

# Efficient Label Tree Structures for Top-k Classification

Arkadiusz Jachnik and Krzysztof Dembczyński

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



IFORS 2014, Barcelona, July 14, 2014

# Introduction

# Introduction

- Large-scale machine learning problems:

# Introduction

- Large-scale machine learning problems:
  - ▶ Large number of training examples (in hundreds of millions),

# Introduction

- Large-scale machine learning problems:
  - ▶ Large number of training examples (in hundreds of millions),
  - ▶ Large number of features (in hundreds of thousands),

# Introduction

- Large-scale machine learning problems:
  - ▶ Large number of training examples (in hundreds of millions),
  - ▶ Large number of features (in hundreds of thousands),
  - ▶ **Large number of class labels (in tens or hundreds of thousands).**

# Introduction

- Large-scale machine learning problems:
  - ▶ Large number of training examples (in hundreds of millions),
  - ▶ Large number of features (in hundreds of thousands),
  - ▶ **Large number of class labels (in tens or hundreds of thousands).**
- A need for efficient algorithms and data structures.

# Introduction

- Large-scale machine learning problems:
  - ▶ Large number of training examples (in hundreds of millions),
  - ▶ Large number of features (in hundreds of thousands),
  - ▶ **Large number of class labels (in tens or hundreds of thousands).**
- A need for efficient algorithms and data structures.
- **Label trees** for problems with a large number of class labels.



# Introduction

- Large-scale machine learning problems:
  - ▶ Large number of training examples (in hundreds of millions),
  - ▶ Large number of features (in hundreds of thousands),
  - ▶ **Large number of class labels (in tens or hundreds of thousands).**
- A need for efficient algorithms and data structures.
- **Label trees** for problems with a large number of class labels.
- Two settings: **multi-class** and multi-label classification.



Image annotation: cloud? sky? tree?

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

411

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

## Supervised learning

- **Example**  $x$  is coming from an unknown input distribution  $P(\mathbf{x})$ .
- **True outcome**  $y \in \mathcal{Y} = \{1, \dots, m\}$  is generated from  $P(y | \mathbf{x})$ .
- **Predicted outcome** is given by  $\hat{y} = h(\mathbf{x})$ .
- The **(task) loss** of a single prediction is  $\ell(y, \hat{y})$ .

## Supervised learning

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

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

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

$$h^* = \arg \min_h L(h)$$

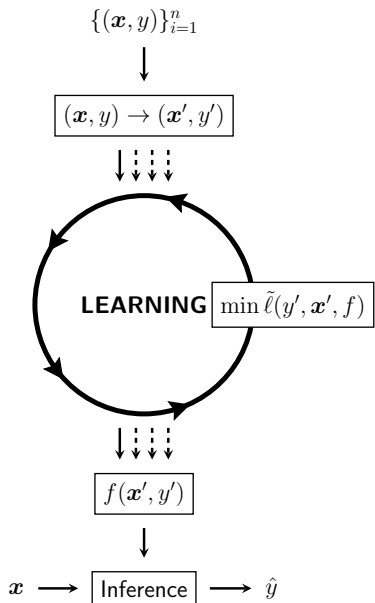
- The **regret** of a classifier  $h$  with respect to  $\ell$  is defined as:

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

## Supervised learning

- We use **training examples**  $\{\mathbf{x}_i, y_i\}_1^n$  to find either:
  - ▶ A good approximation of  $h^*$ .
  - ▶ A scoring function  $f(\mathbf{x}, y)$  being a good estimate of  $P(y | \mathbf{x})$  (or a function of it).
- In the second case, we need to apply an **inference procedure** to approximate  $h^*$ .
- In general, training and/or inference are **hard** problems:
  - ▶ The loss function to be minimized is often neither convex nor differentiable.
- Two approaches to make the learning task easier: surrogate loss minimization and **reduction**.

## Reduction



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

## Statistical consistency

- We say that a surrogate (proxy) loss  $\tilde{\ell}$  is **consistent (calibrated)** with the task loss  $\ell$  when the following holds:

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

- The definition concerns both surrogate loss minimization and reduction algorithms:<sup>1</sup>
  - ▶ Surrogate loss minimization:  $\tilde{\ell} =$  surrogate loss.
  - ▶ Reduction:  $\tilde{\ell} =$  loss used in the reduced problem.

---

<sup>1</sup> P. Bartlett, M. Jordan, and J. McAuliffe. Convexity, classification and risk bounds. *JASA*, 101:138–156, 2006

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

A. Beygelzimer, J. Langford, and B. Zadrozny. Machine Learning Techniques—Reductions Between Prediction Quality Metrics. In *Performance Modeling and Engineering*, pages 3–28. Springer, 2008



## Top- $k$ classification

- Training data:  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ .
- **Predict** top  $k$  classes  $\{y_j\}_{j=1}^k$  for a given  $\mathbf{x}$ .

	$x_1$	$x_2$	$x_3$	$y$
$\mathbf{x}_1$	5.0	4.5	1.2	1
$\mathbf{x}_2$	5.0	4.5	4.3	5
$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\mathbf{x}_n$	3.0	3.5	0.5	16
$\mathbf{x}$	4.0	2.5	0.8	?

## Top- $k$ classification

- Training data:  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ .
- **Predict** top  $k$  classes  $\{y_j\}_{j=1}^k$  for a given  $\mathbf{x}$ .

