

Aggregating Independent and Dependent Models to Learn Multi-label Classifiers

E. Montañés, J.R. Quevedo, J.J. del Coz

Artificial Intelligence Center – University of Oviedo



Outline

- 1 The goal
- 2 Multi-label classification
- 3 Previous approaches
- 4 Our proposal
- 5 Experiments
- 6 Conclusions

Outline

- 1 The goal
- 2 Multi-label classification
- 3 Previous approaches
- 4 Our proposal
- 5 Experiments
- 6 Conclusions

The goal

- To study **dependence and independence on their own**
- To show that **both help** to improve the multi-label classification
- To propose **aggregating vs. stacking**
- To compare **actual labels vs. predicted ones**
- To compare **binary vs. probabilistic outputs**

Outline

- 1 The goal
- 2 Multi-label classification**
- 3 Previous approaches
- 4 Our proposal
- 5 Experiments
- 6 Conclusions

Multi-label vs. mono-label classification

Binary classification

X_1	X_2	X_3	X_4	C
3	T	1.5	-2	1
2	F	1.5	-4	0
9	T	2.3	-2	1
4	F	7.8	-1	0
3	F	1.5	-9	1

Mono-label
classification

Multi-class classification

X_1	X_2	X_3	X_4	C
3	T	1.5	-2	1
2	F	1.5	-4	2
9	T	2.3	-2	2
4	F	7.8	-1	3
3	F	1.5	-9	1

Multi-label classification

X_1	X_2	X_3	X_4	C_1	C_2	C_3	C_4	C_5
3	T	1.5	-2	1	0	0	1	0
2	F	1.5	-4	0	0	0	0	1
9	T	2.3	-2	0	1	1	1	1
4	F	7.8	-1	1	0	0	1	0
3	F	1.5	-9	1	1	0	1	0

Formal statement

Point of departure

$\mathcal{L} = \{\ell_1, \ell_2, \dots, \ell_m\}$ a set of labels

\mathcal{X} an input space

$\mathcal{Y} = \mathcal{P}(\mathcal{L}) \sim \{0, 1\}^m$ (the power set of \mathcal{L})

\Downarrow

$S = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\} \in \mathcal{X} \times \mathcal{Y}$

\approx

$\mathbf{P}(\mathbf{X}, \mathbf{Y})$

The target

Induce $\mathbf{h} : \mathcal{X} \rightarrow \mathcal{Y}$ from S

$\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_m(\mathbf{x}))$

$h_j : \mathcal{X} \rightarrow \{0, 1\}$ predicts if ℓ_j is attached to \mathbf{x}

Evaluating multi-label classification

Example-based measures

- **classification** rather than ranking
- capture correlations among labels at **example level**
- used in **stacking approaches**

Biased measures

- Jaccard index
- Precision, Recall, F_1

Other measures

- Hamming loss
- 0/1 loss

Outline

- 1 The goal
- 2 Multi-label classification
- 3 Previous approaches**
- 4 Our proposal
- 5 Experiments
- 6 Conclusions

Binary relevance

- Each label is learned **independently** of the rest
- **Linear complexity** respect to the number of labels
- Does **not** consider **label dependence**

$$\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_m(\mathbf{x}))$$
$$h_j : \mathcal{X} \longrightarrow \{0, 1\}$$

Stacking approaches

Two groups of classifiers are learned

- The **independent** ones

$$\mathbf{h}^1(\mathbf{x}) = (h_1^1(\mathbf{x}), \dots, h_m^1(\mathbf{x}))$$
$$h_j^1 : \mathcal{X} \longrightarrow \{0, 1\}$$

- The **dependent** ones

$$\mathbf{h}^2(\mathbf{x}, \mathbf{h}^1(\mathbf{x})) = (h_1^2(\mathbf{x}, \mathbf{h}^1(\mathbf{x})), \dots, h_m^2(\mathbf{x}, \mathbf{h}^1(\mathbf{x})))$$
$$h_j^2 : \mathcal{X} \times \{0, 1\}^m \longrightarrow \{0, 1\}$$



$$\mathbf{h}(\mathbf{x}) = \mathbf{h}^2(\mathbf{x}, \mathbf{h}^1(\mathbf{x}))$$

Classifier chains

- **One classifier per label**
- **Same complexity** as binary relevance
- A **chain of classifiers** is built according to certain order of the labels

$$h(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}, h_1(\mathbf{x})), h_3(\mathbf{x}, h_1(\mathbf{x}), h_2(\mathbf{x}, h_1(\mathbf{x}))), \dots)$$
$$h_j : \mathcal{X} \times \{0, 1\}^{j-1} \longrightarrow \{0, 1\}$$

Ensemble version

- Several orders of labels are ensembled
- Diminishes the effect of the label order

Probabilistic version

- The probability product rule is used
- The complexity increases considerably in testing

Other approaches

MLkNN

Instance-based learner

- Posterior (over neighbors) and prior probability based on frequency counting
- Bayes rule gives the labels' probability

Instance-Based Learning by Logistic Regression (IBLR)

Unifies instance-based learning and logistic regression

RANdom k-labELsets (RAkEL)

Ensemble of Label Power set (LP) classifiers

Other approaches

MLkNN

Instance-based learner

Instance-Based Learning by Logistic Regression (IBLR)

Unifies instance-based learning and logistic regression

- Labels of the neighbors as additional features
- Classification by logistic regression

RANdom k-labELsets (RAkEL)

