

Apprentissage Multi-cibles

Applications, Défis, et Modèles

Jesse Read



December 3, 2020

Crédite Agricole, Mont Rouge via CiscoWebex

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	?	?	?	?

Pour une entrée x on obtient la prediction d'un vecteur

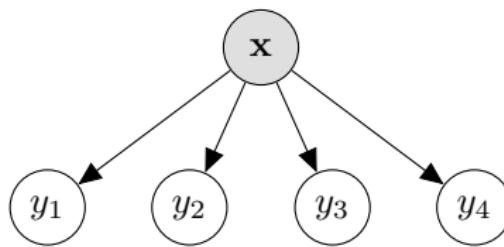
$$\hat{\mathbf{y}} = \mathbf{h}(x) = \mathbf{h}(\underbrace{[x_1, \dots, x_d]}_{\text{inputs}}) = [\underbrace{y_1, \dots, y_L}_{\text{outputs}}]$$

N.B. Not multi-class!

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

Reduction #1 (to binary): Binary Relevance Method

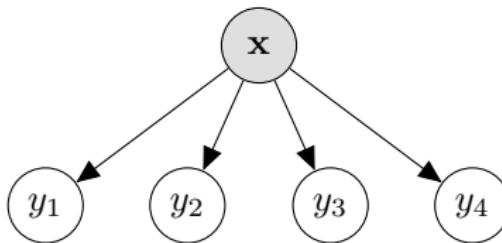
\mathbf{X}	Y_1	Y_2	Y_3	Y_4
$x^{(1)}$	0	1	1	0
$x^{(2)}$	1	0	0	0
$x^{(3)}$	0	1	0	0
$x^{(4)}$	1	0	0	1
$x^{(5)}$	0	0	0	1
\tilde{x}	?	?	?	?



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

X	Y
$x^{(1)}$	0 1 1 0
$x^{(2)}$	1 0 0 0
$x^{(3)}$	0 1 0 0
$x^{(4)}$	1 0 0 1
$x^{(5)}$	0 0 0 1
\tilde{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)!
 - **On fait quoi exactement ? Et pourquoi ?**
- ... 2020 :
 - Models get deep, deeper, ... ; (CNNs, LSTM, ...)
 - Problems get big, bigger, ...
 - **Do we really need label dependence models ? (BR seems OK)**
- Aujourd'hui :
 - **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
	?	?	?	?

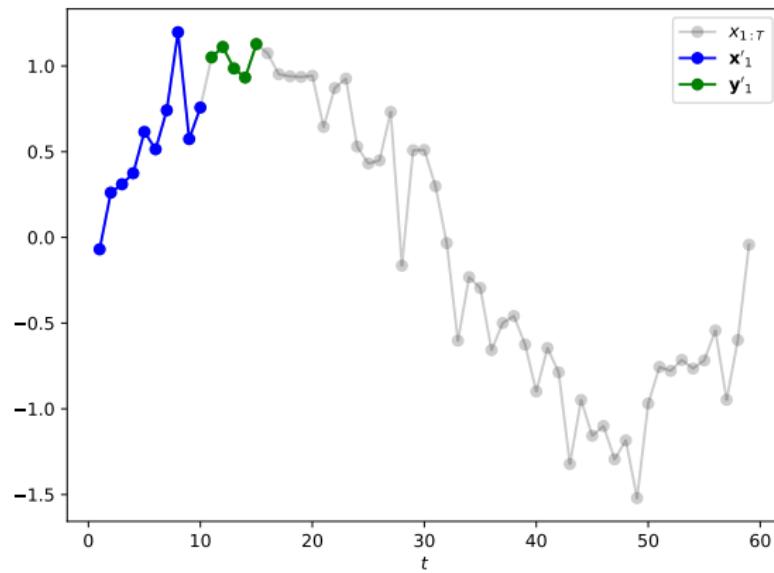
Missing-data Imputation / Recommender Systems

	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

Prévision : séries temporelles



(dont séries multidimensionnelles).