	$x_1$	$x_2$	$x_3$	$y$
$\mathbf{x}_1$	5.0	4.5	1.2	1
$\mathbf{x}_2$	5.0	4.5	4.3	5
$\vdots$	$\vdots$	$\vdots$	$\vdots$	
$\mathbf{x}_n$	3.0	3.5	0.5	16
$\mathbf{x}$	4.0	2.5	0.8	$\{1, 3, 16\}$

## Top- $k$ classification

- In top- $k$  classification we define the prediction function as:

$$h(\mathbf{x}) = \text{top}_f(k),$$

where  $\text{top}_f(k)$  returns top  $k$  classes with the highest score  $f$  for a given  $\mathbf{x}$ .

- The learning problem is then defined as minimization of the  $1-\text{prec}@k$  loss:

$$\ell_{@k}(y, \mathbf{x}, f) = \mathbb{1}[y \notin \text{top}_f(k)].$$

## Top- $k$ classification

- The **conditional risk** for the  $1-\text{prec}@k$  loss is:

$$L_{@k}(h | \mathbf{x}) = 1 - \sum_{y \in \text{top}_f(k)} P(y | \mathbf{x})$$

- Therefore, the **Bayes classifier** for the  $1-\text{prec}@k$  loss predicts top  $k$  classes with largest conditional probabilities  $P(y | \mathbf{x})$ , i.e.:

$$h^*(\mathbf{x}) = \text{top}_P(k).$$

- The **conditional regret** for the  $1-\text{prec}@k$  is:

$$\text{Reg}_{@k}(h | \mathbf{x}) = \sum_{y \in \text{top}_P(k)} P(y | \mathbf{x}) - \sum_{y \in \text{top}_f(k)} P(y | \mathbf{x}).$$

## Top-1 classification

- Top-1 classification corresponds to minimization of 0/1 loss:

$$\ell_{0/1} = \llbracket y \neq h(\mathbf{x}) \rrbracket.$$

- The simplest reduction, referred to as **1vsAll**, relies on solving  $m$  binary classification problems of the form:

$$(\mathbf{x}, y) \longrightarrow (\mathbf{x}, y' = \llbracket y = j \rrbracket) \quad j \in \mathcal{Y}.$$

- Training and prediction is **linear** in  $m$ .
- **Consistent** if used with probabilistic classifiers.

## Label tree structure

- Organize class labels into a tree to improve computational complexity:<sup>2</sup>
  - ▶ Binary coding of class labels.
  - ▶ Hierarchical clustering.
  - ▶ Tree structure learning.

---

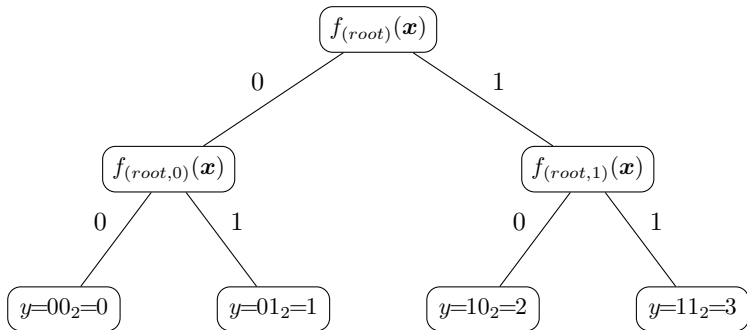
<sup>2</sup> A. Beygelzimer, J. Langford, Y. Lifshits, G. B. Sorkin, and A. L. Strehl. Conditional probability tree estimation analysis and algorithms. In *UAI*, pages 51–58, 2009

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

## Label tree structure

- Assign each class an integer from 0 to  $m - 1$  and use a binary code of length  $t$ .
- Leafs correspond to class labels.
- Internal nodes (including the root) contain binary classifiers.



## Label tree structure

- Different training schemes possible:
  - ▶ Train a classifier in each node.
  - ▶ Train a classifier on each level.
  - ▶ Train one global binary classifier (in several loops).
- Complexity of learning is in fact  $O(\log m)$ :
  - ▶ We need only  $\log m$  copies of a given training example (on a path from the root to the corresponding leaf).
- Prediction can be done in  $O(\log m)$ .



## Greedy prediction

- Greedy prediction follows one path from the root to the leaf:

$$\begin{aligned}f_{(\text{root})} &: \mathbf{x} \mapsto \hat{y}_1 \\f_{(\text{root}, \hat{y}_1)} &: \mathbf{x} \mapsto \hat{y}_2 \\ \dots & \dots \\f_{(\text{root}, \hat{y}_1, \dots, \hat{y}_{t-1})} &: \mathbf{x} \mapsto \hat{y}_t\end{aligned}$$

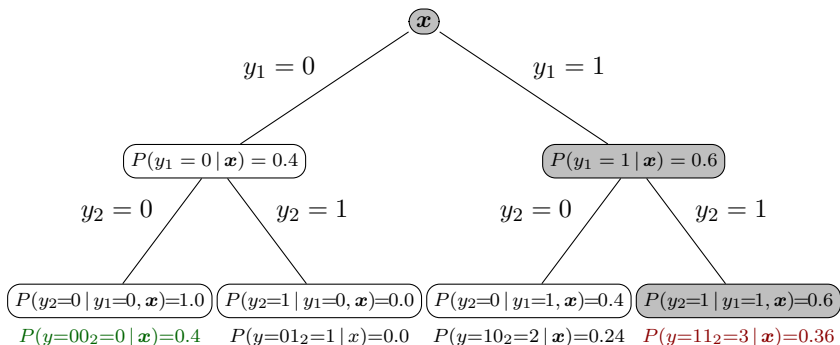
- The final prediction is

$$\hat{y} = \text{decode}(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_t)$$

- Greedy search is fast ( $O(\log m)$ ).
- Any binary classifier can be used.
- However, this approach is **inconsistent**.

