

A Separate-and-Conquer Algorithm for Learning Multi-Label Head Rules



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Master-Thesis by Michael Rapp

1. Referee: Prof. Dr. Johannes Fürnkranz
2. Referee: Dr. Eneldo Loza Mencía
3. Referee: Dr. Frederik Janssen

1 Foundations

2 Learning Single-Label Head Rules

3 Learning Multi-Label Head Rules


4 Evaluation

1 Foundations

1.1 Motivation

- ▶ **Multi-label classification** is the task of assigning a subset of (relevant) **labels** to an object. This is in contrast to
 - ▶ **Binary classification**, where each object is associated with one out of two classes
 - ▶ **Multi-class classification**, where each object is associated with one out of several classes
- ▶ **Typical application areas** are
 - ▶ Tagging books, news, music,... with topics, genres, authors,...
 - ▶ Detecting presence or absence of objects in images
 - ▶ Classifying function of genomes or proteins in biology
 - ▶ ...

Pink Floyd



Pink Floyd in January 1968, from the only known photoshoot of all five members. Clockwise from bottom: Gilmour, Mason, Barrett, Waters, Wright

Background information

| | |
|---------------------|---|
| Origin | London, England |
| Genres | Progressive rock · art rock · psychedelic rock · acid rock · proto-prog · psychedelic pop |
| Years active | 1965–1995 · 2012–2014 (one-off reunion: 2005) |
| Labels | EMI Columbia · Tower · Harvest · Capitol · Columbia · EMI · Parlophone |
| Website | pinkfloyd.com |
| Past members | Nick Mason Roger Waters Richard Wright Syd Barrett David Gilmour |

1 Foundations

1.2 Label-dependencies



- ▶ Studies have shown, that the predictive performance of classification approaches may benefit from exploiting correlations between labels¹.
- ▶ Possible correlations are for example
 - ▶ **Exclusions:** The presence of certain labels might imply, that other labels are irrelevant.
e.g. the music genres “progressive rock” and “easy listening” are contrary
 - ▶ **Subsumptions:** Labels can be ordered hierarchical.
e.g. “progressive rock” is a sub-genre of “rock”

¹Krzysztof Dembczyński, Willem Waegeman, Weiwei Cheng, and Eyke Hüllermeier. On label dependence and loss minimization in multi-label classification. *Machine learning*, 88(1-2):5–45, 2012.

1 Foundations

1.3 Example data set

| Attributes | | | | Labels | | | |
|------------|-------|----------|-------|--------|----------|-------------|----------|
| outlook | temp. | humidity | windy | play | dontplay | eaticecream | drinktee |
| sunny | 85 | 85 | FALSE | 0 | 1 | 1 | 0 |
| sunny | 80 | 90 | TRUE | 0 | 1 | 1 | 0 |
| overcast | 83 | 86 | FALSE | 1 | 0 | 0 | 0 |
| rainy | 70 | 96 | FALSE | 1 | 0 | 0 | 0 |
| rainy | 68 | 80 | FALSE | 1 | 0 | 0 | 0 |
| rainy | 65 | 70 | TRUE | 0 | 1 | 0 | 1 |
| overcast | 64 | 65 | TRUE | 1 | 0 | 0 | 0 |
| sunny | 72 | 95 | FALSE | 0 | 1 | 1 | 0 |
| sunny | 69 | 70 | FALSE | 1 | 0 | 0 | 0 |
| rainy | 75 | 80 | FALSE | 1 | 0 | 0 | 0 |
| sunny | 75 | 70 | TRUE | ? | ? | ? | ? |
| overcast | 72 | 90 | TRUE | ? | ? | ? | ? |
| overcast | 81 | 75 | FALSE | ? | ? | ? | ? |
| rainy | 71 | 91 | TRUE | ? | ? | ? | ? |

Training examples

Test examples

1 Foundations

1.4 Types of Rules



- ▶ Rules contain conditions in their bodies and label assignments in their heads. They are notated as

$$\textit{head} \leftarrow \textit{body}$$

- ▶ **Single-label head rules** predict exactly one label, e.g.

$$\textit{dontplay} = 1 \leftarrow \textit{outlook} = \textit{rainy} \wedge \textit{windy} = \textit{TRUE}$$

