ ), but efficient variants for multi-label classification exist.<sup>21</sup>
- Prediction is always linear in the number of labels ( $O(m)$ ).

---

<sup>21</sup> Chun-Liang Li and Hsuan-Tien Lin. Condensed filter tree for cost-sensitive multi-label classification. In *ICML*, pages 423–431, 2014

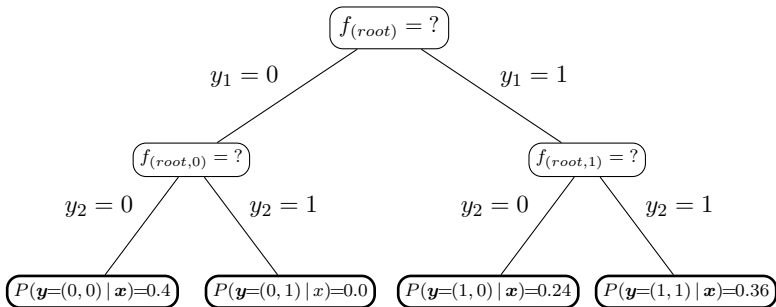
## Filter trees

- Filter trees for the subset 0/1 loss **filter out** all examples that are misclassified by the lower-level classifiers.
- Therefore, training in this case is also linear in the number of labels ( $O(m)$ ).
- $f_{(root, y_1, \dots, y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent labels correctly:

$$f_{(root, y_1, \dots, y_i)} : \mathbf{x} \mapsto (y_{i+1} \mid y_{j+1} = f_{(root, y_1, \dots, y_j)} : j = i + 1, \dots, m - 1)$$

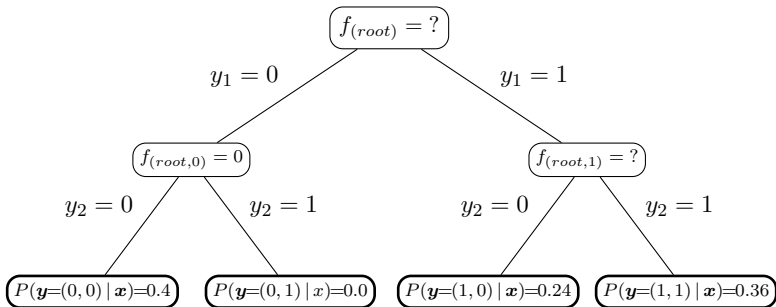
## Filter trees: Consistency

- Consistency of FT for a single  $x$ :



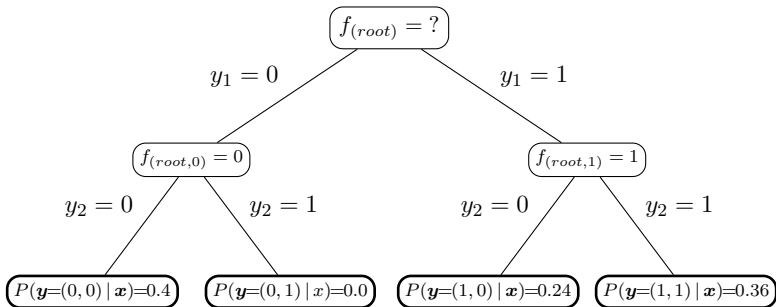
## Filter trees: Consistency

- Consistency of FT for a single  $x$ :



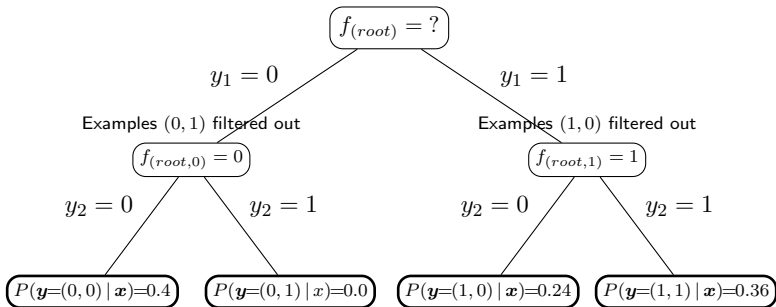
## Filter trees: Consistency

- Consistency of FT for a single  $x$ :



