

# Multi-label Learning

Jesse Read



July 4, 2023

Hi! PARIS Summer School

# Multi-label Learning, Part I (Lecture, 90 min)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)
- 3 Adaptation of Classic ML Methods (5 mins)
- 4 Model-Agnostic Methods and Graphical Models (20 mins)
- 5 Deep Multi-label Learning (20 mins)
- 6 Modern Applications, Trends, and Open Areas (20 mins)
- 7 Summary and Questions (5 mins)

# Introduction and Motivation (10 mins)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)
- 3 Adaptation of Classic ML Methods (5 mins)
- 4 Model-Agnostic Methods and Graphical Models (20 mins)
- 5 Deep Multi-label Learning (20 mins)
- 6 Modern Applications, Trends, and Open Areas (20 mins)
- 7 Summary and Questions (5 mins)

# Multi-label Classification

**Multi-label classification:** a subset/vector of labels is be assigned to each input instance.



$x =$

$$y = [1, 0, 1, 0] \Leftrightarrow \{\text{Beach, Foliage}\}$$

# Multi-label Classification

**Multi-label classification:** a subset/vector of labels is assigned to each input instance.



$$\mathbf{y} = [1, 0, 1, 0] \Leftrightarrow \{\text{Beach}, \text{Foliage}\}$$

And **Multi-label Learning:** Learn the model  $h : x \mapsto \mathbf{y}$  for any  $x$ .

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
	?	?	?	?

The task (of the **model**) is to make **predictions**:

$$\hat{y} = [?, ?, ?, ?] = h(x) \quad \hat{y} \in \{0, 1\}^m$$

# Multi-Label Text (and Media) Classification



**The Lord of the Rings: The Fellowship of the Ring** (2001)

PG-13 | 178 min | **Adventure, Fantasy** | 19 December 2001 (USA) Top 500

**8.8** Your rating: ★★★★★★☆☆ -/10  
Ratings: 8.8/10 from **1,110,948 users** Metascore: 92/100  
Reviews: 4,988 user | 294 critic | 34 from Metacritic.com

A meek hobbit of the Shire and eight companions set out on a journey to Mount Doom to destroy the One Ring and the dark lord Sauron.

Director: [Peter Jackson](#)

Writers: [J.R.R. Tolkien](#) (novel), [Fran Walsh](#) (screenplay), [2 more credits](#) »

Stars: [Elijah Wood](#), [Ian McKellen](#), [Orlando Bloom](#) | [See full cast and crew](#) »

Image Source: [\[1\]](#)

The set of all possible labels (*genres*, in this case) is usually predefined.

# Labels as Keywords



Image Source: [\[2\]](#)



# Another Image-Classification Example

clear;primary














artisanal\_mine;clear;primary;water



agriculture;clear;cultivation;habitation;primary;road

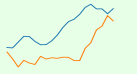
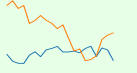


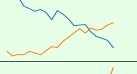
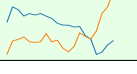


# Missing-Value Imputation and Recommender Systems

					
	0	0	1	1	0
	1	1	?	0	?
	0	0	1	0	0
	1	1	?	0	1
	0	0	0	?	?
	1	0	?	1	?

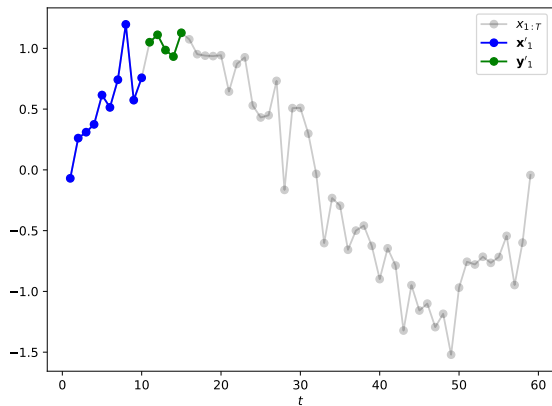
i.e., assign item-labels to users (or user-labels to items).