Ensemble of Label Power set (LP) classifiers

Other approaches

MLkNN

Instance-based learner

Instance-Based Learning by Logistic Regression (IBLR)

Unifies instance-based learning and logistic regression

RANdom k-labELsets (RAkEL)

Ensemble of Label Power set (LP) classifiers

- It randomly selects a k -labelset Y_i from \mathcal{L} without replacement
- It learns a LP classifier of the form $\mathcal{X} \rightarrow \mathcal{P}(Y_i)$
- A voting process determines the final classification

Outline

- 1 The goal
- 2 Multi-label classification
- 3 Previous approaches
- 4 Our proposal**
- 5 Experiments
- 6 Conclusions

Aggregating Independent and Dependent classifiers (AID)

Our hypothesis

Both approaches are **not exclusive**, but **complementary**

- The **independent** ones

$$\mathbf{h}^1(\mathbf{x}) = (h_1^1(\mathbf{x}), \dots, h_m^1(\mathbf{x}))$$

$$h_j^1 : \mathcal{X} \longrightarrow \{0, 1\}$$

- The **dependent** ones

$$\mathbf{h}^2(\mathbf{x}, \mathbf{y}) = (h_1^2(\mathbf{x}, y_2, \dots, y_m), \dots, h_m^2(\mathbf{x}, y_1, \dots, y_{m-1}))$$

$$h_j^2 : \mathcal{X} \times \{0, 1\}^{m-1} \longrightarrow \{0, 1\}$$



$$\mathbf{h}(\mathbf{x}) = \oplus((h_1^1(\mathbf{x}), \dots, h_m^1(\mathbf{x})),$$

$$(h_1^2(\mathbf{x}, h_1^1(\mathbf{x}), \dots, h_m^1(\mathbf{x})), \dots, h_m^2(\mathbf{x}, h_1^1(\mathbf{x}), \dots, h_{m-1}^1(\mathbf{x}))))$$

Comparing AID with other methods (I)

With regard to ...

... binary relevance

- It **also considers correlations** among labels

... stacking approaches

- The outputs of independent classifiers are **additionally employed** to decide the predicted labels
- **Actual labels** rather than predicted labels (more reliable information)

Comparing AID with other methods (II)

With regard to ...

... chain classifiers

- **Free of in-chain dependence**
- **Richer estimations** since all correlations are considered
- Although it **only** offers **greedy approximations** of the entire joint distribution

... MLkNN, IBLR & RAKEL

- **Interpretability** of different kinds of labels predicted
 - Those coming just from the description of the examples
 - Those coming from other labels

About the complexity

	BR	STA	CC	PCC	ECC	EPCC	AID
Models	m	$2m$	m	m	Nm	Nm	$2m$
Testing complexity	Linear	Linear	Linear	Exp.	Linear	Exp.	Linear

Outline

- 1 The goal
- 2 Multi-label classification
- 3 Previous approaches
- 4 Our proposal
- 5 Experiments**
- 6 Conclusions

Settings

The learning process

- For binary relevance, CC, stacking & AID
 - **logistic regression** as binary base learner
 - A **grid search for C** over $\{10^p \mid p \in [-3, \dots, 3]\}$ optimizing the accuracy through a balanced 2-fold cross validation repeated 5 times
- **Default parameters** for MLkNN, IBLR & RAKEL

The evaluation

Using a 10-fold cross validation, we estimate

- Jaccard index
- Precision, Recall, & F_1
- Hamming loss & 0/1 loss

AID vs. Stacking

Average ranks over all data sets

	BR	AID	STA ^y	AID ^{y'}	STA	AID ^p	STA ^p
Precision	4.45	2.45	3.59	4.91	4.82	3.09	4.68
Recall	6.18	1.45	3.00	3.55	4.27	3.55	6.00
F1	5.82	1.36	2.95	4.45	4.91	3.18	5.32
Jaccard	5.64	1.45	2.95	4.64	4.91	3.18	5.23
Hamming	2.27	4.91	5.00	4.59	4.41	2.77	4.05
0/1loss	4.91	2.64	2.73	5.05	4.77	3.27	4.64

- 1 **STA** or **AID**?
- 2 **Actual label** data or **predictions** of independent models?
- 3 In the testing phase, **binary** or **probabilistic** features?
- 4 To **aggregate** or to **stack**?

AID vs. other methods

Average ranks over all data sets

	BR	MLkNN	IBLR	RAKEL	ECC	AID	AID ^P
Precision	4.05	5.55	4.55	3.45	4.14	3.05	3.23
Recall	5.14	6.45	4.55	2.73	4.86	1.41	2.86
F1	4.95	6.36	4.55	2.82	4.86	1.59	2.86
Jaccard	4.77	6.27	4.00	2.91	4.86	2.05	3.14
Hamming	3.05	4.36	4.18	4.36	3.00	5.18	3.86
0/1loss	4.50	5.14	3.73	3.55	4.23	2.91	3.95

- **AID is the best** for all measures except for Hamming loss
- In **Hamming loss**, **ECC is the best** and AID is the worst
- **AID^P** is quite **steady** for all measures

Outline

- 1 The goal
- 2 Multi-label classification
- 3 Previous approaches
- 4 Our proposal
- 5 Experiments
- 6 Conclusions**

Conclusions

- **Interpretability** of two kinds of labels
- **Actual labels** better than **predicted ones**
- **Aggregating** better than **stacking**
- AID exhibits **competitive results**, but not for Hamming loss
- AID has **linear complexity** in both training and testing stages