## Filter trees: Consistency

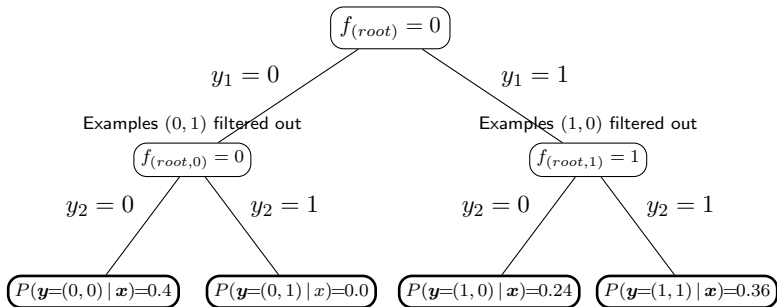
- Consistency of FT for a single  $\mathbf{x}$ :





## Filter trees: Consistency

- Consistency of FT for a single  $\mathbf{x}$ :



## Regret bound for filter trees

- Let  $f_v$  be a classifier for the binary classification problem induced at node  $v$ .
- The average binary regret is defined as:

$$\overline{\text{Reg}}_{0/1}(f, P^{\text{FT}}) = \frac{1}{\sum_v W_v} \sum_v \text{Reg}_{0/1}(f_v, P_v^{\text{FT}}) W_v,$$

where

$$W_v = \mathbb{E}_{(\mathbf{x}, \mathbf{y})} w_v(\mathbf{x}, \mathbf{y}).$$

- **Theorem:**<sup>22</sup> For all distributions and all FT classifiers trained with a binary classifier  $f$ , and any cost-matrix-based task loss  $\ell$ ,

$$\text{Reg}_\ell(\text{FT}(f)) \leq \overline{\text{Reg}}_{0/1}(f, P^{\text{FT}}) \sum_v W_v.$$

---

<sup>22</sup> A. Beygelzimer, J. Langford, and P. D. Ravikumar. Error-correcting tournaments. In *ALT*, pages 247–262, 2009

## Regret bound for filter trees

- For subset 0/1 loss, we have

$$\sum_v w_v(\mathbf{x}, \mathbf{y}) \leq m,$$

since each training example  $(\mathbf{x}, \mathbf{y})$  will appear in training at most once per level with importance weight 1.

- The regret bound has then the form:

$$\text{Reg}_\ell(\text{FT}(f)) \leq m \overline{\text{Reg}}_{0/1}(f, P^{\text{FT}}).$$

# Outline

- 1 Multi-label classification
- 2 Simple approaches to multi-label classification
- 3 Beyond simple approaches
- 4 Other task losses**
- 5 Rank loss minimization
- 6 Summary

# Maximization of the F-measure

- Applications: Information retrieval, document tagging, and NLP.
- JRS 2012 Data Mining Competition: Indexing documents from MEDLINE or PubMed Central databases with concepts from the Medical Subject Headings ontology.

The screenshot shows the NLM website header with the logo and navigation tabs: Databases, Find, Read, Learn, Explore NLM, Research at NLM, and NLM for You. The main heading is "MEDLINE®/PubMed® Resources Guide". Below this, it states: "MEDLINE® contains journal citations and abstracts for biomedical literature from around the world. Pub possible." It then says: "The following resources provide detailed information about MEDLINE data and searching PubMed. If you Service." A horizontal menu lists: News | Overviews | Journals | Data Structure & Content | Data Policies | Searching PubMed | Toc. The "NEWS" section is expanded, showing four items: "PubMed New and Noteworthy" (List of changes to PubMed by date, with links to the Technical Bulletin), "NLM Technical Bulletin" (The NLM Technical Bulletin is your main source for detailed information about changes and updates), "NLM-Announces" (NLM e-mail list for announcing important information and changes to NLM systems including PubMed), and "PubMed-Alerts" (An announcements-only e-mail list that notifies subscribers of major system problems with PubMed 8:30am to 5:00pm ET). The "OVERVIEWS" section is also visible at the bottom.