# Time Series Classification

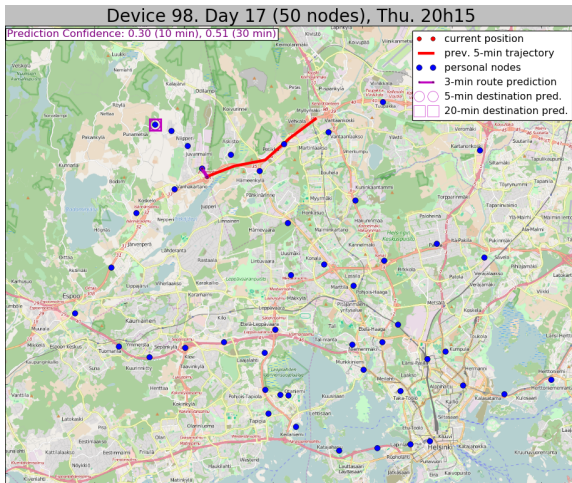
Input	Forecast 1	Forecast 2	Prescribe A	Prescribe B
	-1.01	-0.03	1	0
	0.47	-0.15	0	0
	-0.33	-0.70	1	0
	-1.39	1.57	0	1
	-0.96	1.82	0	1
	?	?	?	?

For example, ECG, EEG, signals. How will a patient's state evolve?  
Which diagnoses? Which treatments?

# Time Series Forecasting

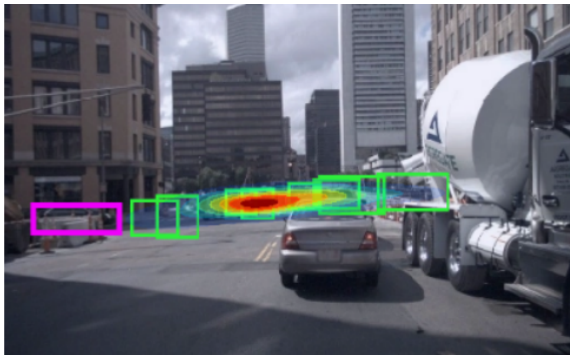


# Trajectory Prediction









Trajectory prediction in urban environment using mobile phone data

# Structured Output Prediction



Object prediction [\[3\]](#)

# Drug Design

	Mol1	Mol2	Mol3	Mol4	Mol5	Mol6
	1,3	0,2	1,4	1,7	3,5	1,3
	2	1,7	1,5	7,5	8,2	7,6
	0,2	0	0,3	0,4	1,2	2,2
	3,1	1,1	1,3	1,1	1,7	5,2
	4,7	2,1	2,5	1,5	2,3	8,5
	?	?	?	?	?	?

Molecule design prediction (binding affinities ( $Y$ ) of molecules ( $X$ ) to new proteins): [4]

# Formalization: Loss Metrics and Label Dependence (10 mins)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)**
- 3 Adaptation of Classic ML Methods (5 mins)
- 4 Model-Agnostic Methods and Graphical Models (20 mins)
- 5 Deep Multi-label Learning (20 mins)
- 6 Modern Applications, Trends, and Open Areas (20 mins)
- 7 Summary and Questions (5 mins)



# A Standard Machine Learning Setup

We are given **data set**  $\mathbf{X}, \mathbf{Y}$ . We want to build **model**  $h$  in order to obtain **predictions**

$$\hat{\mathbf{y}} = h(\mathbf{x})$$

That minimize **expected loss** where the **loss metric**

$$L(\mathbf{y}, \hat{\mathbf{y}})$$

i.e., our model  $h$  should produce

$$\min_{\hat{\mathbf{y}} \in \{0,1\}^m} \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [L(\mathbf{y}, \hat{\mathbf{y}})]$$

We might also be interested in estimating **distribution**  $p(\mathbf{y} | \mathbf{x})$ .

# Multi-label Specificities

$$\min_{\hat{\mathbf{y}} \in \{0,1\}^m} \mathbb{E}_{\mathbf{y} \sim p(\mathbf{y}|\mathbf{x})} [L(\mathbf{y}, \hat{\mathbf{y}})]$$

- Exponential complexity, wrt  $m$  labels!
- Label dependence (joint distribution)

# Important Background: Label Dependence

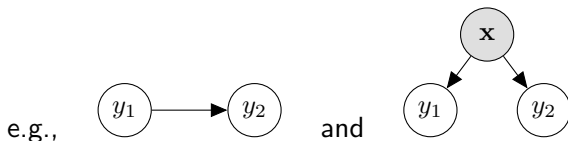
Often one considers **marginal dependence**:

$$P(y_1, y_2) \neq P(y_1)P(y_2)$$

Actually, we should be interested in **conditional dependence**:

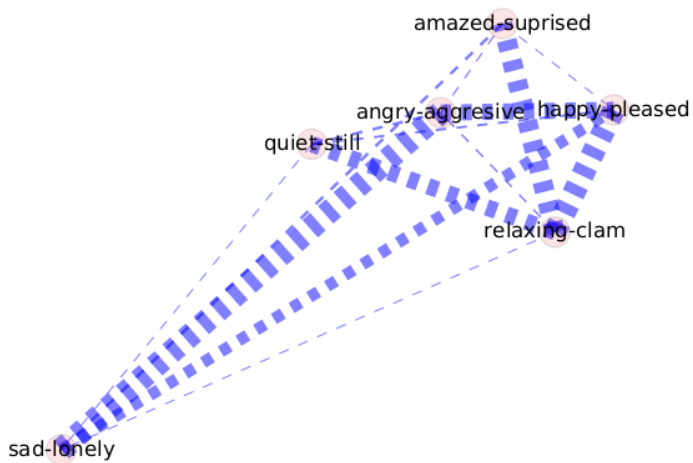
$$P(y_1, y_2 | \mathbf{x}) \neq P(y_1 | \mathbf{x})P(y_2 | \mathbf{x})$$

which is more difficult to measure (requires building models). It's not the same, e.g.,



may be equivalent!

# Example Representation of Label Dependence



Graph of correlation among the labels of the *Music-Emotions* data

## Loss Metrics ( $L$ ): How Bad is a Prediction $\hat{y}$

Example (Music/Emotions Dataset): We predict sad-lonely and angry-aggressive, but true label set is *only* sad-lonely. **How bad is this prediction?** In other words: what is the loss?

- **Hamming loss** (decomposable; **average**):

$$L_H([1, 0, 0, 0, 0, 0], [1, 0, 1, 0, 0, 0]) = 1/6$$

(not too bad)

- **0/1 loss** (non-decomposable; **exact match**):

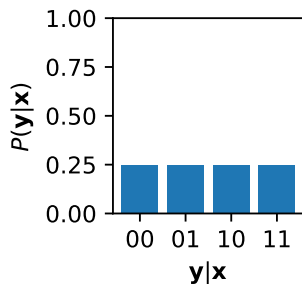
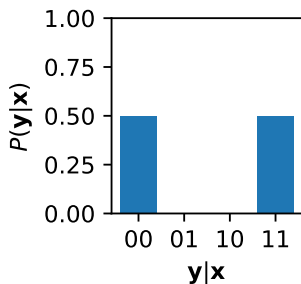
$$L_{0/1}([1, 0, 0, 0, 0, 0], [1, 0, 1, 0, 0, 0]) = 1$$

(worst case)

**The minimizer is not (necessarily) the same!** If 0/1 loss, then we need to consider the joint (predictive posterior) distribution  $p(\mathbf{y} | \mathbf{x})$ .

## Examples of $p(\mathbf{y} | \mathbf{x})$ (Predictive Posterior)

where  $\mathbf{y} \in \{0, 1\}^2$  ( $m = 2$ ), given input instance  $\mathbf{x}$ :



$P(y_1, y_2   \mathbf{x})$	$y_1 = 0$	$y_1 = 1$	$P(y_1, y_2   \mathbf{x})$	$y_1 = 0$	$y_1 = 1$
$y_2 = 0$	0.00	0.50	$y_2 = 0$	0.25	0.25
$y_2 = 1$	0.50	0.00	$y_2 = 1$	0.25	0.25

The marginal probabilities  $p(y_j | \mathbf{x})$  are the same.

## A Closer Look: Hamming Loss

Hamming loss is the averaged **sum of errors**,

$$L_{HL} = \frac{1}{m} \sum_{j=1}^m L(y_j, \hat{y}_j)$$

where, for a given label, e.g.,  $y_2$ ,

$$L(y_2, \hat{y}_2) = \begin{cases} 1 & y_2 \neq \hat{y}_2, \\ 0 & y_2 = \hat{y}_2 \end{cases}$$

i.e., it is **decomposable** across labels;

$$P(y_2 | \mathbf{x}) = \sum_{y_1 \in \{0,1\}} P(y_1 | \mathbf{x}) P(y_2 | \mathbf{x}, y_1)$$

To minimize this loss<sup>1</sup>:  **$P(\mathbf{y} | \mathbf{x})$  is not required!**  $P(y_j | \mathbf{x})$  is sufficient;

$$\hat{y}_j = h_j(\mathbf{x}) = \operatorname{argmax}_{y \in \{0,1\}} p(y_j | \mathbf{x})$$

---

<sup>1</sup>And others based upon in, like ranking loss

## A Closer Look: 0/1 Loss

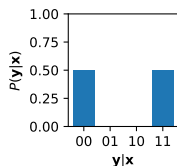
Subset 0/1 loss, is an **exact match**,

$$L_{0/1}(\mathbf{y}, \hat{\mathbf{y}}) = \begin{cases} 1 & \mathbf{y} \neq \hat{\mathbf{y}}, \\ 0 & \mathbf{y} = \hat{\mathbf{y}} \end{cases} \quad (\text{exactly, i.e., } L_{HL}(\mathbf{y}, \hat{\mathbf{y}}) = 0)$$

We need to model label dependence! We need to know  $p(\mathbf{y} | \mathbf{x})$ .