## Greedy prediction

- Greedy prediction fails for finding the most probable class:



## Filter trees

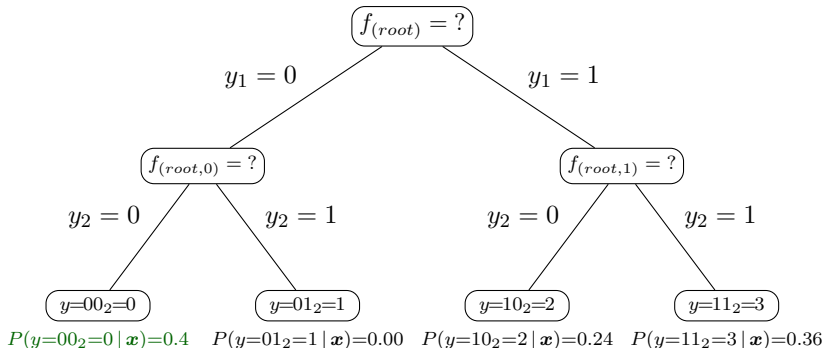
- Filter trees<sup>3</sup> (FT) use a **bottom-up** learning algorithm to train the label tree  $\rightarrow$  a single **elimination tournament** on the set of class labels.
- FT implicitly transform the underlying distribution  $P$  over multi-class examples into a **specific distribution**  $P^{\text{FT}}$  over weighted binary examples.
- FT use the greedy prediction.
- FT are **consistent** for 0/1 loss and cost-sensitive classification.

---

<sup>3</sup> A. Beygelzimer, J. Langford, and P. D. Ravikumar. Error-correcting tournaments. In *ALT*, pages 247–262, 2009

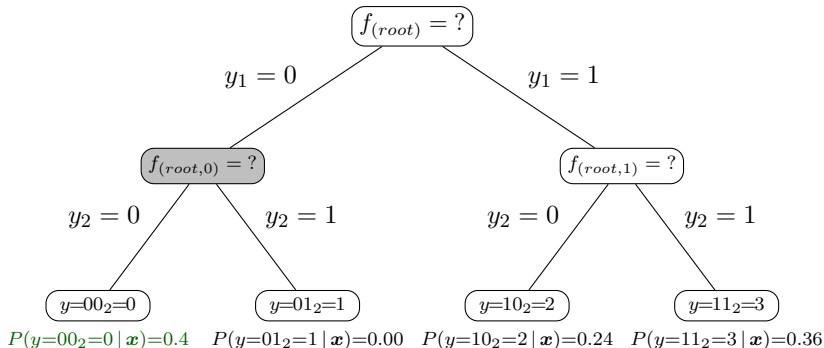
## Filter trees

- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.



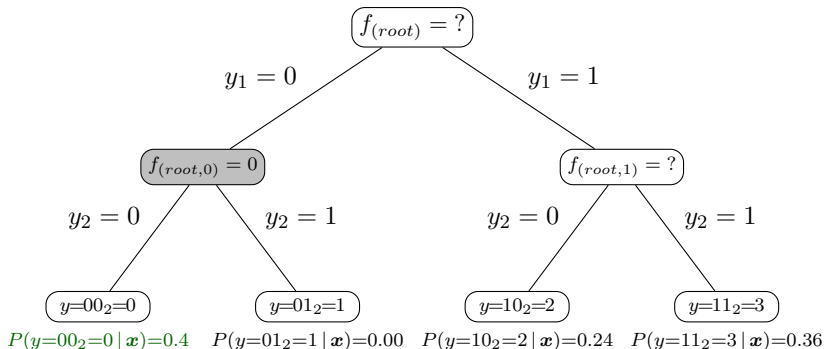
## Filter trees

- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.



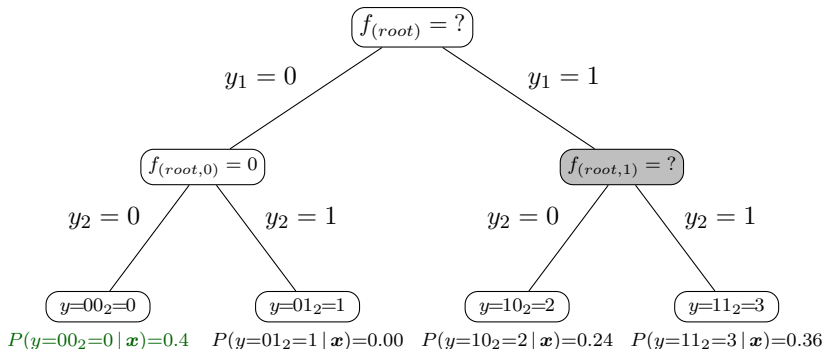
## Filter trees

- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.



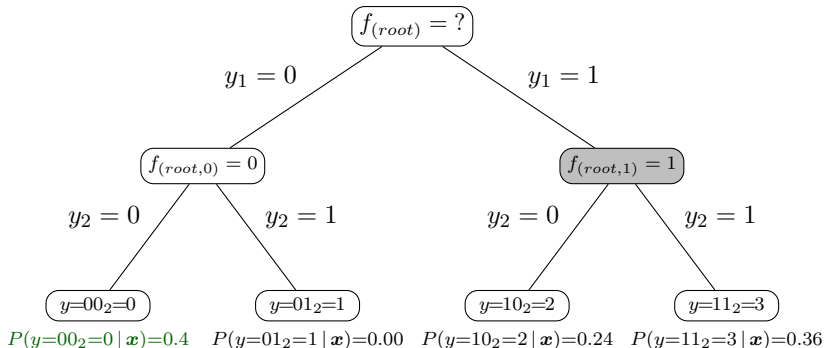
## Filter trees

- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.