## Maximization of the F-measure

- 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.

## SSVMs for $F_\beta$ -based loss

- SSVMs can be used to minimize  $F_\beta$ -based loss.
- Rescale the margin by  $\ell_F(\mathbf{y}, \mathbf{y}')$ .
- Two algorithms:<sup>23</sup>

### RML

No label interactions:

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

Quadratic learning and linear prediction

### SML

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)

- Both are inconsistent.

---

<sup>23</sup> J. Petterson and T. S. Caetano. Reverse multi-label learning. In *NIPS*, pages 1912–1920, 2010  
J. Petterson and T. S. Caetano. Submodular multi-label learning. In *NIPS*, pages 1512–1520, 2011

## Plug-in rule approach

- Plug estimates of required parameters into the Bayes classifier:<sup>24</sup>

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

- **No closed form** solution for this optimization problem.
- The problem **cannot** be solved **naively** by brute-force search:
  - ▶ This would require to check all possible combinations of labels ( $2^m$ )
  - ▶ To sum over  $2^m$  number of elements for computing the expected value.
  - ▶ The number of parameters to be estimated ( $P(\mathbf{y})$ ) is  $2^m$ .

---

<sup>24</sup> W. Waegeman, K. Dembczynski, W. Cheng A. Jachnik, and E. Hüllermeier. On the Bayes-optimality of F-measure maximizers. *Minor revision*, 2014



## Plug-in rule approach

- Approximation needed?

---

<sup>25</sup> N. Ye, K. Chai, W. Lee, and H. Chieu. Optimizing F-measures: a tale of two approaches. In *ICML*, 2012

<sup>26</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An exact algorithm for F-measure maximization. In *NIPS*, volume 25, 2011

<sup>27</sup> K. Dembczynski, A. Jachnik, W. Kotlowski, W. Waegeman, and E. Hüllermeier. Optimizing the F-measure in multi-label classification: Plug-in rule approach versus structured loss minimization. In *ICML*, 2013

## Plug-in rule approach

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

---

<sup>25</sup> N. Ye, K. Chai, W. Lee, and H. Chieu. Optimizing F-measures: a tale of two approaches. In *ICML*, 2012

<sup>26</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An exact algorithm for F-measure maximization. In *NIPS*, volume 25, 2011

<sup>27</sup> K. Dembczynski, A. Jachnik, W. Kotlowski, W. Waegeman, and E. Hüllermeier. Optimizing the F-measure in multi-label classification: Plug-in rule approach versus structured loss minimization. In *ICML*, 2013

## Plug-in rule approach

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

### LFP:

Assumes label independence.

Linear number of parameters:  
 $P(y_i = 1)$ .

Inference based on dynamic programming.<sup>25</sup>

Reduction to LR for each label.

### EFP:

No assumptions.

Quadratic number of parameters:  
 $P(y_i = 1, s = \sum_i y_i)$ .

Inference based on matrix multiplication and top  $k$  selection.<sup>26</sup>

Reduction to multinomial LR for each label.

- EFP is consistent.<sup>27</sup>

---

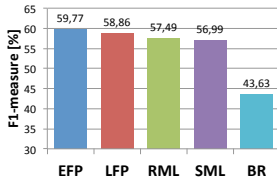
<sup>25</sup> N. Ye, K. Chai, W. Lee, and H. Chieu. Optimizing F-measures: a tale of two approaches. In *ICML*, 2012

<sup>26</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An exact algorithm for F-measure maximization. In *NIPS*, volume 25, 2011