To minimize this loss:

$$\hat{\mathbf{y}} = h(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^m} p(\mathbf{y} | \mathbf{x})$$



$P(y_1, y_2   \mathbf{x})$	$y_1 = 0$	$y_1 = 1$
$y_2 = 0$	0.00	0.50
$y_2 = 1$	0.50	0.00

$P(y_2 = 1 | \mathbf{x}) = 0.5$ , but  $P(y_2 = 1, y_1 = 0 | \mathbf{x}) = 0$ ! Best case (without joint model):  $\mathbb{E}[L_{0/1}] = 0.75$  loss. Best case (with joint model):  $\mathbb{E}[L_{0/1}] = 0.5$  loss.



# Adaptation of Classic ML Methods (5 mins)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)
- 3 Adaptation of Classic ML Methods (5 mins)**
- 4 Model-Agnostic Methods and Graphical Models (20 mins)
- 5 Deep Multi-label Learning (20 mins)
- 6 Modern Applications, Trends, and Open Areas (20 mins)
- 7 Summary and Questions (5 mins)

# A Typical Offering

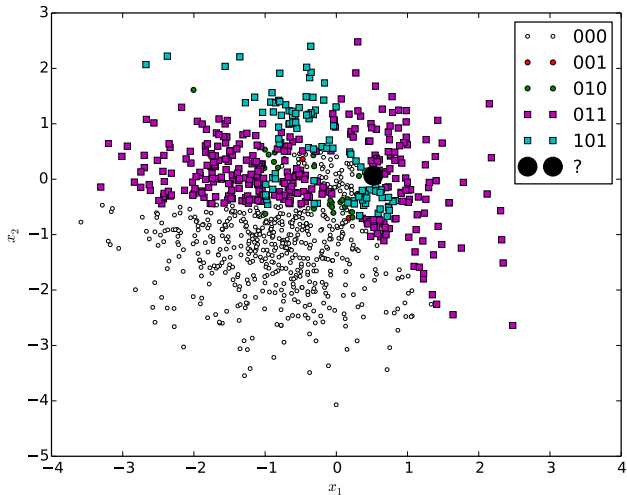
For example, **algorithm adapted** methods in SCIKITLEARN:

- `sklearn.tree.DecisionTreeClassifier`
- `sklearn.tree.ExtraTreeClassifier`
- `sklearn.ensemble.ExtraTreesClassifier`
- `sklearn.neighbors.KNeighborsClassifier`
- `sklearn.neural_network.MLPClassifier`
- `sklearn.ensemble.RandomForestClassifier`
- `sklearn.linear_model.RidgeClassifierCV`
- `sklearn.multiclass.OneVsRestClassifier`
- `sklearn.multioutput.ClassifierChain`

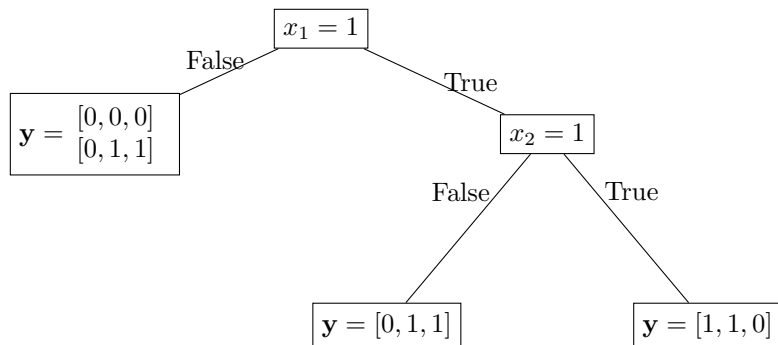
i.e., [Decision Trees](#), [Nearest-Neighbours](#), [Neural Networks](#).

... and some **task adaptation** / problem transformation / model agnostic methods – we come back to these soon!

# k-Nearest Neighbours

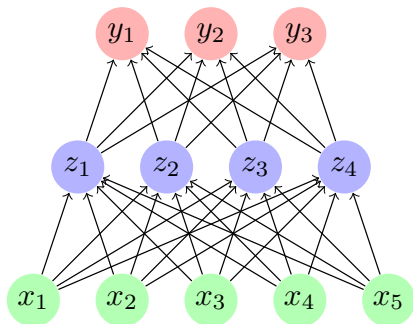


## Decision Tree Methods



Multi-labelled examples at the leaves; summation over (labels) wrt impurity criteria when inducing the tree.