- ▶ **Label-dependent** rules are able to model correlations between labels, e.g.

$$\textit{dontplay} = 1 \leftarrow \textit{play} = 0$$

- ▶ **Multi-label head rules** may predict multiple labels at once, e.g.

$$\textit{play} = 1, \textit{dontplay} = 0 \leftarrow \textit{temperature} \geq 66.5$$



1 Foundations

2 Learning Single-Label Head Rules

3 Learning Multi-Label Head Rules

4 Evaluation

2 Learning Single-Label Head Rules

2.1 Basic Idea

- ▶ Loza Mencía and Janssen proposed a **separate-and-conquer algorithm** for learning single-label head rules².
- ▶ New rules are learned iteratively. They are added to a **decision list** in the order of their induction.
- ▶ After each iteration, training examples, whose labels are predicted to some extent, are removed from the training data set.

²Eneldo Loza Mencía and Frederik Janssen. Learning rules for multi-label classification: A stacking and a separate-and-conquer approach. *Knowledge Engineering Group, Technische Universität Darmstadt, Germany, 2015.*

2 Learning Single-Label Head Rules

2.2 Structure of the Separate-and-Conquer Algorithm

Require: Training data set T , evaluation function δ

- 1 $R = \emptyset$ ▷ Initialize empty decision list
- 2 $T_{current} = T$ with all label vectors set to $(?, \dots, ?)$
- 3 **while** examples are left in T **do** ▷ e.g. until only 5% are left
- 4 $r = \text{FINDBESTGLOBALRULE}(T, T_{current}, \delta)$ ▷ See next slide
- 5 $R = R \cup r$ ▷ Add rule to decision list
- 6 apply head of r to $T_{current}$
- 7 remove covered examples from T ▷ If enough of their labels are predicted
- 8 **return** decision list R

- ▶ The copied data set $T_{current}$ is used to keep track of the labels, which are predicted by already learned rules.
- ▶ These labels can be used to induce label-dependent rules in later iterations.

2 Learning Single-Label Head Rules

2.3 Algorithm `FINDBESTGLOBALRULE`

Require: Training data sets T and $T_{current}$, evaluation function δ

```
1  $r_{best} = \emptyset \leftarrow \emptyset$ 
2 for each label  $y_i$  do ▷ Find best rule for each label
3    $T_i = T \setminus$  examples for which label  $y_i$  is set in  $T_{current}$ 
4    $r = \text{FINDBESTRULE}(T_i, y_i, \delta)$ 
5   if  $r$  outperforms  $r_{best}$  according to  $\delta$  then
6      $r_{best} = r$ 
7 return best rule  $r_{best}$ 
```

- ▶ The subroutine `FINDBESTRULE` can be implemented by using any rule learner for solving binary classification problems.

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification



Decision list:

- 1 $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$
- 2 $eaticecream = 1 \leftarrow dontplay = 1$
- 3 $play = 1 \leftarrow windy = FALSE$
- 4 $dontplay = 1 \leftarrow outlook = rainy$
- 5 $drinktee = 1 \leftarrow dontplay = 1$
- 6 $play = 1 \leftarrow \emptyset$

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|-------------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticecream | drinktea |
| rainy | 71 | 91 | TRUE | ? | ? | ? | ? |

All labels are initially unset

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification



Decision list:

- 1 $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$
- 2 $eaticream = 1 \leftarrow dontplay = 1$
- 3 $play = 1 \leftarrow windy = FALSE$
- 4 $dontplay = 1 \leftarrow outlook = rainy$
- 5 $drinktee = 1 \leftarrow dontplay = 1$
- 6 $play = 1 \leftarrow \emptyset$

does not cover

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|-----------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticream | drinktea |
| rainy | 71 | 91 | TRUE | ? | ? | ? | ? |

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification

Decision list:

- | | | |
|---|---|----------------|
| 1 | $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$ | does not cover |
| 2 | $eaticecream = 1 \leftarrow dontplay = 1$ | does not cover |
| 3 | $play = 1 \leftarrow windy = FALSE$ | |
| 4 | $dontplay = 1 \leftarrow outlook = rainy$ | |
| 5 | $drinktee = 1 \leftarrow dontplay = 1$ | |
| 6 | $play = 1 \leftarrow \emptyset$ | |