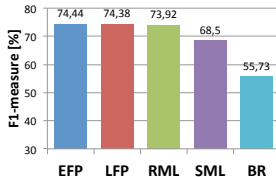
<sup>27</sup> K. Dembczynski, A. Jachnik, W. Kotlowski, W. Waegeman, and E. Hüllermeier. Optimizing the F-measure in multi-label classification: Plug-in rule approach versus structured loss minimization. In *ICML*, 2013

# Maximization of the F-measure

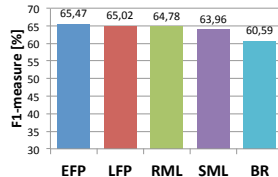
## IMAGE



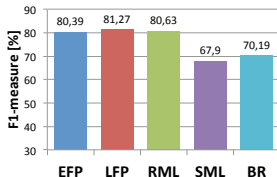
## SCENE



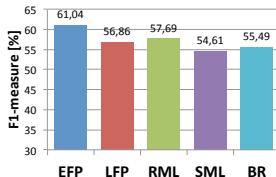
## YEAST



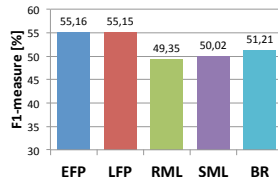
## MEDICAL



## ENRON



## MEDIAMILL



# Outline

- 1 Multi-label classification
- 2 Simple approaches to multi-label classification
- 3 Beyond simple approaches
- 4 Other task losses
- 5 Rank loss minimization**
- 6 Summary

# Multi-label ranking

## Serena romps to fifth Wimbledon title against brave Radwanska

By Paul Gittings, CNN

July 7, 2012 -- Updated 2220 GMT (0620 HKT)



### STORY HIGHLIGHTS

- Serena Williams wins fifth Wimbledon crown
- American beats Agnieszka Radwanska of Poland 6-1 5-7 6-2
- Radwanska battles respiratory

(CNN) -- Serena Williams fended off a stirring fightback from Agnieszka Radwanska to win her fifth Wimbledon singles title with a 6-1 5-7 6-2 victory Saturday.

It was the 30-year-old American's 14th grand slam crown and her first since winning at the All England Club in 2010, but Poland's Radwanska made her fight every inch of the way.

## Multi-label classification

politics	0
economy	0
business	0
sport	1
tennis	1
soccer	0
show-business	0
celebrities	1
:	:

England	1
USA	1
Poland	1
Lithuania	0

# Multi-label ranking

## Multi-label ranking

### Serena romps to fifth Wimbledon title against brave Radwanska

By Paul Gittings, CNN

July 7, 2012 -- Updated 2220 GMT (0620 HKT)



#### Women's singles Wimbledon Championship

##### STORY HIGHLIGHTS

- Serena Williams wins fifth Wimbledon crown
- American beats Agnieszka Radwanska of Poland 6-1 5-7 6-2
- Radwanska battles respiratory

(CNN) -- Serena Williams fended off a stirring fightback from Agnieszka Radwanska to win her fifth Wimbledon singles title with a 6-1 5-7 6-2 victory Saturday.

It was the 30-year-old American's 14th grand slam crown and her first since winning at the All England Club in 2010, but Poland's Radwanska made her fight every inch of the way.

tennis



sport



England



Poland



USA



politics

## Multi-label ranking

- Ranking loss:

$$\ell_{\text{rnk}}(\mathbf{y}, \mathbf{f}) = w(\mathbf{y}) \sum_{(i,j): y_i > y_j} \left( \mathbb{I}[f_i(\mathbf{x}) < f_j(\mathbf{x})] + \frac{1}{2} \mathbb{I}[f_i(\mathbf{x}) = f_j(\mathbf{x})] \right),$$

where  $w(\mathbf{y}) < w_{\max}$  is a weight function.