# Neural Networks (Multi-Layer Perceptrons)



# Limitations of Algorithm-Adaptations

Much of the multi-label literature (and industry application) is dominated by these methods. However,

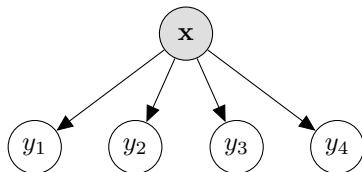
- you get stuck with a particular class of model (inflexible)
- in many cases, a reliable probabilistic interpretation is missing
- a bit 'old fashioned'; not well adapted to image or text input

# Model-Agnostic Methods and Graphical Models (20 mins)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)
- 3 Adaptation of Classic ML Methods (5 mins)
- 4 Model-Agnostic Methods and Graphical Models (20 mins)**
- 5 Deep Multi-label Learning (20 mins)
- 6 Modern Applications, Trends, and Open Areas (20 mins)
- 7 Summary and Questions (5 mins)

# Transformation to Independent Binary Classification

$\mathbf{X}$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	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
$\tilde{x}$	?	?	?	?

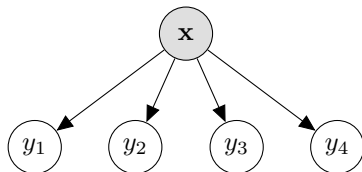


The **binary relevance method** (BR transformation) = *one binary classifier trained for each label*, i.e., **independent models**.



# Transformation to Independent Binary Classification

$\mathbf{X}$	$Y_1$	$\mathbf{X}$	$Y_2$	$\mathbf{X}$	$Y_3$	$\mathbf{X}$	$Y_4$
$\mathbf{x}^{(1)}$	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**.

## Transformation to Multi-Class (Meta-Labels)

e.g., beach+sunset considered a single label.

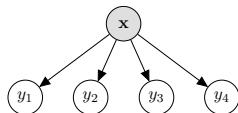
$\mathbf{X}$	$Y$
$\mathbf{x}^{(1)}$	0 1 1 0
$\mathbf{x}^{(2)}$	1 0 0 0
$\mathbf{x}^{(3)}$	0 1 0 0
$\mathbf{x}^{(4)}$	1 0 0 1
$\mathbf{x}^{(5)}$	0 0 0 1
$\tilde{\mathbf{x}}$	?



The **label powerset method** (or meta-label classifier) = *a single target multi-class classifier*. Labels are modeled together, but  $(\mathbf{y} \in \{0, 1\}^m)$  ...

- Overfitting
- Complexity.

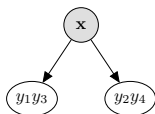
# Probabilistic Graphical Models



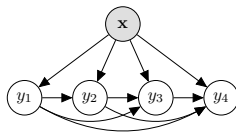
Binary Relevance



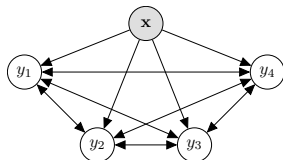
Label Powerset



Random  $k$ -Label Sets and  
Meta Labels



Prob. Classifier Chains and  
Bayesian Networks



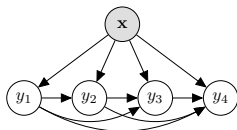
Conditional Dependency  
Networks

Arrows represent  $P(\text{child} \mid \text{parents})$  and more generally (bending the rules a bit) a prediction  $\text{output} = h(\text{input})$  where  $h$  is any **base classifier**.

# Classifier Chains: An Example of 'Problem Transformation'

A chain (structure, graph) over the output variables;

- Cascaded prediction across the chain/graph
- Motivation: Model label dependence



$X$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	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

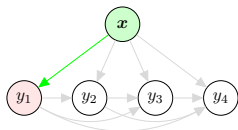
  

$\tilde{x}$	$\hat{y}_1$	$\hat{y}_2$	$\hat{y}_3$	$\hat{y}_4$
-------------	-------------	-------------	-------------	-------------

# Classifier Chains: An Example of 'Problem Transformation'

A chain (structure, graph) over the output variables;

- Cascaded prediction across the chain/graph
- Motivation: Model label dependence



$X$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	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

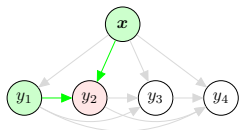
  

$\tilde{x}$	$\hat{y}_1$			
-------------	-------------	--	--	--

# Classifier Chains: An Example of 'Problem Transformation'

A chain (structure, graph) over the output variables;