## Filter trees

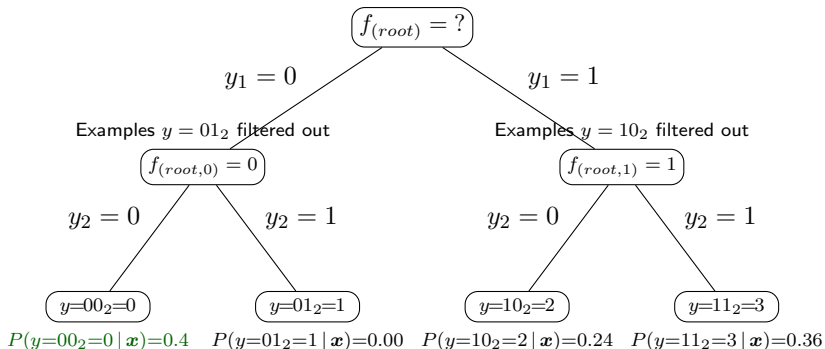
- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.





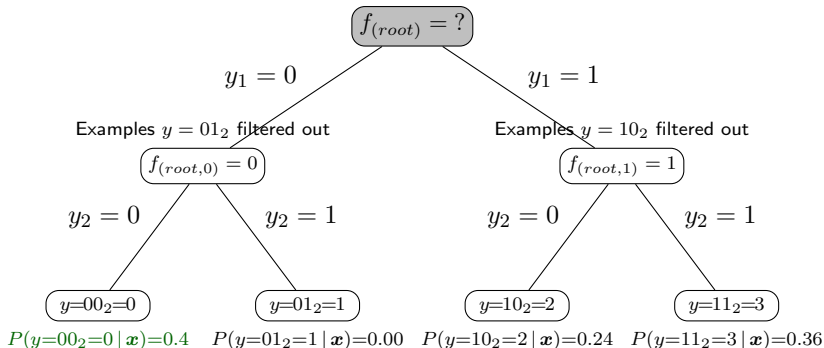
## Filter trees

- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.



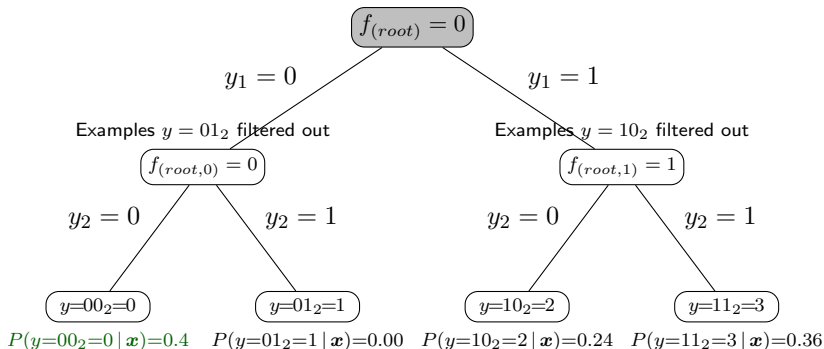
## Filter trees

- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.



## Filter trees

- FT **filter out** all examples that are misclassified by the lower-level classifiers.
- $f_{(root,y_1,\dots,y_i)}(\mathbf{x})$  predicts  $y_{i+1}$  given that all classifiers below predict the subsequent decisions correctly.



## Filter trees

- FT have strong guarantees in the form of the regret bound:

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

- However, they cannot be easily adapted to the problem of top- $k$  classification.

## Probabilistic classifier trees

- Let the node classifiers  $f_{(root, y_1, \dots, y_{i-1})}$  estimate  $P(y_i | \mathbf{x}, y_1, \dots, y_{i-1})$ :

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

$$Q(y_i = 0 | \mathbf{x}, y_1, \dots, y_{i-1}) = 1 - f_{(root, y_1, \dots, y_{i-1})}$$

- Then, we can easily estimate  $P(y | \mathbf{x})$  by:<sup>4</sup>

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

- We refer to this approach as probabilistic classifier trees (PCT).

---

<sup>4</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  
K. Dembczyński, W. Cheng, and E. Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In *ICML*, pages 279–286, 2010

## Advanced search techniques

- Finding the most probable class label relies on finding the most probable path in the tree.
- Advanced search techniques: beam search,<sup>5</sup> uniform-cost search.<sup>6</sup>
- The use of a priority queue and a cut point gives a fast algorithm with provable guarantees.

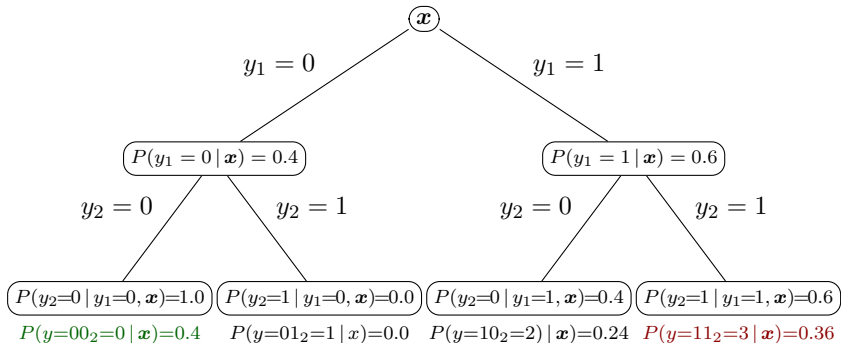
---

<sup>5</sup> A. Kumar, S. Vembu, A.K. Menon, and C. Elkan. Beam search algorithms for multilabel learning. In *Machine Learning*, 2013