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|-------------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticecream | drinktea |
| rainy | 71 | 91 | TRUE | ? | ? | ? | ? |

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification

Decision list:

- | | | |
|---|---|----------------|
| 1 | $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$ | does not cover |
| 2 | $eaticecream = 1 \leftarrow dontplay = 1$ | does not cover |
| 3 | $play = 1 \leftarrow windy = FALSE$ | does not cover |
| 4 | $dontplay = 1 \leftarrow outlook = rainy$ | |
| 5 | $drinktee = 1 \leftarrow dontplay = 1$ | |
| 6 | $play = 1 \leftarrow \emptyset$ | |

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|-------------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticecream | drinktea |
| rainy | 71 | 91 | TRUE | ? | ? | ? | ? |

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification

Decision list:

- | | | |
|---|---|----------------|
| 1 | $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$ | does not cover |
| 2 | $eaticecream = 1 \leftarrow dontplay = 1$ | does not cover |
| 3 | $play = 1 \leftarrow windy = FALSE$ | does not cover |
| 4 | $dontplay = 1 \leftarrow outlook = rainy$ | covers |
| 5 | $drinktee = 1 \leftarrow dontplay = 1$ | |
| 6 | $play = 1 \leftarrow \emptyset$ | |

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|-------------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticecream | drinktea |
| rainy | 71 | 91 | TRUE | ? | 1 | ? | ? |

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification



Decision list:

- | | | |
|---|---|----------------|
| 1 | $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$ | does not cover |
| 2 | $eaticecream = 1 \leftarrow dontplay = 1$ | does not cover |
| 3 | $play = 1 \leftarrow windy = FALSE$ | does not cover |
| 4 | $dontplay = 1 \leftarrow outlook = rainy$ | covers |
| 5 | $drinktee = 1 \leftarrow dontplay = 1$ | covers |
| 6 | $play = 1 \leftarrow \emptyset$ | |

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|-------------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticecream | drinktea |
| rainy | 71 | 91 | TRUE | ? | 1 | ? | 1 |

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification



Decision list:

| | | |
|---|---|----------------|
| 1 | $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$ | does not cover |
| 2 | $eaticcream = 1 \leftarrow dontplay = 1$ | does not cover |
| 3 | $play = 1 \leftarrow windy = FALSE$ | does not cover |
| 4 | $dontplay = 1 \leftarrow outlook = rainy$ | covers |
| 5 | $drinktee = 1 \leftarrow dontplay = 1$ | covers |
| 6 | $play = 1 \leftarrow \emptyset$ | covers |

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|------------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticcream | drinktea |
| rainy | 71 | 91 | TRUE | 1 | 1 | ? | 1 |

2 Learning Single-Label Head Rules

2.4 Using decision lists for classification



Decision list:

| | | |
|---|---|----------------|
| 1 | $dontplay = 1 \leftarrow humidity \geq 82.5 \wedge outlook = sunny$ | does not cover |
| 2 | $eaticcream = 1 \leftarrow dontplay = 1$ | does not cover |
| 3 | $play = 1 \leftarrow windy = FALSE$ | does not cover |
| 4 | $dontplay = 1 \leftarrow outlook = rainy$ | covers |
| 5 | $drinktee = 1 \leftarrow dontplay = 1$ | covers |
| 6 | $play = 1 \leftarrow \emptyset$ | covers |

Test example to classify:

| Attributes | | | | Labels | | | |
|------------|-------------|----------|-------|--------|----------|------------|----------|
| outlook | temperature | humidity | windy | play | dontplay | eaticcream | drinktea |
| rainy | 71 | 91 | TRUE | 1 | 1 | 0 | 1 |

Assume remaining labels to be irrelevant!



1 Foundations

2 Learning Single-Label Head Rules

3 Learning Multi-Label Head Rules

4 Evaluation

3 Learning Multi-Label Head Rules

3.1 Basic Idea



- ▶ The structure of the separate-and-conquer algorithm – including the removal of training examples – remains unchanged.
- ▶ For classification, the rules of a decision list are successively applied as before. However, several labels might be set by multi-label head rules.
- ▶ To induce multi-label head rules, the algorithm `FINDBESTGLOBALRULE` must be modified, because rules cannot be induced for each label individually.

