

# Multi-label and Multi-target Learning

## Applications, Challenges, and Models

Jesse Read



January 5, 2021  
Zoom

# Outline

- 1 Multi-Label and Multi-Target Learning
- 2 Algorithm Adaptations
- 3 Classifier Chains
- 4 Regressor Chains
- 5 Modern Multi-Output Topics
- 6 Summary

# Multi-Label and Multi-Target Learning

- 1 Multi-Label and Multi-Target Learning
- 2 Algorithm Adaptations
- 3 Classifier Chains
- 4 Regressor Chains
- 5 Modern Multi-Output Topics
- 6 Summary

## Classification Multi-label

$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
1	0.1	3	1	0	0	1	1	0
0	0.9	1	0	1	1	0	0	0
0	0.0	1	1	0	0	1	0	0
1	0.8	2	0	1	1	0	0	1
1	0.0	2	0	1	0	0	0	1
0	0.0	3	1	1	?	?	?	?

For input  $x$  we get a vector output

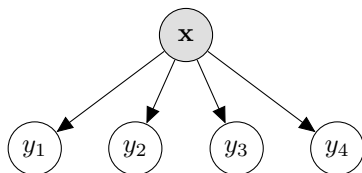
$$\hat{y} = \mathbf{h}(x) = \mathbf{h}(\underbrace{[x_1, \dots, x_d]}_{\text{inputs}}) = \underbrace{[y_1, \dots, y_L]}_{\text{outputs}}$$

N.B. Not multi-class, but multi-class **multi-label**!

$y = [0, 1, 1, 0] \Leftrightarrow$  labels  $\{2, 3\}$  are relevant to corresponding  $x$ .

## Reduction #1 (to binary): Binary Relevance Method

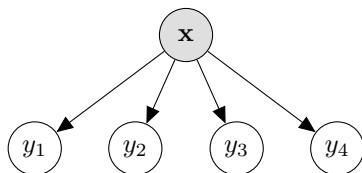
$\mathbf{X}$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$\mathbf{x}^{(1)}$	0	1	1	0
$\mathbf{x}^{(2)}$	1	0	0	0
$\mathbf{x}^{(3)}$	0	1	0	0
$\mathbf{x}^{(4)}$	1	0	0	1
$\mathbf{x}^{(5)}$	0	0	0	1
$\tilde{\mathbf{x}}$	?	?	?	?



The **binary relevance method** (BR transformation) = *one binary classifier trained for each label*, i.e., **independent models**.

## Reduction #1 (to binary): Binary Relevance Method

$\mathbf{X}$	$Y_1$	$\mathbf{X}$	$Y_2$	$\mathbf{X}$	$Y_3$	$\mathbf{X}$	$Y_4$
$\mathbf{x}^{(1)}$	0	$\mathbf{x}^{(1)}$	1	$\mathbf{x}^{(1)}$	0	$\mathbf{x}^{(1)}$	1
$\mathbf{x}^{(2)}$	1	$\mathbf{x}^{(2)}$	0	$\mathbf{x}^{(2)}$	1	$\mathbf{x}^{(2)}$	0
$\mathbf{x}^{(3)}$	0	$\mathbf{x}^{(3)}$	1	$\mathbf{x}^{(3)}$	0	$\mathbf{x}^{(3)}$	1
$\mathbf{x}^{(4)}$	1	$\mathbf{x}^{(4)}$	0	$\mathbf{x}^{(4)}$	1	$\mathbf{x}^{(4)}$	0
$\mathbf{x}^{(5)}$	0	$\mathbf{x}^{(5)}$	0	$\mathbf{x}^{(5)}$	0	$\mathbf{x}^{(5)}$	0
$\tilde{\mathbf{x}}$	?	$\tilde{\mathbf{x}}$	?	$\tilde{\mathbf{x}}$	?	$\tilde{\mathbf{x}}$	?



The **binary relevance method** (BR transformation) = *one binary classifier trained for each label*, i.e., **independent models**.