- Cascaded prediction across the chain/graph
- Motivation: Model label dependence



$X$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	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

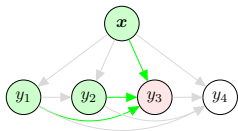
---

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

# Classifier Chains: An Example of 'Problem Transformation'

A chain (structure, graph) over the output variables;

- Cascaded prediction across the chain/graph
- Motivation: Model label dependence



$X$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$x^{(1)}$	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
$\tilde{x}$	$\hat{y}_1$	$\hat{y}_2$	$\hat{y}_3$	

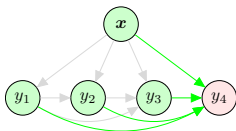
For example,  $\hat{y}_3 = h_3(x, \hat{y}_1, \hat{y}_2)$  with **base classifier** (or regressor)  $h_3$  (e.g., decision tree, logistic regression, ...).

This is a **greedy approximation** of  $\operatorname{argmax} p(\mathbf{y} | \mathbf{x})$ .

# Classifier Chains: An Example of 'Problem Transformation'

A chain (structure, graph) over the output variables;

- Cascaded prediction across the chain/graph
- Motivation: Model label dependence



$\mathbf{X}$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
$\mathbf{x}^{(1)}$	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
$\tilde{\mathbf{x}}$	$\hat{y}_1$	$\hat{y}_2$	$\hat{y}_3$	$\hat{y}_4$

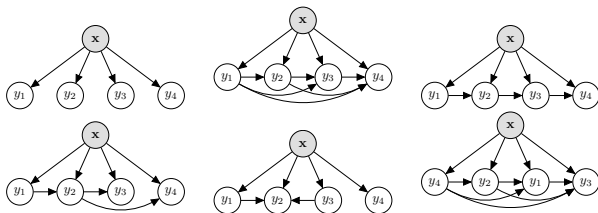
For example,  $\hat{y}_3 = h_3(\mathbf{x}, \hat{y}_1, \hat{y}_2)$  with **base classifier** (or regressor)  $h_3$  (e.g., decision tree, logistic regression, ...).

This is a **greedy approximation** of  $\operatorname{argmax} p(\mathbf{y} | \mathbf{x})$ .

FAQ. "Why this order in particular, could another one work better?"



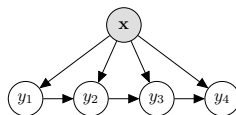
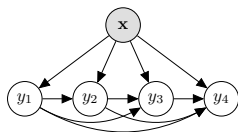
# Structure Search: Some Options



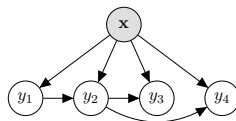
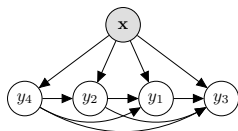
- 1 Random structure (often in an ensemble).
- 2 Use an existing hierarchy (expert knowledge)
- 3 Impose a full/complete structure
- 4 Search for a structure, based on (heuristic)
  - marginal label dependence;
  - conditional label dependence,
  - accuracy of individual models
  - accuracy of overall structure

# Structure Search is Difficult

These models perform well:



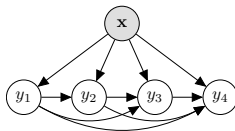
These ones perform not so well:



- Difficult to associate accuracy to a particular structure
- Considerations of measurements of dependence, time order, or 'inherent' hierarchy, are *at best* a rough guide
- A super-exponential number of possible structures
- Can never know (without uncertainty) which is the 'ground truth'

## Probabilistic Inference: Also difficult

Even with a single chosen structure,



Recall (to minimize 0/1-loss), we want:

$$\hat{\mathbf{y}} = \underset{\mathbf{y} \in \{0,1\}^m}{\operatorname{argmax}} P(\mathbf{y} | \mathbf{x}) = h(\mathbf{x})$$

$$= \underset{\mathbf{y} \in \{0,1\}^m}{\operatorname{argmax}} P(y_1 | \mathbf{x}) \prod_{j=1}^m P(y_j | \mathbf{x}, y_1, \dots, y_{j-1}) \quad \triangleright \text{from the graph}$$

e.g., (when 4 labels)