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 multilabel / multioutput in SCIKITLEARN:

Des méthodes adaptés à multi-output :

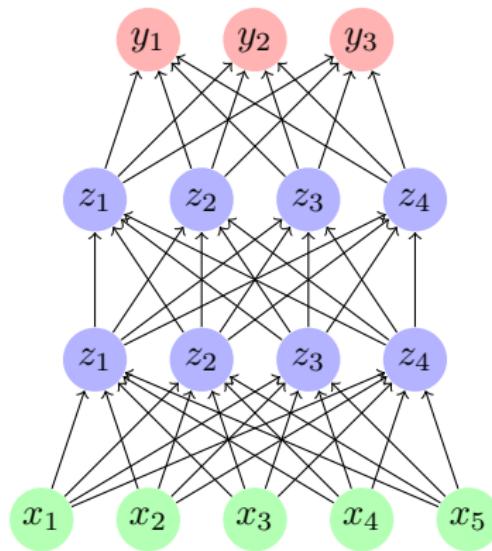
- `sklearn.tree.DecisionTreeClassifier`
- `sklearn.tree.ExtraTreeClassifier`
- `sklearn.ensemble.ExtraTreesClassifier`
- `sklearn.neighbors.KNeighborsClassifier`
- `sklearn.neural_network.MLPClassifier`
- `sklearn.neighbors.RadiusNeighborsClassifier`
- `sklearn.ensemble.RandomForestClassifier`
- `sklearn.linear_model.RidgeClassifierCV`

i.e., `Decision Trees`, `Nearest-Neighbours`, `Neural Networks`.

Les méthodes classifieur-agnostic :

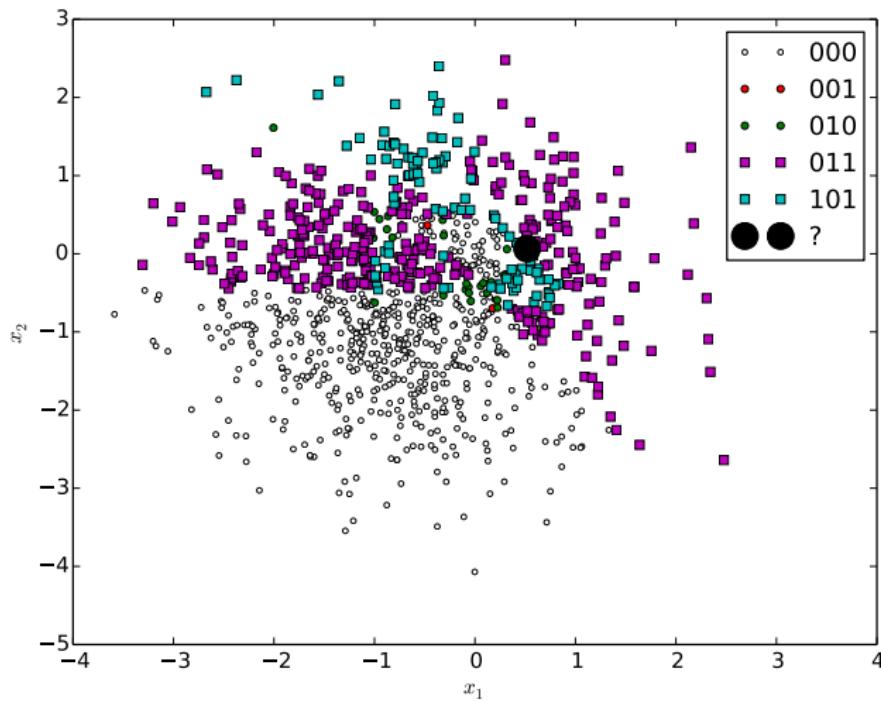
- `sklearn.multiclass.OneVsRestClassifier` ← BR
- `sklearn.multioutput.ClassifierChain` ← Coming to this soon

Réseaux de neurones

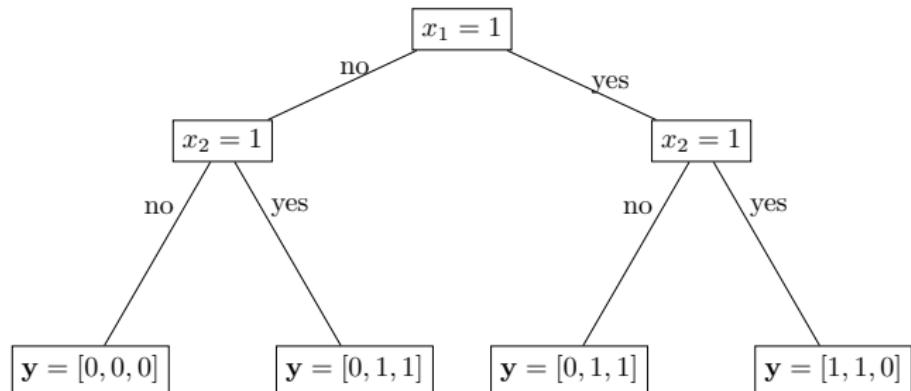


(on revient à ce sujet plus tard !)