## Reduction #2 (to multi-class): Label Powerset Method

$\mathbf{X}$	$\mathbf{Y}$
$\mathbf{x}^{(1)}$	0 1 1 0
$\mathbf{x}^{(2)}$	1 0 0 0
$\mathbf{x}^{(3)}$	0 1 0 0
$\mathbf{x}^{(4)}$	1 0 0 1
$\mathbf{x}^{(5)}$	0 0 0 1
$\tilde{\mathbf{x}}$	?



The **label powerset method** (LP transformation) = a *single target multi-class classifier*. Labels are modeled together, mais ...

- Overfitting
- $y \in \{0, 1\}^L$ .

# A Brief Timeline of Multi-label Learning in Academia






- < 2000s : Use (baseline) reduction #1 (BR), or #2 (LP)
- ... 2010 :
  - We beat BR (using label dependence)!
  - Many applications!
- ... 2015 :
  - We beat the methods that beat BR (using label dependence in a more sophisticated way)!
  - Wait – what are we doing? And why?
- ... 2020 :
  - Models get deep, deeper, ... ; (CNNs, LSTM, ...)
  - Problems get big, bigger, ...
  - Do we actually need label dependence models? (BR seems to work well!)
- Recently/Currently:
  - New tasks and applications: partial labels, weak labels, label ambiguity, imprecise prediction/with abstention, ...
  - Models: neural, graph embeddings, adversarial, attention, ...



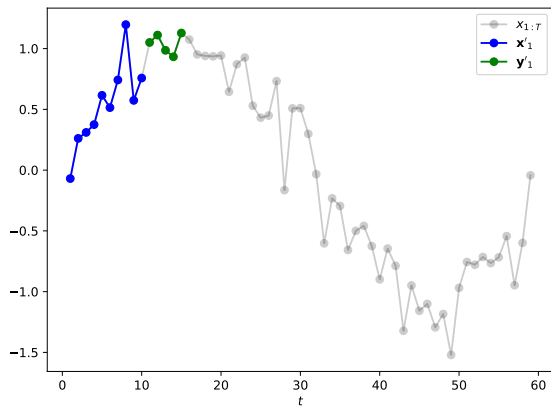
# Example Application: Multi-Label Classification

Input	Beach	Sunset	Foliage	Urban
	1	0	1	0
	0	1	0	0
	0	1	0	1
	0	1	1	0
	0	0	1	1
	?	?	?	?

# Missing-data Imputation / Recommender Systems

	Film 2 $X_2$	Film 3 $X_4$	Book 1 $X_1$	Book 2 $X_3$	Song 5 $X_5$
	0	0	1	1	0
	1	1	?	0	?
	0	0	1	0	0
	1	1	?	0	1
	0	0	0	?	?
	1	0	?	1	?

# Time Series Forecasting / Trajectory Prediction



(including multi-dimensional time series).

# Algorithm Adaptations

- 1 Multi-Label and Multi-Target Learning
- 2 Algorithm Adaptations**
- 3 Classifier Chains
- 4 Regressor Chains
- 5 Modern Multi-Output Topics
- 6 Summary

Support for multilabel / multioutput in SCIKITLEARN.

Adapted methods:

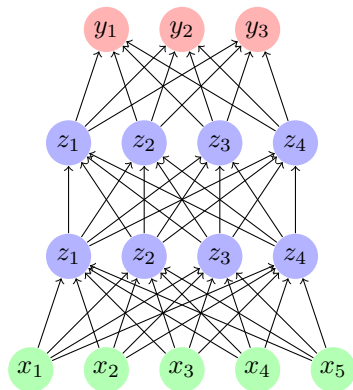
- `sklearn.tree.DecisionTreeClassifier`
- `sklearn.tree.ExtraTreeClassifier`
- `sklearn.ensemble.ExtraTreesClassifier`
- `sklearn.neighbors.KNeighborsClassifier`
- `sklearn.neural_network.MLPClassifier`
- `sklearn.neighbors.RadiusNeighborsClassifier`
- `sklearn.ensemble.RandomForestClassifier`
- `sklearn.linear_model.RidgeClassifierCV`

i.e., [Decision Trees](#), [Nearest-Neighbours](#), [Neural Networks](#).