3 Learning Multi-Label Head Rules

3.2 Algorithm `FINDBESTGLOBALRULE`

Require: Training data sets T and $T_{current}$, evaluation function δ

```
1  $r_{best} = \emptyset \leftarrow \emptyset$ 
2 improved = true
3 while improved do                                     ▷ Until no refinements possible
4    $r = \text{REFINERULE}(T, T_{current}, r_{best}, \delta)$           ▷ See next slide
5   if  $r$  outperforms  $r_{best}$  according to  $\delta$  then
6      $r_{best} = r$ 
7   else
8     improved = false
9 return best rule  $r_{best}$ 
```

- ▶ Starting with the most-general rule, rules are specialized by adding additional conditions to their bodies.

3 Learning Multi-Label Head Rules

3.3 Algorithm REFINERULE



Require: Training data sets T and $T_{current}$, rule r , evaluation function δ

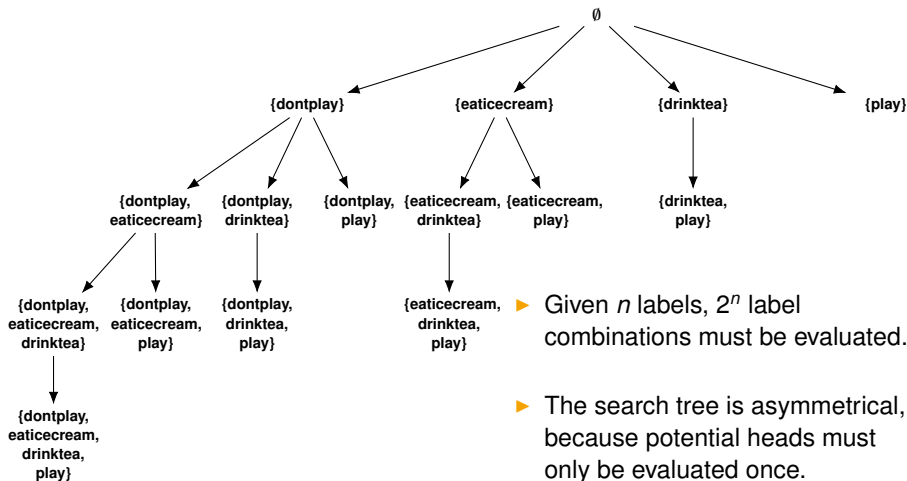
```
1  $r_{best} = r$ 
2 for each possible condition  $c$  not already in  $r$  do
3    $r_{refined} = r$ 
4   add  $c$  to body of  $r_{refined}$ 
5    $r_{refined} = \text{FINDBESTHEAD}(T, T_{current}, r_{refined}, \delta)$ 
6   if  $r_{refined}$  outperforms  $r_{best}$  according to  $\delta$  then
7      $r_{best} = r_{refined}$ 
8 return best refined rule  $r_{best}$ 
```

▷ See next slide

- ▶ For each potential refinement, the best multi-label head must be found.
- ▶ The algorithm FINDBESTHEAD performs (pruned) searches through the label space.

3 Learning Multi-Label Head Rules

3.4 Exhaustive Search through the Label Space



3 Learning Multi-Label Head Rules

3.5 Evaluation Functions (1/2)

- ▶ **Precision:** Percentage of correctly predicted labels among all predicted labels

$$\delta_{prec} = \frac{TP}{TP + FP}$$

- ▶ **Recall:** Percentage of correctly predicted labels among all relevant labels

$$\delta_{rec} = \frac{TP}{TP + FN}$$

- ▶ **Hamming Accuracy:** Percentage of correctly predicted labels among all labels

$$\delta_{hamm} = \frac{TP + TN}{TP + TN + FP + FN}$$

3 Learning Multi-Label Head Rules

3.5 Evaluation Functions (2/2)



- ▶ **F-Measure:** Weighted harmonic mean between precision and recall

$$\delta_f = \frac{(\beta^2 + 1) \cdot \delta_{prec} \cdot \delta_{rec}}{\beta^2 \cdot \delta_{prec} + \delta_{rec}}, \text{ with } \beta \in [0, \infty]$$