k -Nearest Neighbours (k NN)



Arbres de décision – Classification

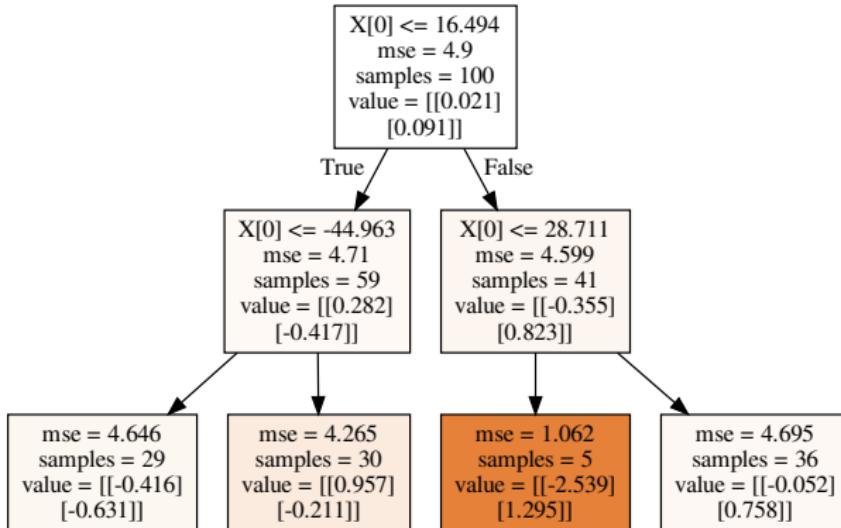


Using multi-label entropy,

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

Typical advantages/disadvantages of decision trees apply.

Arbres de décision – Regression



Using redefined impurity measure:

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

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

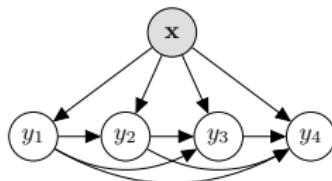
Classifier Chains

- 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

Une ‘chaîne’ (structure) entre les variables de sortie ;

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



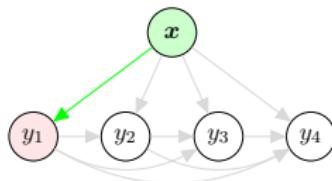
\mathbf{X}	Y_1	Y_2	Y_3	Y_4
$\mathbf{x}^{(1)}$	0	1	1	1
$\mathbf{x}^{(2)}$	1	0	0	0
$\mathbf{x}^{(3)}$	0	1	0	1
$\mathbf{x}^{(4)}$	1	0	0	0
$\mathbf{x}^{(5)}$	0	0	0	0

$\tilde{\mathbf{x}}$	\hat{y}_1	\hat{y}_2	\hat{y}_3	\hat{y}_4

(Greedy) Classifier Chains

Une ‘chaîne’ (structure) entre les variables de sortie ;

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



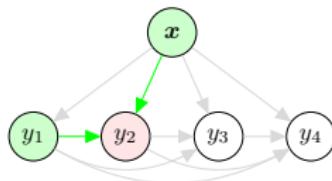
X	Y ₁	Y ₂	Y ₃	Y ₄
$\tilde{x}^{(1)}$	0	1	1	1
$\tilde{x}^{(2)}$	1	0	0	0
$\tilde{x}^{(3)}$	0	1	0	1
$\tilde{x}^{(4)}$	1	0	0	0
$\tilde{x}^{(5)}$	0	0	0	0

\tilde{x}	\hat{y}_1
-------------	-------------

(Greedy) Classifier Chains

Une ‘chaîne’ (structure) entre les variables de sortie ;