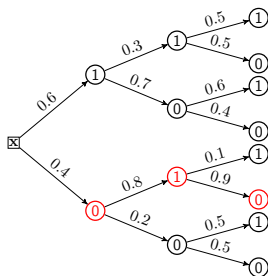
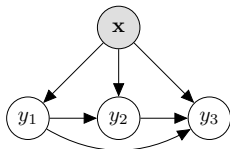
$$\mathbf{y} \in \{[0, 0, 0, 0], [0, 0, 0, 1], \dots, [1, 1, 1, 1]\}$$

and, in general,  $\mathbf{y} \in \{0, 1\}^m$  for  $m$  labels; **exponential complexity!**

# Probabilistic Classifier Chains: Inference as Tree-Search

$$\hat{y}_j = h_j(\mathbf{x}) = \operatorname{argmax}_{y_j \in \{0,1\}} P(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$$

e.g., logistic regression, then:



$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^m} P(\mathbf{y}_1 | \mathbf{x}) \prod_{j=2}^m P(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$$

This is not the same as  $\hat{y}_1, \hat{y}_2, \hat{y}_3$  obtained greedily. We now have  $p(\mathbf{y} | \mathbf{x})$ . Expensive, but many approximations via **tree search**.

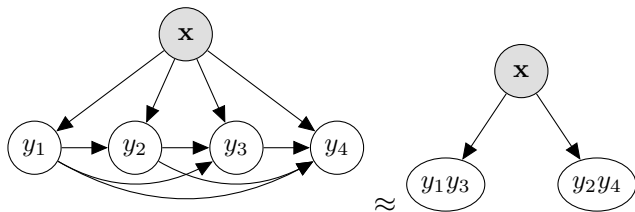
---

Dembczyński, Cheng, and Hüllermeier, "Bayes optimal multilabel classification via probabilistic classifier chains", 2010; Mena et al., "An Overview of Inference Methods in Probabilistic Classifier Chains for Multilabel Classification", 2016

# Meta Labels (e.g., RakEL) vs Probabilistic Classifier Chains

e.g., (recall) beach+sunset considered a meta label  
(transformation to multi-class as a special case).

$$\operatorname{argmax}_{\mathbf{y} \in \{0,1\}^L} P(y_1|\mathbf{x}) \prod_{j=2}^L P(y_j|\mathbf{x}, y_1, \dots, y_{j-1}) \approx \operatorname{argmax}_{\mathbf{y} \in \mathcal{S}_L \times \mathcal{S}_R} P(\mathbf{y}|\mathbf{x})$$



(more efficient search vs smaller space(s) to search through)

# Summary of Problem-Transformation Methods

We have a

- Principled way to minimize 0/1 loss (exact match);
- A flexible and interpretable (and probabilistic) structure; and
- Can use our favourite off-the-shelf classifiers (model agnostic)

But:

- (to make a long story short) sometimes the gain results from 'black magic' rather than owing to the principled approach
- Methodology tends to be scale poorly
- Still a bit old fashioned, perhaps?

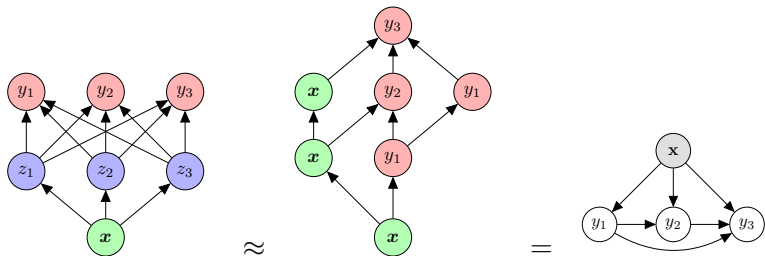
What next? Deep learning provides black magic, scalability and is fashionable!

# Deep Multi-label Learning (20 mins)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)
- 3 Adaptation of Classic ML Methods (5 mins)
- 4 Model-Agnostic Methods and Graphical Models (20 mins)
- 5 Deep Multi-label Learning (20 mins)**
- 6 Modern Applications, Trends, and Open Areas (20 mins)
- 7 Summary and Questions (5 mins)

# Graphical Models = Deep Neural Networks

We **already have this** (from graphical models): Structure among labels  $\Rightarrow$  'deep'; base classifiers as transfer functions  $\Rightarrow$  'neural'.



(' $\approx$ ' in terms of capacity; ' $=$ ' in terms of greedy inference)

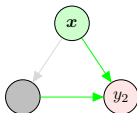
But previously, we didn't have deep *learning*:

- No back propagation
- The hidden nodes are not 'hidden'.