- ▶ **Subset Accuracy:** Percentage of perfectly predicted label vectors

$$\delta_{acc} = \frac{1}{m} \sum_{j=1}^m [Y_j = \hat{Y}_j], \text{ with } [x] = \begin{cases} 1, & \text{if } x \text{ is true} \\ 0, & \text{otherwise} \end{cases}$$

3 Learning Multi-Label Head Rules

3.6 Confusion matrices

- ▶ **True positives:** Correctly predicted labels of covered examples.
- ▶ **False positives:** Incorrectly predicted labels of covered examples.
- ▶ **True negatives:** Irrelevant labels of uncovered examples.
- ▶ **False negatives:** Relevant labels of uncovered examples.

Candidate rule:

$$play = 0, dontplay = 1 \leftarrow outlook = sunny$$

Training examples:

| Attributes | | | | Labels | | | | |
|------------|-------|----------|-------|--------|----------|------------|----------|---------------|
| outlook | temp. | humidity | windy | play | dontplay | eaticcream | drinktee | |
| sunny | 72 | 80 | TRUE | 0 | 1 | 1 | 0 | } covered |
| sunny | 68 | 72 | FALSE | 0 | 1 | 1 | 0 | |
| overcast | 65 | 83 | FALSE | 1 | 0 | 0 | 0 | } not covered |
| rainy | 64 | 75 | FALSE | 1 | 0 | 0 | 0 | |

3 Learning Multi-Label Head Rules

3.7 Different Averaging Strategies (1/4)

- ▶ **Micro-Averaging:** Evaluation function is applied to a global confusion matrix.

$$h = \delta \left(\begin{array}{c|c} 6 & 2 \\ \hline 2 & 6 \end{array} \right)$$

| Attributes | | | | Labels | | | | |
|------------|-------|----------|-------|--------|----------|-------------|----------|---------------|
| outlook | temp. | humidity | windy | play | dontplay | eaticecream | drinktee | |
| sunny | 72 | 80 | TRUE | 0 | 1 | 1 | 0 | } covered |
| sunny | 68 | 72 | FALSE | 0 | 1 | 1 | 0 | |
| overcast | 65 | 83 | FALSE | 1 | 0 | 0 | 0 | } not covered |
| rainy | 64 | 75 | FALSE | 1 | 0 | 0 | 0 | |

3 Learning Multi-Label Head Rules

3.7 Different Averaging Strategies (2/4)

- ▶ **Label-based Averaging:** A confusion matrix is created per label. The evaluation function is applied to each one and the results are averaged.

$$h = \frac{1}{4} \cdot \left(\underbrace{\delta \begin{pmatrix} 2 & 2 \\ 0 & 0 \end{pmatrix}}_{\text{play}} + \underbrace{\delta \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}}_{\text{dontplay}} + \underbrace{\delta \begin{pmatrix} 0 & 0 \\ 2 & 2 \end{pmatrix}}_{\text{eaticecream}} + \underbrace{\delta \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}}_{\text{drinktee}} \right)$$

| Attributes | | | | Labels | | | | |
|------------|-------|----------|-------|--------|----------|-------------|----------|---------------|
| outlook | temp. | humidity | windy | play | dontplay | eaticecream | drinktee | |
| sunny | 72 | 80 | TRUE | 0 | 1 | 1 | 0 | } covered |
| sunny | 68 | 72 | FALSE | 0 | 1 | 1 | 0 | |
| overcast | 65 | 83 | FALSE | 1 | 0 | 0 | 0 | } not covered |
| rainy | 64 | 75 | FALSE | 1 | 0 | 0 | 0 | |

3 Learning Multi-Label Head Rules

3.7 Different Averaging Strategies (3/4)

- ▶ **Example-based Averaging:** A confusion matrix is created per example. The evaluation function is applied to each one and the results are averaged.