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



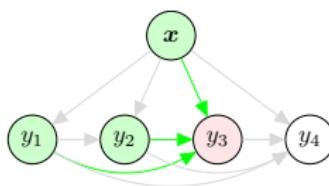
\mathbf{X}	Y_1	Y_2	Y_3	Y_4
$\mathbf{x}^{(1)}$	0	1	1	1
$\mathbf{x}^{(2)}$	1	0	0	0
$\mathbf{x}^{(3)}$	0	1	0	1
$\mathbf{x}^{(4)}$	1	0	0	0
$\mathbf{x}^{(5)}$	0	0	0	0

$\tilde{\mathbf{x}}$	\hat{y}_1	\hat{y}_2		
----------------------	-------------	-------------	--	--

(Greedy) Classifier Chains

Une ‘chaîne’ (structure) entre les variables de sortie ;

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



\mathbf{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

$\tilde{\mathbf{x}}$	\hat{y}_1	\hat{y}_2	\hat{y}_3	
			\hat{y}_3	

For example,

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

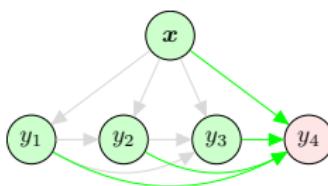
Use training data to fit base classifier (or regressor) h_3 (e.g., decision tree, logistic regression, ...).

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

(Greedy) Classifier Chains

Une ‘chaîne’ (structure) entre les variables de sortie ;

- 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

Prédictions Multi-label: On fait quoi ? Pourquoi ?

Exemples de fonctions de coût :

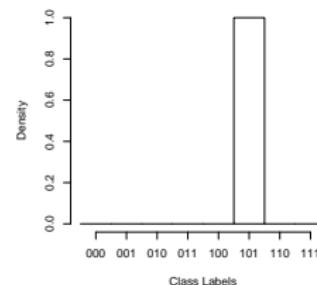
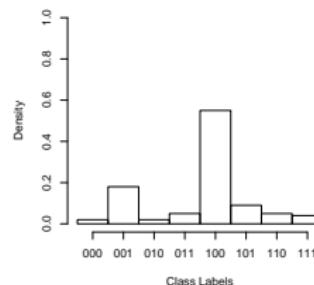
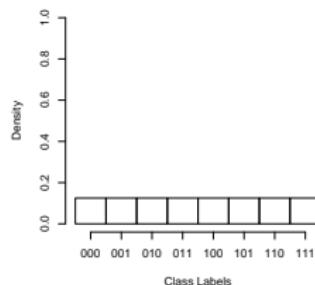
- Hamming loss: decomposable

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

- 0/1 loss: non-decomposable

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

The minimizer is not (necessarily) the same!

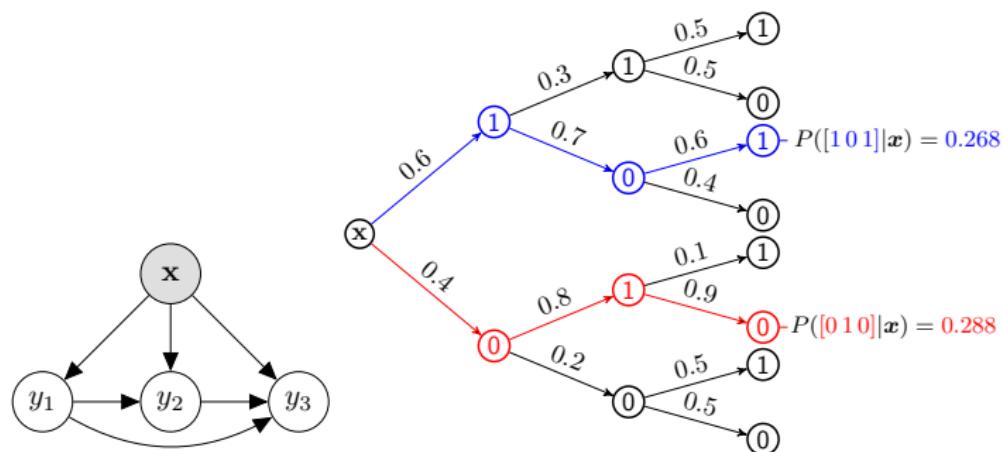


Under uncertainty, left is optimal for Hamming loss, right for 0/1-loss

Probabilistic Classifier Chains

Remarques : Au lieu d'utiliser les predictions \hat{y} comme entrées (façon *greedy*) ; on peut mettre $y_1, \dots, y_j = \mathbf{y} \in \{0, 1\}^j$; afin de minimiser **0/1-loss** :