Classifier-agnostic (transformation/reduction methods):

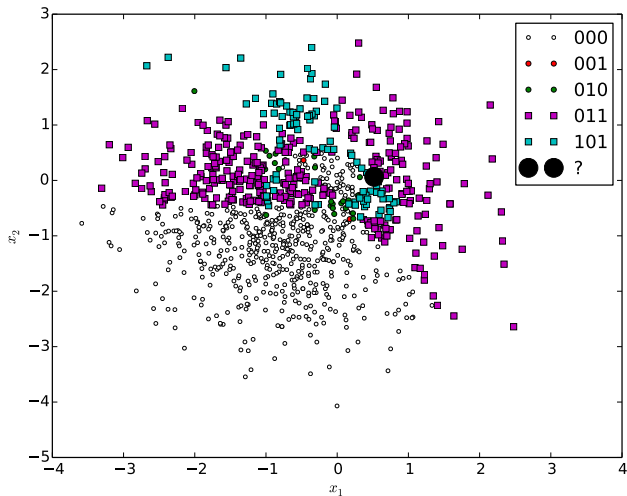
- `sklearn.multiclass.OneVsRestClassifier` ← Baseline BR
- `sklearn.multioutput.ClassifierChain` ← Coming to this soon

# Neural Networks

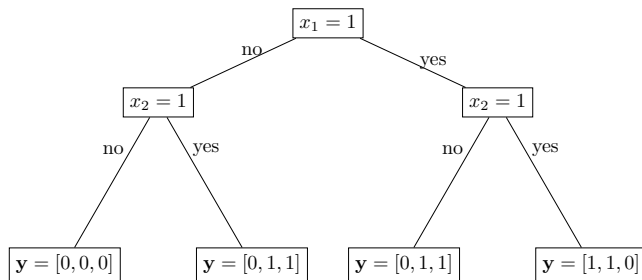


(we're coming back to this later ...)

# $k$ -Nearest Neighbours ( $k$ NN)



# Decision Tree Classifiers



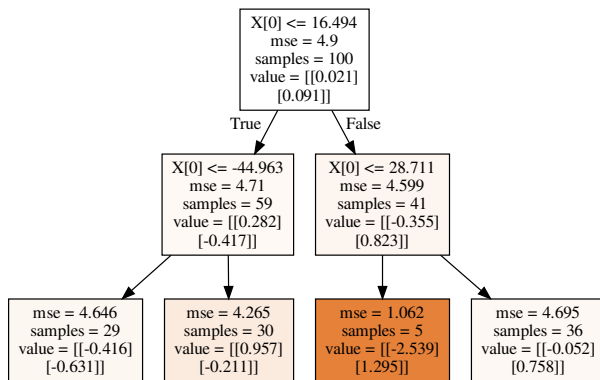
Using **multi-label entropy**,

$$H_{\text{ML}}(S) = - \sum_{j=1}^L \sum_{k \in \{0,1\}} P(y_j = k) \log_2 P(y_j = k)$$

Typical advantages/disadvantages of decision trees apply.



# Decision Tree Regression



Using redefined impurity measure:

$$\sum_{i=1}^N \sum_{j=1}^L (y_{ij} - \bar{y}_j)^2$$

where  $\bar{y}_j$  is the mean of  $Y_j$  in the node.

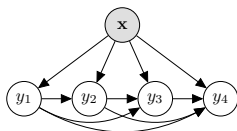
# Classifier Chains

- 1 Multi-Label and Multi-Target Learning
- 2 Algorithm Adaptations
- 3 Classifier Chains**
- 4 Regressor Chains
- 5 Modern Multi-Output Topics
- 6 Summary

# (Greedy) Classifier Chains

A chain (**structure**) over the output variables;

- Cascaded prediction across a chain/graph
- Motivation: Model label dependence with structure

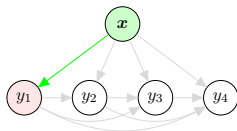


$X$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0
$\tilde{x}$	$\hat{y}_1$	$\hat{y}_2$	$\hat{y}_3$	$\hat{y}_4$

# (Greedy) Classifier Chains

A chain (**structure**) over the output variables;