	$X_1$	$X_2$	$Y_1$	$Y_2$	$\dots$	$Y_m$			
$x$	4.0	2.5	1	0		0			
			$h_2$	$>$	$h_1$	$>$	$\dots$	$>$	$h_m$



## Multi-label ranking

- **Ranking loss:**

$$\ell_{\text{rnk}}(\mathbf{y}, \mathbf{f}) = w(\mathbf{y}) \sum_{(i,j): y_i > y_j} \left( \mathbb{I}[f_i(\mathbf{x}) < f_j(\mathbf{x})] + \frac{1}{2} \mathbb{I}[f_i(\mathbf{x}) = f_j(\mathbf{x})] \right),$$

where  $w(\mathbf{y}) < w_{\max}$  is a weight function.

The weight function  $w(\mathbf{y})$  is usually used to normalize the range of rank loss to  $[0, 1]$ :

$$w(\mathbf{y}) = \frac{1}{n_+ n_-},$$

i.e., it is equal to the inverse of the total number of pairwise comparisons between labels.

## Pairwise surrogate losses

- The most intuitive approach is to use pairwise **convex surrogate** losses of the form

$$\tilde{\ell}_{\phi}(\mathbf{y}, \mathbf{f}) = \sum_{(i,j): y_i > y_j} w(\mathbf{y}) \phi(f_i - f_j),$$

where  $\phi$  is

- ▶ an exponential function (BoosTexter)<sup>28</sup>:  $\phi(f) = e^{-f}$ ,
- ▶ logistic function (LLLR)<sup>29</sup>:  $\phi(f) = \log(1 + e^{-f})$ ,
- ▶ or hinge function (RankSVM)<sup>30</sup>:  $\phi(f) = \max(0, 1 - f)$ .

---

<sup>28</sup> R. E. Schapire and Y. Singer. BoosTexter: A Boosting-based System for Text Categorization. *Machine Learning*, 39(2/3):135–168, 2000

<sup>29</sup> O. Dekel, Ch. Manning, and Y. Singer. Log-linear models for label ranking. In *NIPS*. MIT Press, 2004

<sup>30</sup> A. Elisseeff and J. Weston. A kernel method for multi-labelled classification. In *NIPS*, pages 681–687, 2001

## Multi-label ranking

- This approach is, however, **inconsistent** for the most commonly used convex surrogates.<sup>31</sup>
- The **consistent** classifier can be, however, obtained by using univariate loss functions<sup>32</sup> . . .

---

<sup>31</sup> J. Duchi, L. Mackey, and M. Jordan. On the consistency of ranking algorithms. In *ICML*, pages 327–334, 2010

W. Gao and Z.-H. Zhou. On the consistency of multi-label learning. *Artificial Intelligence*, 199-200:22–44, 2013

<sup>32</sup> K. Dembczynski, W. Kotłowski, and E. Hüllermeier. Consistent multilabel ranking through univariate losses. In *ICML*, 2012

## Reduction to weighted binary relevance

- The Bayes ranker can be obtained by sorting labels according to:

$$\Delta_i^1 = \sum_{\mathbf{y}: y_i=1} w(\mathbf{y})P(\mathbf{y} | \mathbf{x}).$$

- For  $w(\mathbf{y}) \equiv 1$ ,  $\Delta_i^u$  reduces to **marginal probabilities**  $P(y_i = u | \mathbf{x})$ .
- The solution can be obtained with BR or its weighted variant in a general case.

## Reduction to weighted binary relevance

- Consider the sum of **univariate (weighted)** losses:

$$\tilde{\ell}_{\text{exp}}(\mathbf{y}, \mathbf{f}) = w(\mathbf{y}) \sum_{i=1}^m e^{-y' f_i},$$

$$\tilde{\ell}_{\log}(\mathbf{y}, \mathbf{f}) = w(\mathbf{y}) \sum_{i=1}^m \log \left( 1 + e^{-y' f_i} \right).$$

where  $y' = 2y_i - 1$ .