<sup>6</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An analysis of chaining in multi-label classification. In *ECAI*, 2012

## Advanced search techniques

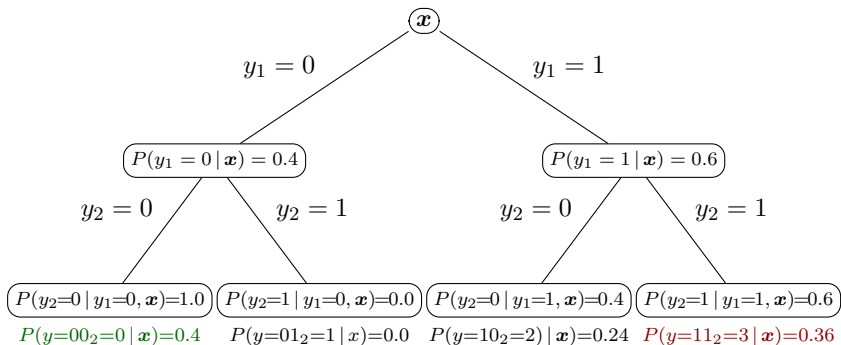
- Uniform-cost search



- Priority list  $Q$ :

## Advanced search techniques

- Uniform-cost search

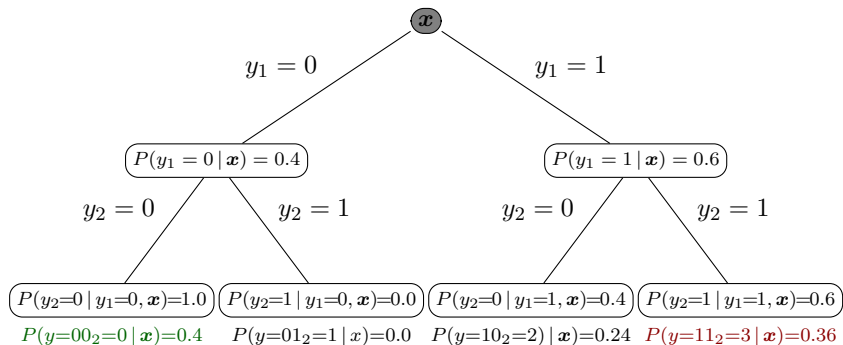


- Priority list  $Q$ : *root*



## Advanced search techniques

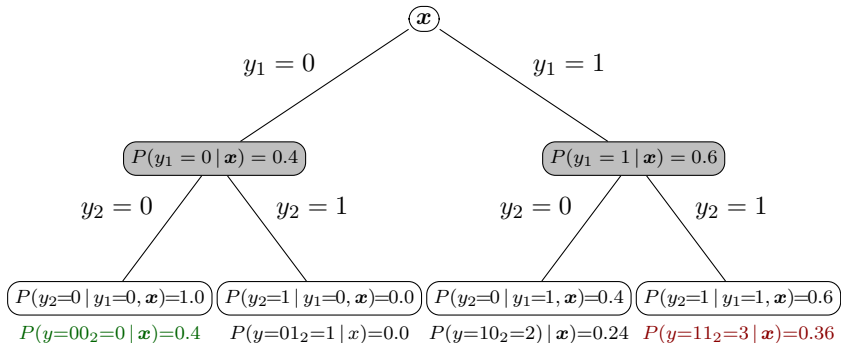
- Uniform-cost search



- Priority list  $Q$ :