$$h = \frac{1}{4} \cdot \left(\underbrace{\delta \begin{pmatrix} 3 & 0 \\ 1 & 0 \end{pmatrix}}_{\text{1st example}} + \underbrace{\delta \begin{pmatrix} 3 & 0 \\ 1 & 0 \end{pmatrix}}_{\text{2nd example}} + \underbrace{\delta \begin{pmatrix} 0 & 1 \\ 0 & 3 \end{pmatrix}}_{\text{3rd example}} + \underbrace{\delta \begin{pmatrix} 0 & 1 \\ 0 & 3 \end{pmatrix}}_{\text{4th example}} \right)$$

| Attributes | | | | Labels | | | | |
|------------|-------|----------|-------|--------|----------|------------|----------|---------------|
| outlook | temp. | humidity | windy | play | dontplay | eaticcream | drinktee | |
| sunny | 72 | 80 | TRUE | 0 | 1 | 1 | 0 | } covered |
| sunny | 68 | 72 | FALSE | 0 | 1 | 1 | 0 | |
| overcast | 65 | 83 | FALSE | 1 | 0 | 0 | 0 | } not covered |
| rainy | 64 | 75 | FALSE | 1 | 0 | 0 | 0 | |

3 Learning Multi-Label Head Rules

3.7 Different Averaging Strategies (4/4)

- ▶ **Macro-Averaging:** A confusion matrix is created per example and label. The evaluation function is applied to each one and the results are averaged example- and label-wise.

$$h = \frac{1}{16} \cdot \left(\underbrace{\delta \left(\begin{array}{c|c} 1 & 0 \\ \hline 0 & 0 \end{array} \right)}_{\text{1st example / play}} + \underbrace{\delta \left(\begin{array}{c|c} 1 & 0 \\ \hline 0 & 0 \end{array} \right)}_{\text{2nd example / play}} + \dots + \underbrace{\delta \left(\begin{array}{c|c} 0 & 0 \\ \hline 0 & 1 \end{array} \right)}_{\text{4th example / drinktee}} \right)$$

| Attributes | | | | Labels | | | | |
|------------|-------|----------|-------|--------|----------|------------|----------|---------------|
| outlook | temp. | humidity | windy | play | dontplay | eaticcream | drinktee | |
| sunny | 72 | 80 | TRUE | 0 | 1 | 1 | 0 | } covered |
| sunny | 68 | 72 | FALSE | 0 | 1 | 1 | 0 | |
| overcast | 65 | 83 | FALSE | 1 | 0 | 0 | 0 | } not covered |
| rainy | 64 | 75 | FALSE | 1 | 0 | 0 | 0 | |

3 Learning Multi-Label Head Rules

3.8 Rule-dependent vs. Rule-independent Evaluation

- ▶ **Rule-dependent:** Only labels, which are predicted by a rule, are taken into account.
- ▶ **Rule-independent:** All labels are taken into account. Labels, which are not contained in the head, are considered as unpredicted.

Candidate rule:

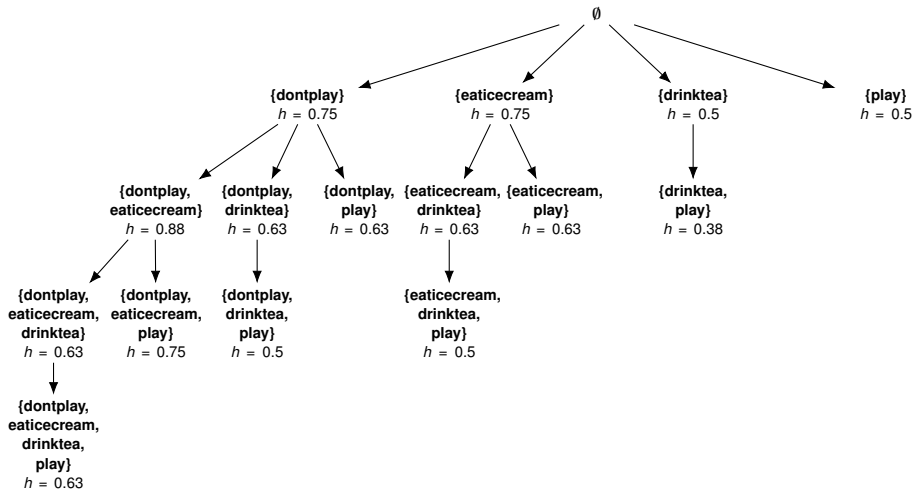
$$play = 0, dontplay = 1 \leftarrow outlook = sunny$$