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

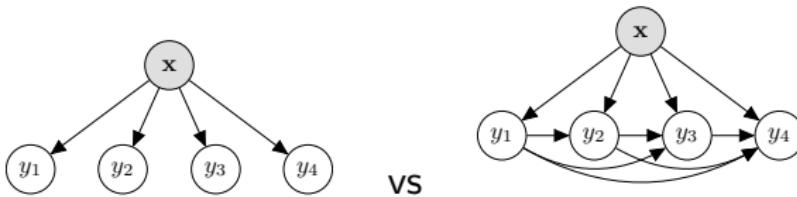


i.e., a path through the probability tree; e.g., $p([0 1 0] | \mathbf{x}) = 0.288$

Motivation for Structure in Multi-Target Learning ?

Common argument: Because label dependence! Yes for 0/1 loss.

But 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 (ainsi que le but de modéliser la probabilité combinée pour minimiser 0/1-loss, etc.):

- Connectivity = **efficiency** (sometimes)
- Connectivity = **interpretation** (sometimes)
- Connectivity = **power** (it's why deep nets or stacking works¹)
- 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)

¹Different 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", 2019 (accepted minor revisions)

Other reasons for modelling targets together (ainsi que le but de modéliser la probabilité combinée pour minimiser 0/1-loss, etc.):

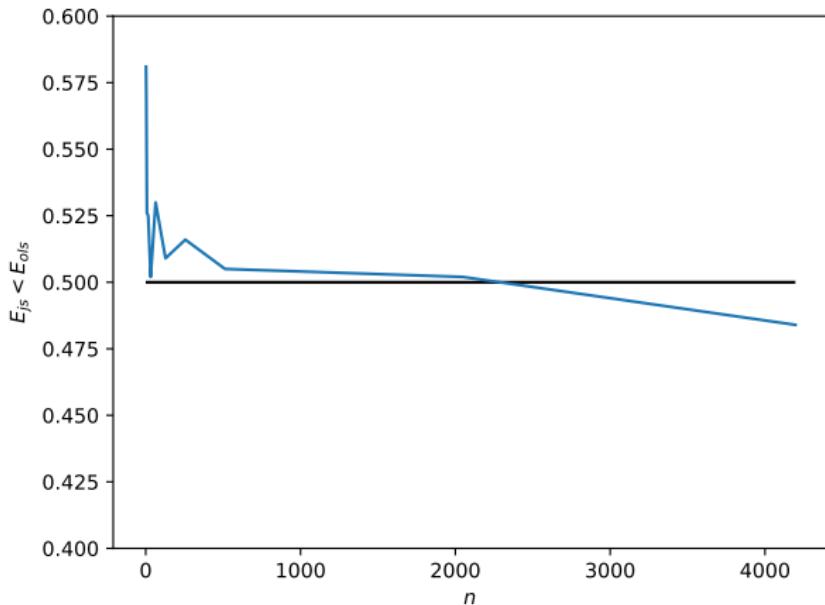
- Connectivity = efficiency (sometimes)
- Connectivity = interpretation (sometimes)
- Connectivity = power (it's why deep nets or stacking works¹)
- In same cases the minimizer *is* the same (e.g., low-noise scenarios / prediction is easy) = surrogates work well.
- Multiple tasks = regularization (regularization is good)

James-Stein Estimator

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

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

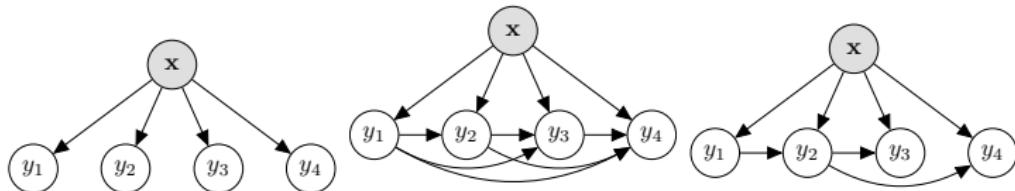
Waegeman, Dembczyński, and Hüllermeier, "Multi-target prediction: a unifying view on problems and methods", 2019; Read et al., "Classifier Chains: A Review and Perspectives", 2019 (accepted minor revisions)