## Advanced search techniques

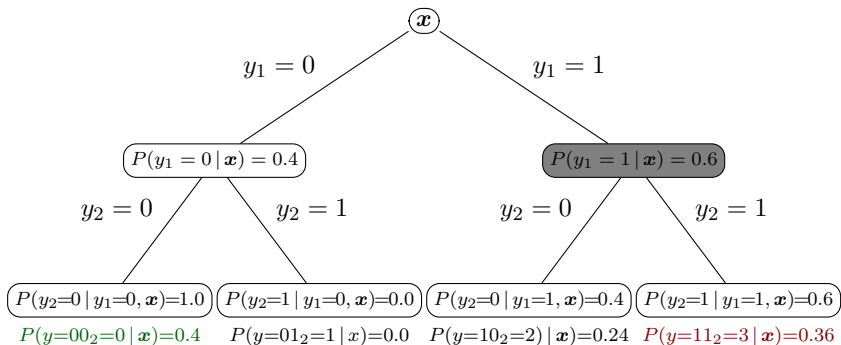
- Uniform-cost search



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

## Advanced search techniques

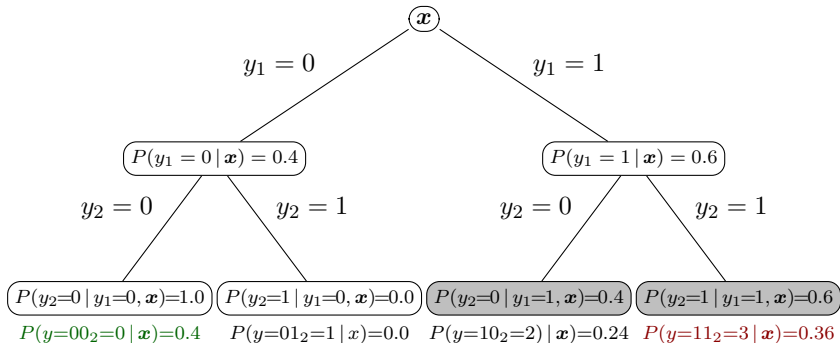
- Uniform-cost search



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

## Advanced search techniques

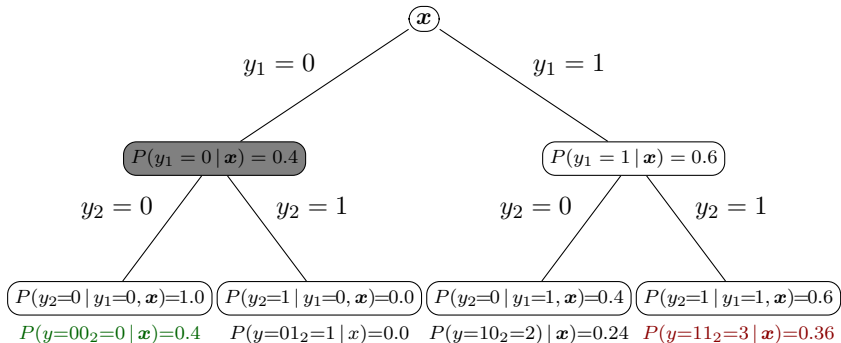
- Uniform-cost search



- Priority list  $Q$ :  $[(root, 0), 0.4]$ ,  $[(root, 1, 1), 0.36]$ ,  $[(root, 1, 0), 0.24]$

## Advanced search techniques

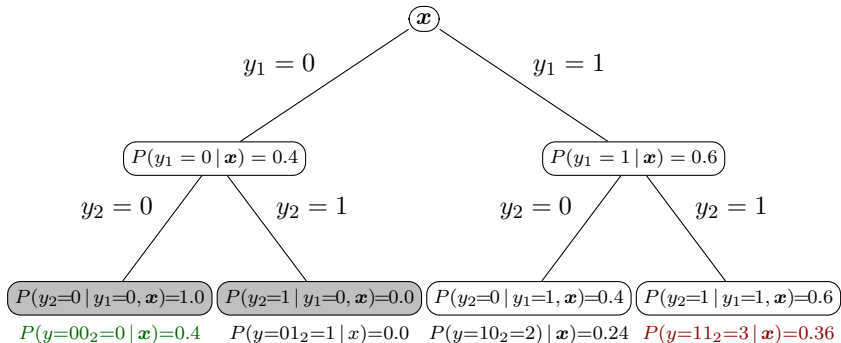
- Uniform-cost search



- Priority list  $Q$ :  $[(1, 1), 0.36], [(1, 0), 0.24]$

## Advanced search techniques

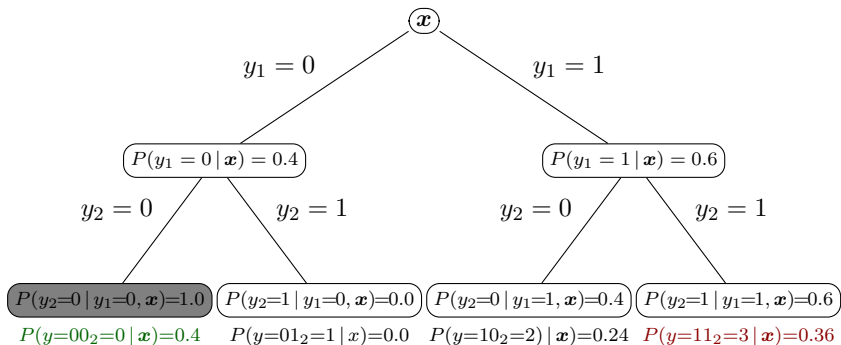
- Uniform-cost search



- Priority list  $Q$ :  $[(root, 0, 0), 0.4]$ ,  $[(root, 1, 1), 0.36]$ ,  $[(root, 1, 0), 0.24]$ ,  $[(root, 0, 1), 0.0]$

## Advanced search techniques

- Uniform-cost search



- Priority list  $Q$ : Solution is found

## Advanced search techniques

- $\epsilon$ -approximation inference:<sup>7</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.

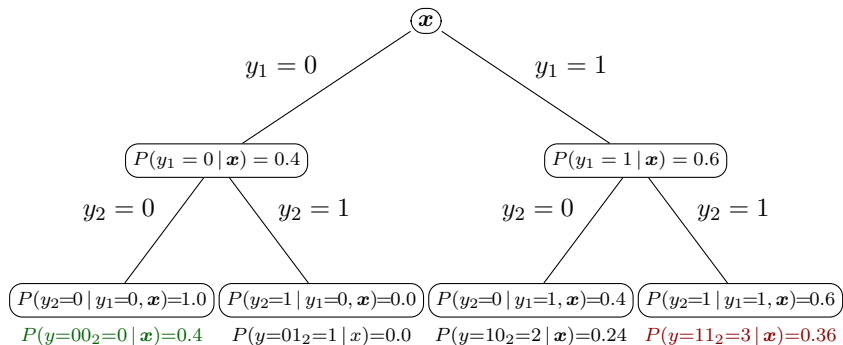
---

<sup>7</sup> K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. An analysis of chaining in multi-label classification. In *ECAI*, 2012