- Cascaded prediction across a chain/graph
- Motivation: Model label dependence with structure

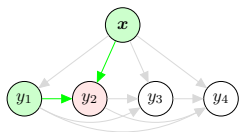


$X$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0
$\tilde{x}$	$\hat{y}_1$			

# (Greedy) Classifier Chains

A chain (**structure**) over the output variables;

- Cascaded prediction across a chain/graph
- Motivation: Model label dependence with structure

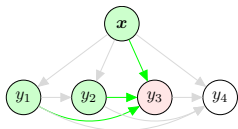


$X$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	0	1	1	1
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	1
$x^{(4)}$	1	0	0	0
$x^{(5)}$	0	0	0	0
$\tilde{x}$	$\hat{y}_1$	$\hat{y}_2$		

# (Greedy) Classifier Chains

A chain (**structure**) over the output variables;

- Cascaded prediction across a chain/graph
- Motivation: Model label dependence with structure



$\mathbf{x}$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$\mathbf{x}^{(1)}$	0	1	1	1
$\mathbf{x}^{(2)}$	1	0	0	0
$\mathbf{x}^{(3)}$	0	1	0	1
$\mathbf{x}^{(4)}$	1	0	0	0
$\mathbf{x}^{(5)}$	0	0	0	0
$\tilde{\mathbf{x}}$	$\hat{y}_1$	$\hat{y}_2$	$\hat{y}_3$	

For example,

$$\hat{y}_3 = h_3(\mathbf{x}, \hat{y}_1, \hat{y}_2)$$

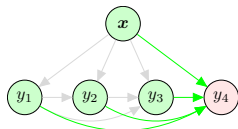
Use training data to fit **base classifier** (or regressor)  $h_3$  (e.g., decision tree, logistic regression, ...).

Inference:  $\hat{y}_1, \hat{y}_2, \dots$  are **greedy predictions** from  $h_1, h_2, \dots$

# (Greedy) Classifier Chains

A chain (**structure**) over the output variables;

- Cascaded prediction across a chain/graph
- Motivation: Model label dependence with structure



$\mathbf{x}$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$\mathbf{x}^{(1)}$	0	1	1	1
$\mathbf{x}^{(2)}$	1	0	0	0
$\mathbf{x}^{(3)}$	0	1	0	1
$\mathbf{x}^{(4)}$	1	0	0	0
$\mathbf{x}^{(5)}$	0	0	0	0
$\tilde{\mathbf{x}}$	$\hat{y}_1$	$\hat{y}_2$	$\hat{y}_3$	$\hat{y}_4$

For example,

$$\hat{y}_3 = h_3(\mathbf{x}, \hat{y}_1, \hat{y}_2)$$

Use training data to fit **base classifier** (or regressor)  $h_3$  (e.g., decision tree, logistic regression, ...).

Inference:  $\hat{y}_1, \hat{y}_2, \dots$  are **greedy predictions** from  $h_1, h_2, \dots$

# Multi-label Inference: What are we doing?; Why?

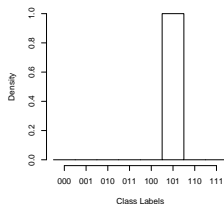
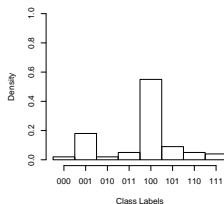
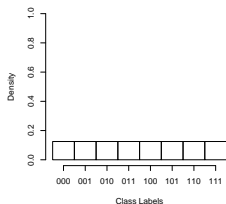
- Hamming loss (decomposable):

$$\ell_H([1, 0, 0], [1, 0, 1]) = 1/3$$

- 0/1 loss (non-decomposable):

$$\ell_{0/1}([1, 0, 0], [1, 0, 1]) = 1$$

The minimizer is not (necessarily) the same!



