

# Classifier Chains for Multi-label Classification

Jesse Read, Bernhard Pfahringer, Geoff Holmes, Eibe Frank

University of Waikato  
New Zealand

ECML PKDD 2009, September 9, 2009. Bled, Slovenia

## Multi-label Classification

- Each instance may be associated with multiple labels
- set of instances  $X = \{x_1, \dots, x_m\}$ ; set of *predefined* labels  $L = \{l_1, \dots, l_n\}$ ; dataset  $(x_1, S_1), (x_2, S_2), \dots$  where each  $S_i \subseteq L$ .
- For example, a film can be labeled  $\{\text{romance, comedy}\}$

## Multi-label Classification

- Each instance may be associated with multiple labels
- set of instances  $X = \{x_1, \dots, x_m\}$ ; set of *predefined* labels  $L = \{l_1, \dots, l_n\}$ ; dataset  $(x_1, S_1), (x_2, S_2), \dots$  where each  $S_i \subseteq L$ .
- For example, a film can be labeled {romance, comedy}

## Applications

- Scene, Video classification
- Text classification
- Medical classification
- Biology, Genomics

## Multi-label Classification

- Each instance may be associated with multiple labels
- set of instances  $X = \{x_1, \dots, x_m\}$ ; set of *predefined* labels  $L = \{l_1, \dots, l_n\}$ ; dataset  $(x_1, S_1), (x_2, S_2), \dots$  where each  $S_i \subseteq L$ .
- For example, a film can be labeled {romance,comedy}

## Applications

- Scene, Video classification
- Text classification
- Medical classification
- Biology, Genomics

## Multi-label Issues

- label correlations: consider {romance,comedy} vs {romance,horror}
- computational complexity

- Binary relevance method (BR): binary problem for each label
  - simple, efficient
  - does not take into account label correlations

# Prior Work

- Binary relevance method (BR): binary problem for each label
  - simple, efficient
  - does not take into account label correlations
- Nearest neighbor approaches based on BR, e.g. MLkNN
- Stacking approaches, e.g. meta level stacking (MS)
- Pairwise approaches, e.g. calibrated label ranking

# Prior Work

- Binary relevance method (BR): binary problem for each label
  - simple, efficient
  - does not take into account label correlations
- Nearest neighbor approaches based on BR, e.g. MLkNN
- Stacking approaches, e.g. meta level stacking (MS)
- Pairwise approaches, e.g. calibrated label ranking
- Label powerset method: label sets are treated as single labels
  - takes into account label correlations
  - computationally complex

# Prior Work

- Binary relevance method (BR): binary problem for each label
  - simple, efficient
  - does not take into account label correlations
- Nearest neighbor approaches based on BR, e.g. MLkNN
- Stacking approaches, e.g. meta level stacking (MS)
- Pairwise approaches, e.g. calibrated label ranking
- Label powerset method: label sets are treated as single labels
  - takes into account label correlations
  - computationally complex
- RAKEL: ensembles of subsets
- EPS: ensembles of pruned sets

- Binary relevance method (BR): binary problem for each label
  - simple, efficient
  - does not take into account label correlations
- Nearest neighbor approaches based on BR, e.g. MLkNN
- Stacking approaches, e.g. meta level stacking (MS)
- Pairwise approaches, e.g. calibrated label ranking
- Label powerset method: label sets are treated as single labels
  - takes into account label correlations
  - computationally complex
- RAKEL: ensembles of subsets
- EPS: ensembles of pruned sets
- Many other methods
  - take into account label correlations
  - complex, prone to overfitting

# Binary Relevance (BR)

$L = \{\text{romance,horror,comedy,drama,action,western}\}$  ( $|L| = 6$ )

Classifiers	Classifications
$C_1 : x \rightarrow \{\text{romance,!romance}\}$	romance
$C_2 : x \rightarrow \{\text{horror,!horror}\}$	!horror
$C_3 : x \rightarrow \{\text{comedy,!comedy}\}$	comedy
$C_4 : x \rightarrow \{\text{drama,!drama}\}$	!drama
$C_5 : x \rightarrow \{\text{action,!action}\}$	!action
$C_6 : x \rightarrow \{\text{western,!western}\}$	!western
$Y \subseteq L$	$\{\text{romance,comedy}\}$

# Binary Relevance (BR)

$L = \{\text{romance,horror,comedy,drama,action,western}\}$  ( $|L| = 6$ )

Classifiers	Classifications
$C_1 : x \rightarrow \{\text{romance,!romance}\}$	romance
$C_2 : x \rightarrow \{\text{horror,!horror}\}$	!horror
$C_3 : x \rightarrow \{\text{comedy,!comedy}\}$	comedy
$C_4 : x \rightarrow \{\text{drama,!drama}\}$	!drama
$C_5 : x \rightarrow \{\text{action,!action}\}$	!action
$C_6 : x \rightarrow \{\text{western,!western}\}$	!western
$Y \subseteq L$	$\{\text{romance,comedy}\}$