- The risk minimizer of these losses is:

$$f_i^*(\mathbf{x}) = \frac{1}{c} \log \frac{\Delta_i^1}{\Delta_i^0} = \frac{1}{c} \log \frac{\Delta_i^1}{W - \Delta_i^1},$$

which is a strictly increasing transformation of  $\Delta_i^1$ , where

$$W = \mathbb{E}_{\mathbf{y}}[w(\mathbf{y}) \mid \mathbf{x}] = \sum_{\mathbf{y}} w(\mathbf{y}) P(\mathbf{y} \mid \mathbf{x}).$$

## Reduction to weighted binary relevance

- **Vertical reduction**: Solving  $m$  independent classification problems.
- Standard algorithms, like AdaBoost and logistic regression, can be adapted to this setting.
- AdaBoost.MH follows this approach for  $w = 1$ .<sup>33</sup>
- Besides its **simplicity** and **efficiency**, this approach is **consistent** (regret bounds have also been derived).<sup>34</sup>

---

<sup>33</sup> R. E. Schapire and Y. Singer. BoosTexter: A Boosting-based System for Text Categorization. *Machine Learning*, 39(2/3):135–168, 2000

<sup>34</sup> K. Dembczynski, W. Kotlowski, and E. Hüllermeier. Consistent multilabel ranking through univariate losses. In *ICML*, 2012

## Weighted binary relevance

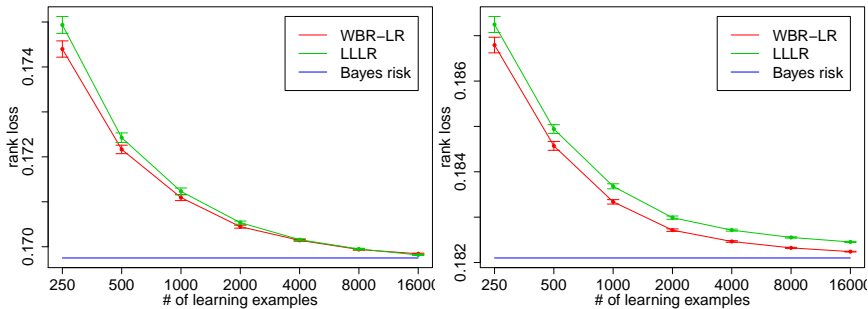


Figure : WBR LR vs. LLLR. Left: independent data. Right: dependent data.

- **Label independence:** the methods perform more or less en par.
- **Label dependence:** WBR shows small but consistent improvements.

## Benchmark data

Table : WBR-AdaBoost vs. AdaBoost.MR (left) and WBR-LR vs LLLR (right).

DATASET	AB.MR	WBR-AB	LLLRLR	WBR-LR
IMAGE	0.2081	0.2041	0.2047	0.2065
EMOTIONS	0.1703	0.1699	0.1743	0.1657
SCENE	0.0720	0.0792	0.0861	0.0793
YEAST	0.2072	0.1820	0.1728	0.1736
MEDIAMILL	0.0665	0.0609	0.0614	0.0472

- WBR is at least competitive to state-of-the-art algorithms defined on pairwise surrogates.



# Outline

- 1 Multi-label classification
- 2 Simple approaches to multi-label classification
- 3 Beyond simple approaches
- 4 Other task losses
- 5 Rank loss minimization
- 6 Summary**

## Summary

- Multi-label classification.
- Simple approaches to multi-label classification.
- Task losses minimized by BR and LP.
- CRFs and SSVMs.
- PCC and Filter trees.
- Approaches for other loss functions: F-measure and rank loss.

## Open challenges

- Learning and inference algorithms for any task loss and output structure.
- Consistency of the algorithms.
- Large-scale datasets: number of instances, features, and labels.

## Conclusions

- Take-away message:
  - ▶ Two main issues: loss minimization and label dependence.
  - ▶ Two main approaches: surrogate loss minimization and reduction.
  - ▶ Consistency of algorithms.
  - ▶ High regret between solutions for different losses.
  - ▶ Proper modeling of label dependence for different loss functions.
  - ▶ Be careful with empirical evaluations.
  - ▶ Independent models can perform quite well.
- For more check:

<http://www.cs.put.poznan.pl/kdembczynski>