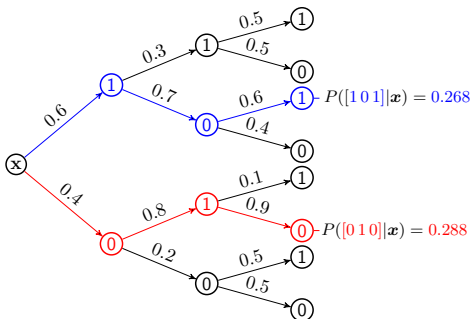
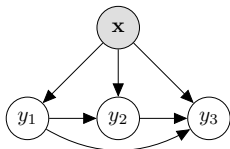
Under total uncertainty, left is optimal for Hamming loss, right for 0/1-loss



# Probabilistic Classifier Chains

We can plug in predictions  $\hat{y}$  (*greedy*); or any  $y_1, \dots, y_j = \mathbf{y} \in \{0, 1\}^j$ ; to minimise 0/1-loss :

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \{0, 1\}^3} P(\mathbf{y} | \mathbf{x}) \quad \text{where} \quad P(\mathbf{y} | \mathbf{x}) = \prod_{j=1}^3 P(y_j | y_1, \dots, y_{j-1}, \mathbf{x})$$



i.e., a path through the probability tree; e.g.,  $p([0 1 0] | \mathbf{x}) = 0.288$

# Motivation for Structure in Multi-Target Learning ?

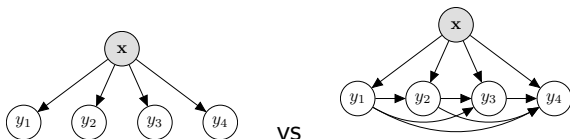
Common argument: Because [label dependence](#)!

# Motivation for Structure in Multi-Target Learning ?

Common argument: Because **label dependence!**

Yes for 0/1 loss.

Hamming loss & other **decomposable metrics**  $\Rightarrow$  **classifier chains are useless?** (and other structure/dependence-based models).



Risk minimization says that **yes** (chains are useless) under Hamming loss,

but empirical results show classifier chains performing well under most metrics (incl. Hamming loss)!

i.e., **structure is generally effective?** – then why?

Other reasons for modelling targets together (other than excuse 'because label dependence [to minimise 0/1-loss]', etc.):

- Connectivity = **efficiency** (sometimes)
- Connectivity = **interpretation** (sometimes)
- Connectivity = **power** (it's why deep nets or stacking works<sup>1</sup>)
- In same cases the minimizer *is* the same (e.g., low-noise scenarios / prediction is easy) = **surrogates** work well.
- Multiple tasks = **regularization** (regularization is good)

---

<sup>1</sup>Different reason if you *train* on  $y_j^{(i)}$  or  $\hat{y}_j^{(i)}$  as inputs

Waegeman, Dembczyński, and Hüllermeier, "Multi-target prediction: a unifying view on problems and methods", 2019; Read et al., "Classifier Chains: A Review and Perspectives", 2021

Other reasons for modelling targets together (other than excuse 'because label dependence [to minimise 0/1-loss]', etc.):

- Connectivity = **efficiency** (sometimes)
- Connectivity = **interpretation** (sometimes)
- Connectivity = **power** (it's why deep nets or stacking works<sup>1</sup>)
- In same cases the minimizer *is* the same (e.g., low-noise scenarios / prediction is easy) = **surrogates** work well.
- Multiple tasks = **regularization** (regularization is good)

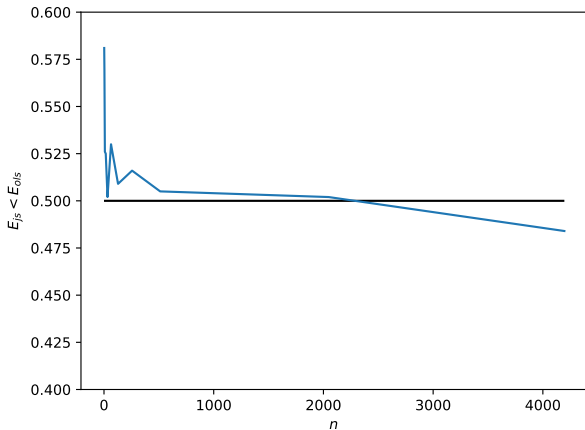
## James-Stein Estimator

Joint-target regularization is beneficial *even if targets are intrinsically independent*.