Training examples:

| Attributes | | | | Labels | | | | |
|------------|-------|----------|-------|--------|----------|-------------|----------|---------------|
| outlook | temp. | humidity | windy | play | dontplay | eaticecream | drinktee | |
| sunny | 80 | 90 | TRUE | 0 | 1 | 1 | 0 | } covered |
| sunny | 72 | 95 | FALSE | 0 | 1 | 1 | 0 | |
| overcast | 83 | 86 | FALSE | 1 | 0 | 0 | 0 | } not covered |
| rainy | 75 | 80 | FALSE | 1 | 0 | 0 | 0 | |

3 Learning Multi-Label Head Rules

3.9 Heuristic Values of Potential Heads



3 Learning Multi-Label Head Rules

3.10 Pruning according to Anti-Monotonicity (1/2)

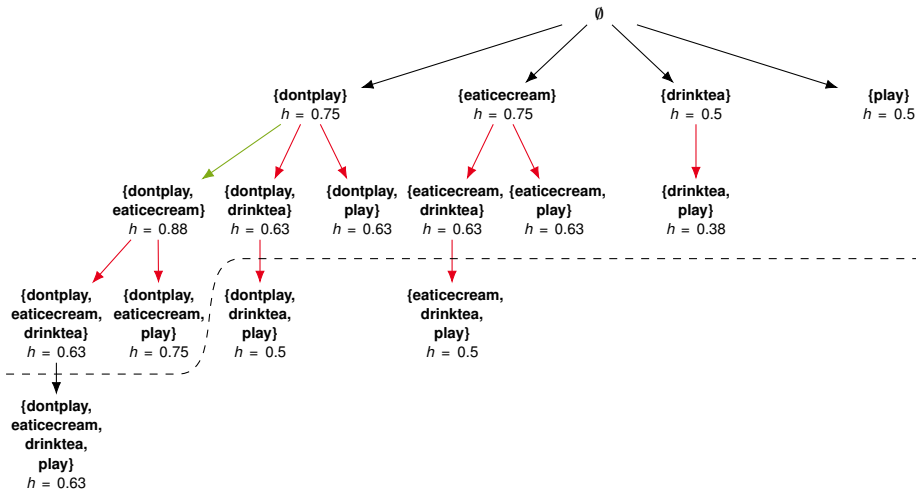


▶ **Definition “Anti-Monotonicity”:**

If adding a label to a rule’s head causes the performance to decrease, by adding more labels the best possible performance cannot be reached anymore.

3 Learning Multi-Label Head Rules

3.10 Pruning according to Anti-Monotonicity (2/2)



3 Learning Multi-Label Head Rules

3.11 Pruning according to Decomposability (1/2)



- ▶ Imagine the following multi-label head rule:

$$\{y_1, \dots, y_n\} \leftarrow \textit{body}$$

- ▶ For each label a corresponding single-label head rule can be created:

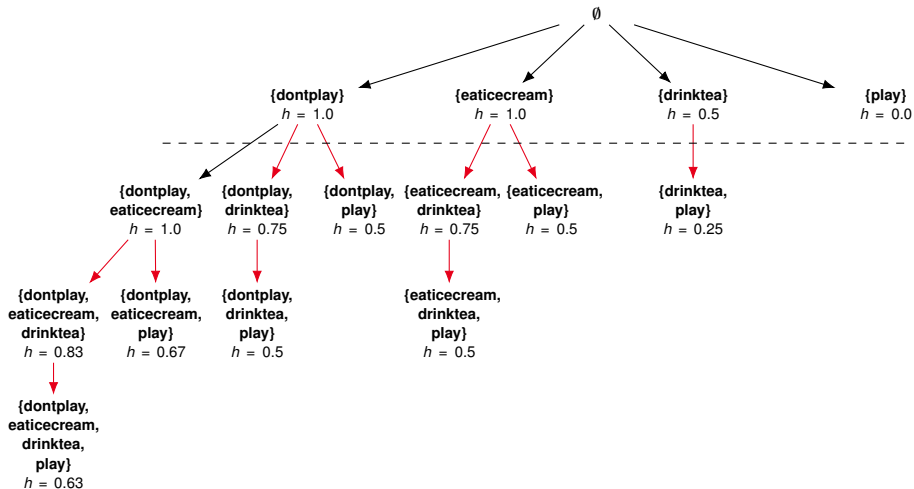
$$y_1 \leftarrow \textit{body}, \dots, y_n \leftarrow \textit{body}$$

- ▶ **Definition “Decomposability”:**

1. If all single-label head rules reach an equal performance, the corresponding multi-label head rule reaches that performance as well.
2. If the performance of at least one single-label head rule is less than the best possible performance, the corresponding multi-label head rule does not reach that performance either.