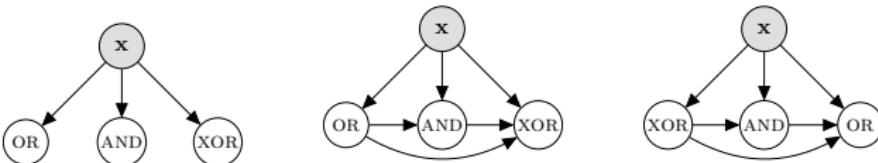
Advantages quickly fade as $n \gg 0$.

Explains reemergence of independent models vs structure debate...

Quelle(s) Structure(s) est la Optimale

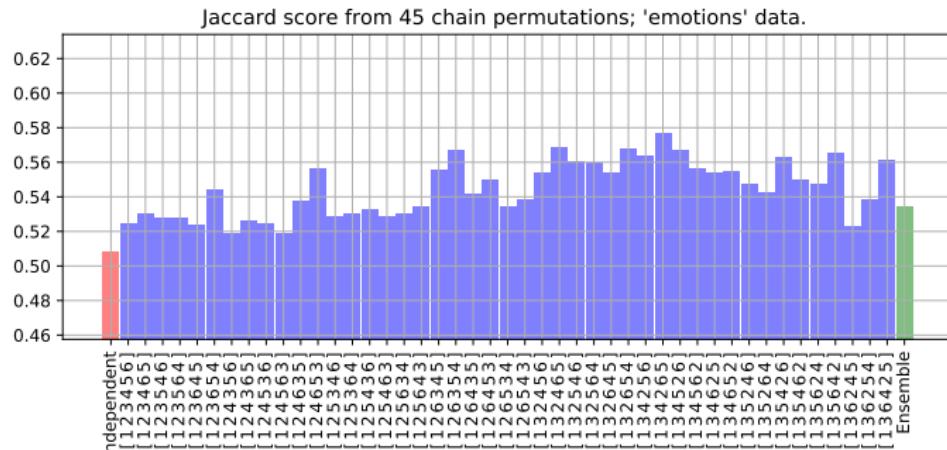


- Different chain *orders* are equivalent in theory if you have P
- Dependence is not the only component to consider (and your hierarchy is probably not better than a random one)
- Weaker base learner/smaller training set \Rightarrow more connectivity
- Weaker (greedy) inference = choose more carefully
- Best structure for loss ℓ_a , may not be the best for loss ℓ_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

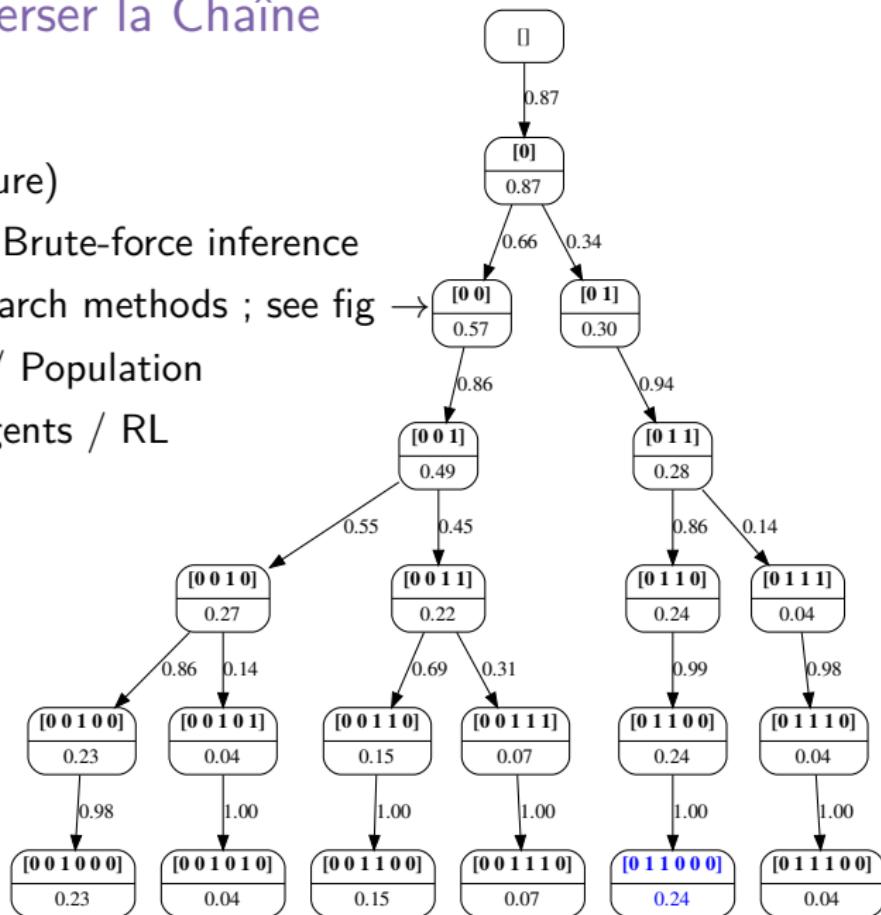
Where $x \in \{0, 1\}^2$; Base-model = Logistic regression; † But not greedy inference!