---

<sup>1</sup>Different reason if you *train* on  $y_j^{(i)}$  or  $\hat{y}_j^{(i)}$  as inputs

Waegeman, Dembczyński, and Hüllermeier, "Multi-target prediction: a unifying view on problems and methods", 2019; Read et al., "Classifier Chains: A Review and Perspectives", 2021



Advantages quickly fade as  $n \gg 0$ .

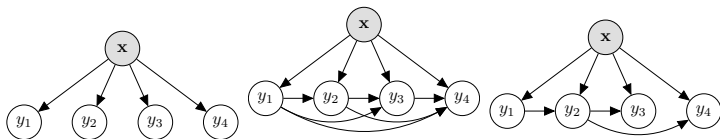
Explains reemergence of independent models vs structure debate. . .

# Chains vs Other Approaches

	$X_1$	$X_2$	$X_3$	$Y_2$	$X_1$	$X_2$	$X_3$	$Y_2$
Basis expansion	$x$	$\phi_1$	$\phi_2$	$y_2$	$\tilde{x}$	$\phi_1$	$\phi_2$	$\hat{y}_2$
Classifier chain	$x$	$y_1$		$y_2$	$\tilde{x}$	$\hat{y}_1$		$\hat{y}_2$
Stacking	$x$	$\hat{y}_1$	$\hat{y}_2$	$y_2$	$\tilde{x}$	$\hat{y}_1^{[1]}$	$\hat{y}_2^{[1]}$	$\hat{y}_2^{[2]}$
Neural network	$x$			$y_2$	$\tilde{x}$	$\hat{z}_1$	$\hat{z}_2$	$\hat{y}_2$

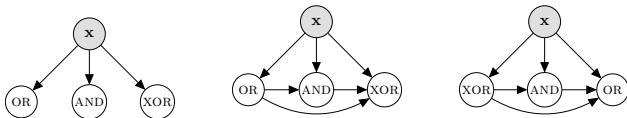
Training (left) vs Testing (right) – wrt  $\hat{y}_2 | \tilde{x}$

# Lessons on Finding a Good Structure



- Different chain *orders* are **equivalent** in theory *if* you have  $P$
- **Dependence** is not the only component to consider (and your hierarchy is probably not better than a random one)
- Weaker **base learner**/smaller training set  $\Rightarrow$  more connectivity
- Weaker (greedy) **inference** = choose more carefully
- Best structure for **loss**  $\ell_a$ , may not be the best for loss  $\ell_b$
- Best structure for  $x$  not the best for  $\tilde{x}$  (you can use a population; do **dynamic selection**)
- **Search**: Space is huge, but **local optimum can be good**

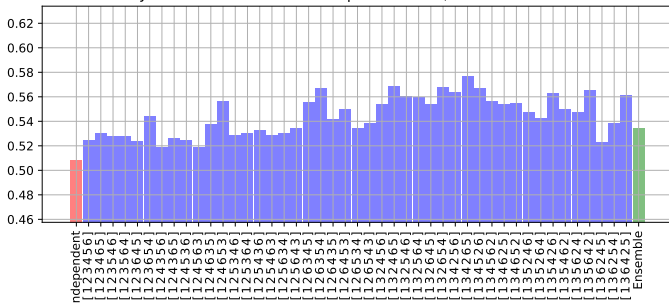




Metric	BR (left)	$CC_1$ (mid)	$CC_2^\dagger$ (right)
HAMMING LOSS	0.17	0	0
0/1 LOSS	0.50	0	0