## $\epsilon$ -approximation inference

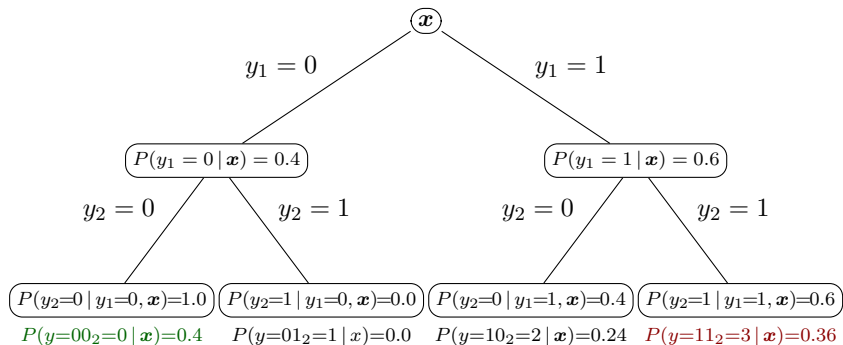
- $\epsilon = 0.5$



- Priority list  $Q$ :

## $\epsilon$ -approximation inference

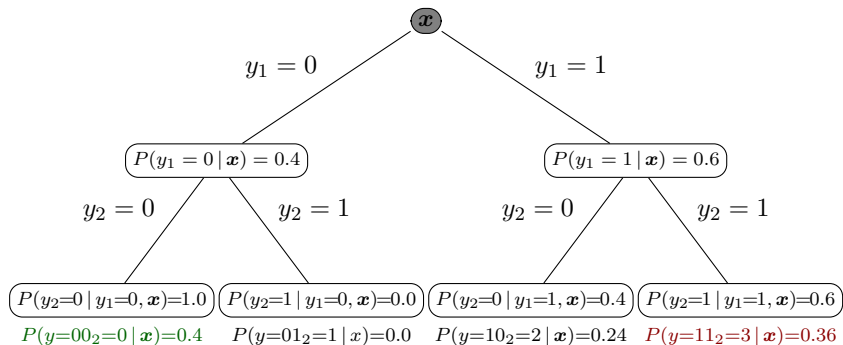
- $\epsilon = 0.5$



- Priority list  $Q$ : *root*

## $\epsilon$ -approximation inference

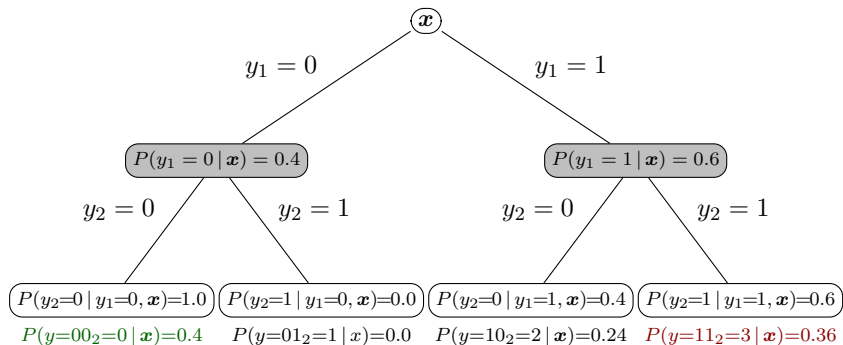
- $\epsilon = 0.5$



- Priority list  $Q$ :  $\epsilon = 0.5$

## $\epsilon$ -approximation inference

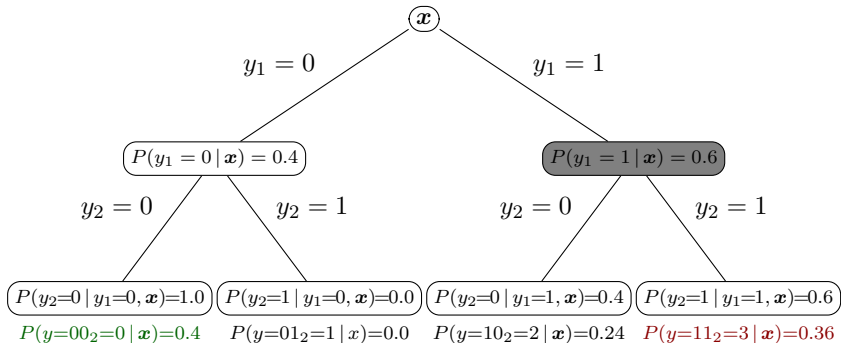
- $\epsilon = 0.5$



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

## $\epsilon$ -approximation inference

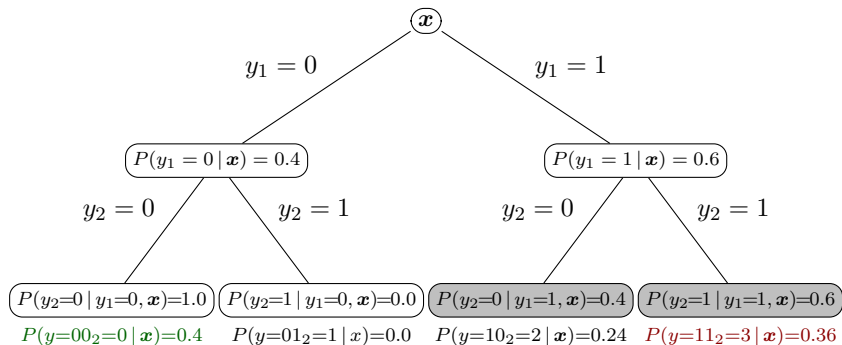
- $\epsilon = 0.5$



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

## $\epsilon$ -approximation inference

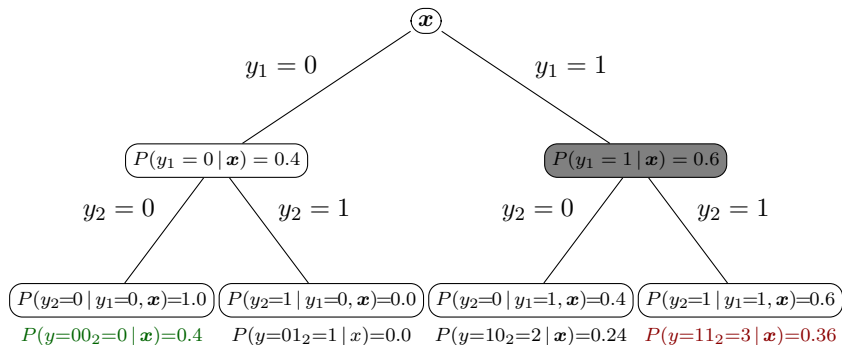
- $\epsilon = 0.5$



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

## $\epsilon$ -approximation inference

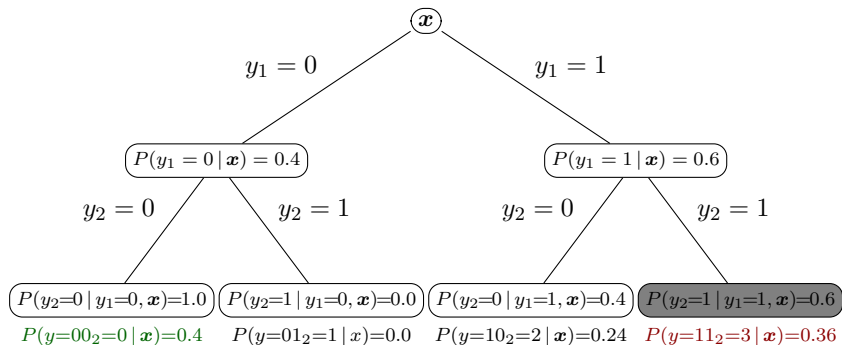
- $\epsilon = 0.5$



- Priority list  $Q$ : Start the greedy search from  $(root, 1)$ .

## $\epsilon$ -approximation inference

- $\epsilon = 0.5$



- Priority list  $\mathcal{Q}$ : Suboptimal solution  $11_2$  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 guarantee can be given:

## $\epsilon$ -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 guarantee can be given:

**Theorem:** Let  $\epsilon = 2^{-c}$ , where  $1 \leq c \leq t$ . The class label  $\hat{y}$  will be returned in time  $\mathcal{O}(t2^c)$  with a guarantee that:

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

## $\epsilon$ -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 prediction has very bad guarantees:

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

## Regret bound for PCT

- The overall guarantees of PCT can be formalized as follows:

## Regret bound for PCT

- The overall guarantees of PCT can be formalized as follows:

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

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

## PCT for top- $k$ classification

- Training of PCT does not change for top- $k$  classification.