Comment Traverser la Chaîne

(given a structure)

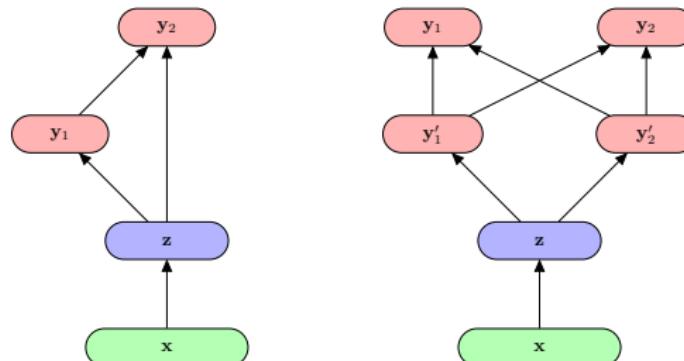
- Greedy vs Brute-force inference
- AI Tree-search methods ; see fig →
- Dynamic / Population
- Generic agents / RL



Vs Réseaux de Neurones Profonds ?

Est-ce qu'on peut faire face aux réseaux de neurones profonds ?

- **par rapport à performance**: oui, mais pas comme avant (parce que les bases de données sont plus grandes)
- **par rapport à l'explicabilité**: peut-être;
 - par exemple si on utilise les arbres de décision
 - connexions directes entre les étiquettes (pas via une couche obscure)
- pas besoin de faire face – on peut combiner :



A combination of chaining and deep-neural architectures

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	?	?	?

On peut utiliser des modèles indépendants (comme avant)

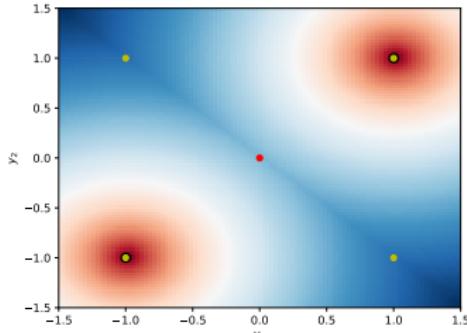
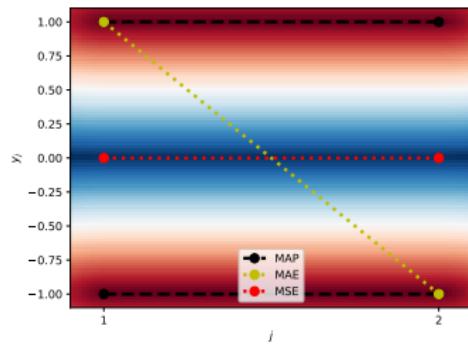
Facile de mettre les modèles dans une chaîne (comme avant; *regression chains*) ; mais typiquement c'est *inutile de le faire* (*regressor chains* \approx *independent classifiers*) , parce que :