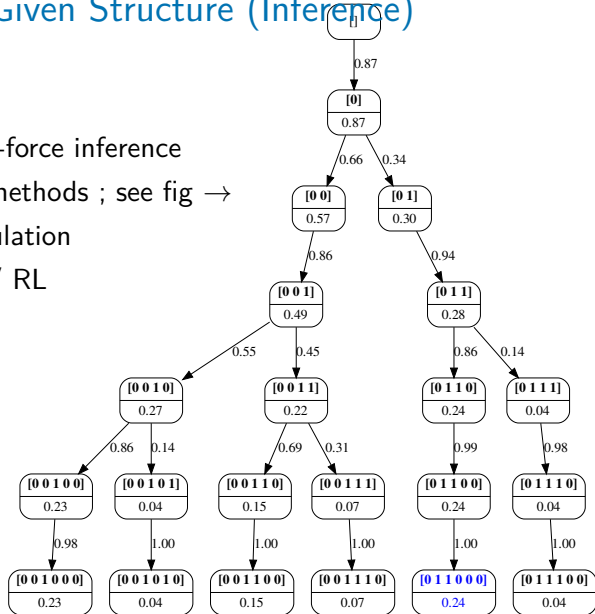
Where  $x \in \{0, 1\}^2$ ; Base-model = Logistic regression;  $\dagger$  But not greedy inference!

Jaccard score from 45 chain permutations; 'emotions' data.



# How to Traverse a Given Structure (Inference)

- Greedy vs Brute-force inference
- AI Tree-search methods ; see fig →
- Dynamic / Population
- Generic agents / RL

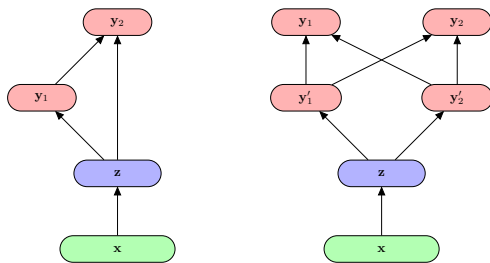


# Chains vs Deep Learning?

Can chains compete against deep architectures?

Some years ago: Yes! Now:

- Maybe **wrt accuracy**, but only on relatively smaller datasets
- Maybe **wrt explainability**:
  - decision trees (etc.) as base model
  - connection among outputs
- Classifier chains *are* deep architectures; can be combined:



A combination of chaining and deep-neural architectures

Read and Hollmén, *Multi-label Classification using Labels as Hidden Nodes*, 2017, Cisse, Al-Shedivat, and Bengio, "ADIOS: Architectures Deep In Output Space", 2016,, "Learning Deep Latent Spaces for Multi-Label Classification", 2017

# Regressor Chains

- 1 Multi-Label and Multi-Target Learning
- 2 Algorithm Adaptations
- 3 Classifier Chains
- 4 Regressor Chains**
- 5 Modern Multi-Output Topics
- 6 Summary

## Regression Multi-Cibles

$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$Y_1$	$Y_2$	$Y_3$
2.12	1.217	-0.675	-0.451	0.342	37.00	25	0.88
-0.717	-0.826	0.064	-0.259	-0.717	-22.88	22	0.22
1.374	0.95	0.175	-0.006	-0.522	19.21	12	0.25
1.392	-0.496	-2.441	-1.012	0.268	88.23	11	0.77
1.591	0.208	0.17	-0.207	1.686	?	?	?

As in classification: We can use independent models ...

As in classification: We can put variables into a chain (**regressor chains**);

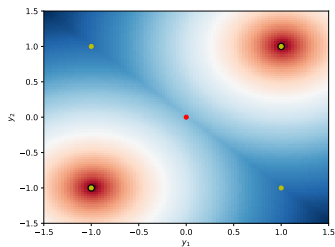
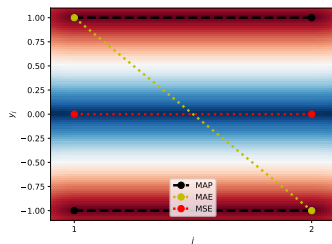
But it's probably useless to do that, because

- Our loss **metric has changed** (probably MSE, MAE, ...)
- We probably chose linear regression; **lost our non-linearity**

# Regressor Chains

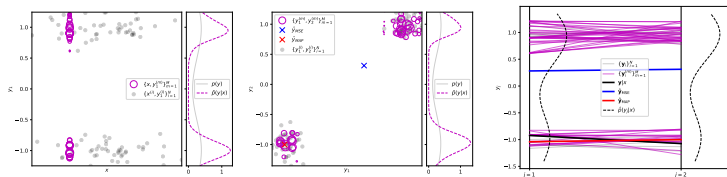
We can . . .