- The  $\epsilon$ -approximate inference can be easily extended to the case of searching top- $k$  class labels.
- The overall guarantees will slightly change by a factor of  $k$ .

## Experimental results

### Top-1 multi-class classification

		0/1[%]	$t_{train}$	$t_{test}$	A	$min$	$max$	1q	2q	3q
ILSVR2010 $m = 1000$	1vsA	<b>91.80</b>	13435	581.43	1000	1000	1000	1000	1000	1000
	FT	95.77	<b>1051</b>	<b>24.50</b>	<b>10</b>	10	10	10	10	10
	PCT $_{\epsilon=0.5}$	96.92	1148	27.18	<b>10</b>	10	10	10	10	10
	PCT $_{\epsilon=0.25}$	94.99	1148	43.60	22	10	30	21	21	21
	PCT $_{\epsilon=0.0}$	92.93	1148	69.56	55	10	192	38	52	69
	PCT* $_{\epsilon=0.0}$	92.66	2379	109.55	94	17	237	75	91	111

### Top-5 multi-class classification

		1-P@5	$t_{train}$	$t_{test}$	A	$min$	$max$	1q	2q	3q
ILSVR2010 $m = 1000$	1vsA	<b>80.25</b>	13435	581.43	1000	1000	1000	1000	1000	1000
	PCT $_{\epsilon=0.5}$	90.1	<b>1148</b>	<b>75.08</b>	<b>41</b>	39	55	42	42	42
	PCT $_{\epsilon=0.25}$	88.39	<b>1148</b>	89.05	55	39	90	52	55	59
	PCT $_{\epsilon=0.0}$	81.39	<b>1148</b>	124.06	100	15	251	79	97	118
	PCT* $_{\epsilon=0.0}$	<b>80.25</b>	2379	106.75	88	14	233	70	85	103

\* – SGD with 25 epochs (for other methods 5 epochs have been used).

## FT vs. PCT

### Filter Trees

- Logarithmic prediction time.
- Can be used with any binary classifier.
- Training is more demanding.
- Filtering may reduce the number of training examples in the top levels of the tree.
- Does not work for top- $k$  classification.

### Probabilistic Classifier Trees

- The prediction scales from logarithmic to linear time.
- Requires probabilistic classifiers.
- Training can be parallelized.
- The  $\epsilon$ -approximate inference may not work properly for noisy problems.
- Can be adapted for top- $k$  classification.



## Conclusions

- Take-away message:
  - ▶ Label tree structures for efficient classification.
  - ▶ Standard approach is inconsistent.
  - ▶ FT are consistent for 0/1 loss and cost-sensitive classification.
  - ▶ PCT are consistent for top- $k$  classification.
- For more check:

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