- simple, intuitive
- efficient
- useful for incremental contexts
- **doesn't account for label correlations**

# Classifier Chains (CC)

$L = \{\text{romance,horror,comedy,drama,action,western}\}$  ( $|L| = 6$ )

Classifiers	Classifications
$C_1 : x \rightarrow \{\text{romance,!romance}\}$	romance
$C_2 : x \cup \text{romance} \rightarrow \{\text{horror,!horror}\}$	!horror
$C_3 : x \cup \text{romance} \cup \text{!horror} \rightarrow \{\text{comedy,!comedy}\}$	comedy
$C_4 : x \cup \text{romance} \cup \text{!horror} \cup \text{comedy} \rightarrow \{\text{drama,!drama}\}$	!drama
$C_5 : x \cup \text{romance} \cup \text{!horror} \cup \text{comedy} \cup \text{!drama} \rightarrow \dots$	!action
$C_6 : x \cup \text{romance} \cup \text{!horror} \cup \text{comedy} \cup \text{!drama} \cup \dots$	!western
$Y \subseteq L = \{\text{romance,comedy}\}$	

# Classifier Chains (CC)

$L = \{\text{romance,horror,comedy,drama,action,western}\}$  ( $|L| = 6$ )

Classifiers	Classifications
$C_1 : x \rightarrow \{\text{romance,!romance}\}$	romance
$C_2 : x \cup \text{romance} \rightarrow \{\text{horror,!horror}\}$	!horror
$C_3 : x \cup \text{romance} \cup \text{!horror} \rightarrow \{\text{comedy,!comedy}\}$	comedy
$C_4 : x \cup \text{romance} \cup \text{!horror} \cup \text{comedy} \rightarrow \{\text{drama,!drama}\}$	!drama
$C_5 : x \cup \text{romance} \cup \text{!horror} \cup \text{comedy} \cup \text{!drama} \rightarrow \dots$	!action
$C_6 : x \cup \text{romance} \cup \text{!horror} \cup \text{comedy} \cup \text{!drama} \cup \dots$	!western
$Y \subseteq L = \{\text{romance,comedy}\}$	

- similar advantages to binary relevance method
- time complexity similar in practice
- takes into account label correlations
- **how to order the chain?**

# Ensembles of Classifier Chains (ECC)

- Ensembles known for augmenting accuracy
- more label correlations learnt, without overfitting
- solves 'chain order' issue: each chain has a random order

# Ensembles of Classifier Chains (ECC)

- Ensembles known for augmenting accuracy
- more label correlations learnt, without overfitting
- solves 'chain order' issue: each chain has a random order
- For  $i \in 1 \dots m$  iterations:
  - $L' \leftarrow$  shuffle label set  $L$
  - $D' \leftarrow$  subset of training set  $D$
  - train a model  $CC_i$ ; given  $L'$  and  $D'$
- Generic vote/score/threshold method for classification:
  - collect votes from models
  - assign a score to each label
  - apply a threshold to determine relevant labels

# Ensembles of Classifier Chains (ECC)

- Ensembles known for augmenting accuracy
- more label correlations learnt, without overfitting
- solves 'chain order' issue: each chain has a random order
- For  $i \in 1 \dots m$  iterations:
  - $L' \leftarrow$  shuffle label set  $L$
  - $D' \leftarrow$  subset of training set  $D$
  - train a model  $CC_i$ ; given  $L'$  and  $D'$
- Generic vote/score/threshold method for classification:
  - collect votes from models
  - assign a score to each label
  - apply a threshold to determine relevant labels
- Can also be *applied to binary relevance method*, i.e. EBR

- WEKA-based framework
- Support Vector Machines as base classifiers
- Multi-label datasets:

	Labels $ L $	Instances $ D $
6 Standard	6 ... 103	2407 ... 6000
6 Large	22 ... 983	7395 ... 95424

- Multi-label evaluation metrics:
  - accuracy, macro F-measure (label set evaluation)
  - log loss,  $AU(\overline{PRC})$  (per-label evaluation)
  - build times, test times
- Method parameters preset to optimise *predictive performance* (ECC requires no additional parameters)
- Experiments:
  - 1 Compare Classifier Chains (CC) to the Binary Relevance method (BR) and related BR-based methods.
  - 2 Compare ECC to EBR and modern methods of proven success: RAKEL, EPS, and MLkNN

# Results 1

Comparing CC to BR and related methods SM<sup>1</sup> and MS<sup>2</sup>.

Table: Standard Datasets: Wins for each evaluation measure.