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

Regressor Chains

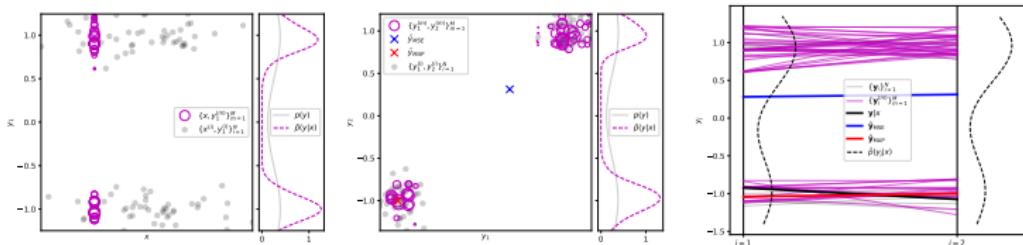
On peut

- Utiliser très bien la structure (non-linearity)
(mais encore difficile de faire face aux arbres multi-output)
- Regarder des autres fonctions de coût – par exemple, la minimisation d'un mode (et pas forcément MSE ou MAE)



Deux trajectoires incertains (donné \boldsymbol{x}) sur $y_1 \in \mathbb{R}$, $y_2 \in \mathbb{R}$: MSE vs MAE vs MAP approx.

Sequential Monte Carlo Methods for tracking modal predictions (i.e., **trajectories**)²:



Related approach : Multi-target regression via output space quantization³

Other options: Multi-target decision trees⁴ and ensembles; and deep learning.

² Read and Martino, Neurocomputing 2020

³ Spyromitros-Xioufis, Sechidis, and Vlahavas, ArXiv preprint 2020

⁴ Stepišnik and Kocev, ArXiv preprint 2020

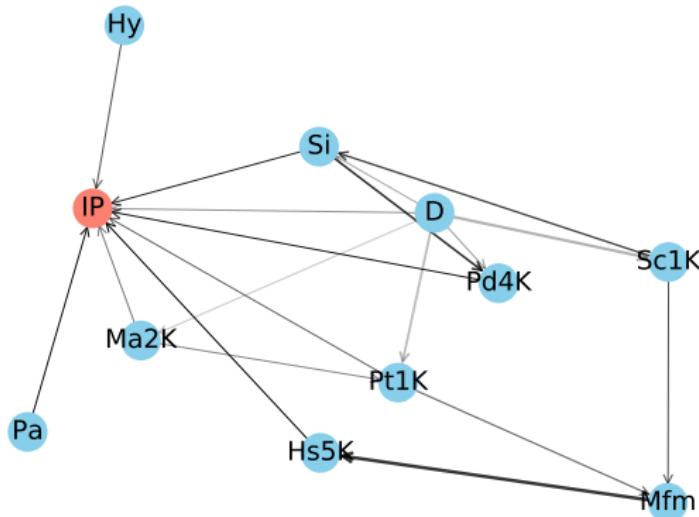
Modern Multi-Output Topics

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

Questions Ouvertes

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?



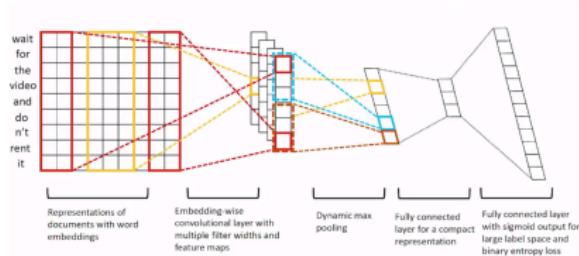
Tendances

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.).

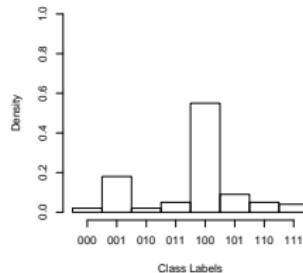
Tendances

- Importer des problématiques ‘supplémentaires’ (utilisées déjà dans des autres sous-domaines de l’apprentissage automatique) : flux de données, fouille de graphes, 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
A photograph of two people standing on a grassy hill overlooking the ocean. This image is used as a training sample with specific ground truth 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. But we argue that there is no single optimal structure, and there is more to multi-target learning than ‘modelling label dependencies’; consider metrics, efficiency, base models,
- Multi-target problems getting bigger and more diverse; intersecting with other areas
- Many applications; Likely to be an important topic in the future; but more theoretical research needed

Apprentissage Multi-cibles

Applications, Défis, et Modèles

Jesse Read



Merci !

<http://jmread.github.io/>