- Work very hard on structure
- Look at other loss functions (other than MSE, MAE, . . .); such as modal estimates.



Two equally un/likely trajectories (given  $\mathbf{x}$ ) over  $y_1 \in \mathbb{R}$ ,  $y_2 \in \mathbb{R}$  : MSE vs MAE vs MAP approx.

## Sequential Monte Carlo Methods for tracking modal predictions (i.e., trajectories)<sup>2</sup>:



Related approach : Multi-target regression via output space quantization<sup>3</sup>

Other options: Multi-target decision trees<sup>4</sup> and ensembles; and deep learning.

<sup>2</sup>Read and Martino, Neurocomputing 2020

<sup>3</sup>Spyromitros-Xioufis, Sechidis, and Vlahavas, ArXiv preprint 2020

<sup>4</sup>Stepišnik and Kocev, ArXiv preprint 2020

# Modern Multi-Output Topics

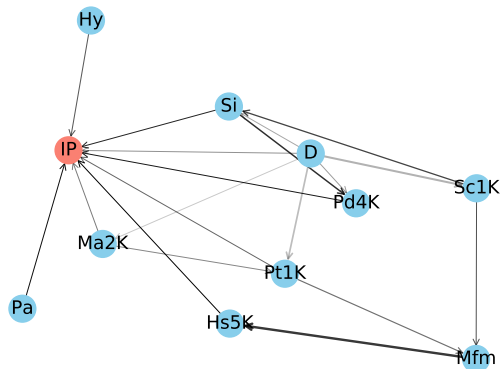
- 1 Multi-Label and Multi-Target Learning
- 2 Algorithm Adaptations
- 3 Classifier Chains
- 4 Regressor Chains
- 5 Modern Multi-Output Topics**
- 6 Summary



# Open Questions

**Loss metrics:** which loss is more appropriate, and given some loss, how to minimize it in a principled way.

**Interpretation/Explainability:** What does label dependence mean wrt the data?



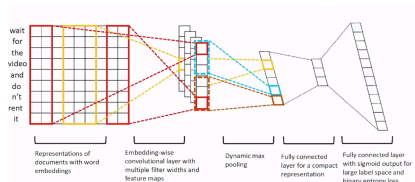
# Trends

Bigger / deeper.

Larger target spaces (4,000—3,000,000 labels), i.e., ‘**extreme multi-label classification**’

Wider range of applications

- tagging
- video recommendation
- ...



Intersection with existing areas (deep learning, multi-task, transfer learning, etc.).

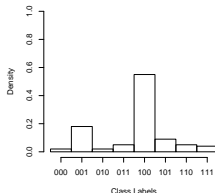
# Trends


- Import extra problems (already known in wider machine learning) : streams, semi-supervised learning, time series classification, etc.
- Learning with **partial labels** (noisy annotators), **weak labels** (lazy annotators), label ambiguity and **imprecise classification** (messy ground truth/partial abstention).



The set of candidate labels

building window  
sky street  
people car  
tree



training image	GroundTruth	Tagged Labels
	people clothing cloud sky water sea nature	people clothing sky

# Summary

- 1 Multi-Label and Multi-Target Learning
- 2 Algorithm Adaptations
- 3 Classifier Chains
- 4 Regressor Chains
- 5 Modern Multi-Output Topics
- 6 Summary

## Summary: Multi-target Learning and Prediction

- Special cases: Multi-label classification; multi-target regression
- A look at methods through the lens of **classifier chains** and **regressor chains** (decision trees and neural networks as alternative/overlap)
- Main question: *If*, and *why*, and *how* to use **structure**
- Not answered (in detail): how to *find* that structure. There is no single optimal structure, and there is more to multi-target learning than 'modelling label dependencies': consider metrics, efficiency, base models, interpretation, . . . .
- Multi-target problems getting bigger and more diverse; intersecting with other areas
- Many applications; More theoretical research needed

# Multi-label and Multi-target Learning

## Applications, Challenges, and Models

Jesse Read



Thank you !

<http://www.lix.polytechnique.fr/~jread/>