	CC	BR	SM	MS
Accuracy	<b>5</b>	0	1	0
Macro F1	<b>5</b>	0	1	0
Micro F1	<b>3</b>	1	0	2
Exact Match	<b>6</b>	0	0	0
Total wins	<b>19</b>	1	2	2

- CC's chaining technique justified over default BR
- CC outperforms other similar methods

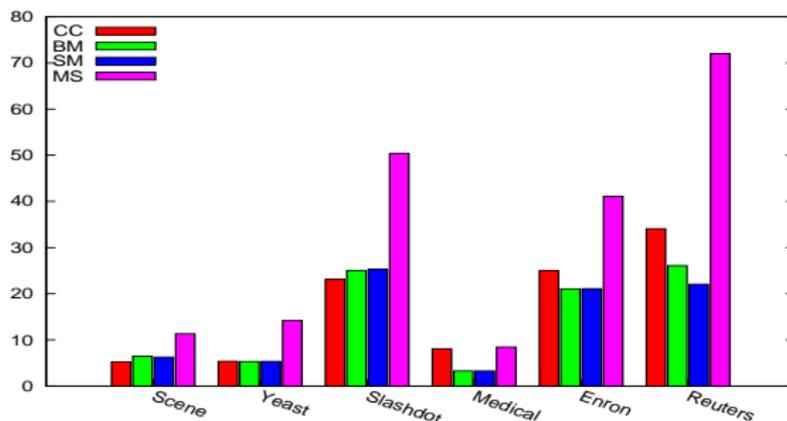
<sup>1</sup>Subset Mapping: maps output of BR to nearest (Hamming distance) known subset

<sup>2</sup>Meta Stacking: stacking the output of BR with meta classifiers

# Results 1

Comparing CC to BR and related methods SM<sup>1</sup> and MS<sup>2</sup>.

Figure: Standard Datasets: Build times (seconds).



- CC's complexity comparable to BR
  - except for special cases like *Medical* (relatively large label set)

<sup>1</sup>Subset Mapping: maps output of BR to nearest (Hamming distance) known subset

<sup>2</sup>Meta Stacking: stacking the output of BR with meta classifiers

## Results 2

Comparing ECC to EBR and methods: RAKEL<sup>3</sup>, EPS<sup>4</sup>, and MLkNN<sup>5</sup>.

Table: Standard Datasets: Wins for each evaluation measure.

	ECC	EBR	RAKEL	EPS	MLkNN
Accuracy	2	0	1	<b>3</b>	0
Macro F1	1	0	1	<b>4</b>	0
Log Loss	<b>3</b>	0	1	1	1
$AU(\overline{PRC})$	<b>3</b>	0	0	0	<b>3</b>
Total wins	<b>9</b>	0	3	8	4

- ECC best at per-label prediction (as a binary method)
- Other methods can sometimes predict better label sets
- ECC rewarded by conservative prediction (log loss)

<sup>3</sup>Tsoumakas and Vlahavas, 2007

<sup>4</sup>Read, Pfahringer, Holmes, 2008

<sup>5</sup>Zhang and Zhou, 2005

## Results 2

Comparing ECC to EBR and methods: RAKEL<sup>3</sup>, EPS<sup>4</sup>, and MLkNN<sup>5</sup>.

Table: Large Datasets: Wins for each evaluation measure.

	ECC	EBR	RAKEL <sup>†</sup>	EPS <sup>†</sup>	MLkNN
Accuracy	<b>4</b>	0	0	1	1
Macro F1	<b>3</b>	0	1	1	1
Log Loss	1	1	0	0	<b>4</b>
$AU(\overline{PRC})$	<b>4</b>	0	0	0	2
Total wins	<b>12</b>	1	1	2	8

<sup>†</sup>Note: 2 DNF for RAKEL and 1 DNF for EPS.

- Binary methods are the best choice for large datasets
- ECC best overall

<sup>3</sup>Tsoumakas and Vlahavas, 2007

<sup>4</sup>Read, Pfahringer, Holmes, 2008

<sup>5</sup>Zhang and Zhou, 2005

## Results 2

Comparing build and test times between ECC, RAKEL, and EPS.

Table: All Datasets: Method with fastest Build, Test time<sup>†</sup>.

Dataset	Build	Test	Dataset	Build	Test
Scene	EPS	RAK	OHSUMED	ECC	ECC
Yeast	ECC	ECC	TMC2007	EPS	ECC
Slashdot	RAK	RAK	Bibtex	ECC	ECC
Medical	RAK	RAK	MediaMill	ECC	ECC
Enron	EPS	ECC	IMDB	RAK	ECC
Reuters	ECC	ECC	Delicious	EPS	EPS

<sup>†</sup>EBR and MLkNN not included

- ECC's efficiency most noticeable on the larger datasets
- RAKEL most efficient on smaller datasets
- EPS can make large gains by pruning, but occasionally too much

- Ensembles of Classifier Chains
  - classifier chains improves on the binary relevance method
  - takes into account label correlations without overfitting
  - flexible, efficient
  - performs well, especially on large data sets

- Ensembles of Classifier Chains
  - classifier chains improves on the binary relevance method
  - takes into account label correlations without overfitting
  - flexible, efficient
  - performs well, especially on large data sets

*Thank you. Any questions?*