3 Learning Multi-Label Head Rules

3.11 Pruning according to Decomposability (2/2)



3 Learning Multi-Label Head Rules

3.12 Results of Formal Examinations (1/2)



| Evaluation Function | Evaluation Strategy | Averaging Strategy | Anti-Monotonicity | Decomposability |
|---------------------|---------------------|-------------------------|-------------------|-----------------|
| Precision | Rule-dependent | Micro-Averaging | Yes | Yes |
| | | Label-based Averaging | Yes | Yes |
| | | Example-based Averaging | Yes | Yes |
| | | Macro-Averaging | Yes | Yes |
| | Rule-independent | Micro-Averaging | Yes | - |
| | | Label-based Averaging | Yes | - |
| | | Example-based Averaging | Yes | - |
| | | Macro-Averaging | Yes | - |
| Recall | Rule-dependent | Micro-Averaging | Yes | Yes |
| | | Label-based Averaging | Yes | Yes |
| | | Example-based Averaging | - | - |
| | | Macro-Averaging | Yes | Yes |
| | Rule-independent | Micro-Averaging | Yes | - |
| | | Label-based Averaging | Yes | - |
| | | Example-based Averaging | - | - |
| | | Macro-Averaging | Yes | - |

3 Learning Multi-Label Head Rules

3.12 Results of Formal Examinations (2/2)



| Evaluation Function | Evaluation Strategy | Averaging Strategy | Anti-Monotonicity | Decomposability |
|---------------------|---------------------|-------------------------|-------------------|-----------------|
| Hamming Accuracy | Rule-dependent | Micro-Averaging | Yes | Yes |
| | | Label-based Averaging | Yes | Yes |
| | | Example-based Averaging | Yes | Yes |
| | | Macro-Averaging | Yes | Yes |
| | Rule-independent | Micro-Averaging | Yes | - |
| | | Label-based Averaging | Yes | - |
| | | Example-based Averaging | Yes | - |
| | | Macro-Averaging | Yes | - |
| F-Measure | Rule-dependent | Micro-Averaging | Yes | Yes |
| | | Label-based Averaging | Yes | Yes |
| | | Example-based Averaging | Yes | Yes |
| | | Macro-Averaging | Yes | Yes |
| | Rule-independent | Micro-Averaging | Yes | - |
| | | Label-based Averaging | Yes | - |
| | | Example-based Averaging | Yes | - |
| | | Macro-Averaging | Yes | - |
| Subset Accuracy | Rule-dependent | Example-based Averaging | Yes | - |
| | Rule-independent | | - | - |

1 Foundations

2 Learning Single-Label Head Rules

3 Learning Multi-Label Head Rules

4 Evaluation

4 Evaluation

4.2 Predictive Performance



- ▶ Using **micro- or label-based averaging** results in the best performance.
- ▶ When using micro- or label-based averaging, the algorithm tends to benefit from **predicting irrelevant labels**.
- ▶ Using the **rule-dependent evaluation** is usually a better choice than using the rule-independent evaluation.
- ▶ The algorithm always **outperforms the “Binary Relevance” approach**, if label dependencies are given.
- ▶ The algorithm **is able to compete with the original approach**. In some cases it is even ranked higher.

4 Evaluation

4.3 Characteristics of Learned Models



- ▶ Using **macro- or example-based averaging** results in less rules being learned.
- ▶ **Predicting irrelevant labels** causes the number of label-dependent rules to increase.
- ▶ The amount of **label-dependent rules** is usually similar to the original approach. Multi-label head rules can be considered as an alternative.
- ▶ Significantly more multi-label head rules are learned when using the precision metric.

Thank You!



Any Questions?