Consider prediction task

$$\tilde{x} \mapsto \hat{y}_2$$



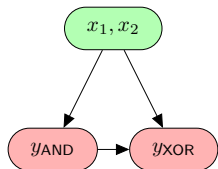
and the data available at training time (left) vs test time (right):

	$X_1$		$Y_2$	$X_1$		$Y_2$
Basis expansion	$x$	$\phi$	$y_2$	$\tilde{x}$	$\phi$	$\hat{y}_2$
Stacking	$x$	$\tilde{y}_2$	$y_2$	$\tilde{x}$	$\tilde{y}_2$	$\hat{y}_2$
Classifier chain	$x$	$y_1$	$y_2$	$\tilde{x}$	$\hat{y}_1$	$\hat{y}_2$
Neural network	$x$		$y_2$	$\tilde{x}$	$z$	$\hat{y}_2$

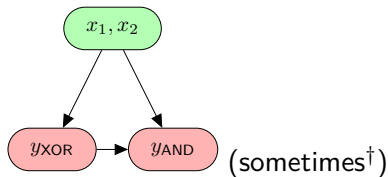
We're talking about **capacity** more than dependency here!

## A 'Logical' Problem: The 'Wrong' Dependence

$X_1$	$X_2$	$Y_{\text{XOR}}$	$Y_{\text{AND}}$
0	0	0	0
0	1	1	1
1	0	1	1
1	1	0	1



outperforms

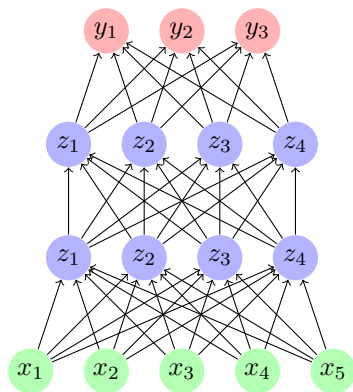


<sup>†</sup>when

$$P_{\star}(y_1, y_2 | \mathbf{x}) \neq \hat{P}(y_1, y_2 | \mathbf{x})$$

where  $\hat{P}$  depends on **base classifier**, **inference**, etc. We measured the 'wrong' dependence; but got extra **capacity** from it!

# Deep Multi-Label Learning



‘Off-the-shelf’ deep multi-label learning.

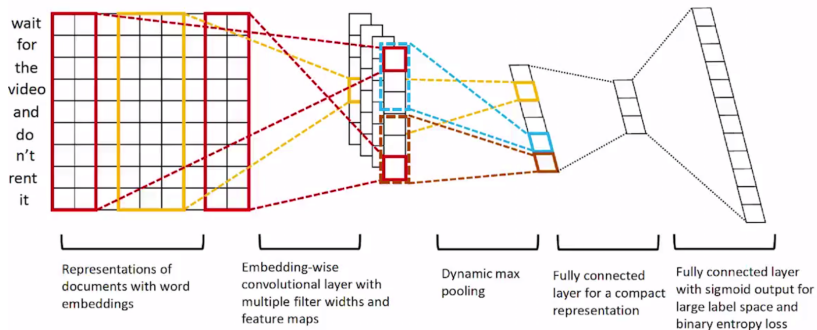
Basic idea: Powerful embeddings/capacity; go nuts with your favourite deep-learning framework (easy to add CNN, etc. layers).

---

Nam et al., “Large-Scale Multi-label Text Classification - Revisiting Neural Networks”, 2014; Read and Perez-Cruz, *Deep Learning for Multi-label Classification*, 2013; Wang et al., “CNN-RNN: A unified framework for multi-label image classification”, 2016

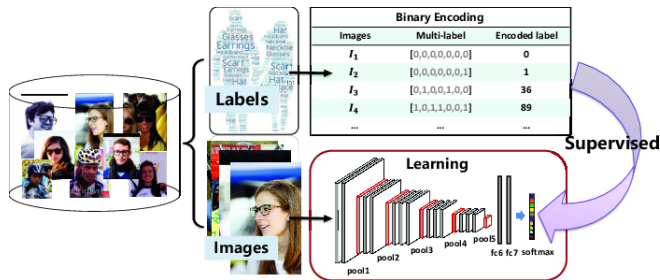
# Extreme Multi-label Classification (XMC)

An example<sup>2</sup>

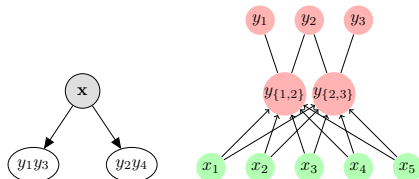


<sup>2</sup>e.g., Jasinska-Kobus et al., "Probabilistic Label Trees for Extreme Multi-label Classification", 2020

# Deep Multi-Label via Multi-Class Transformation



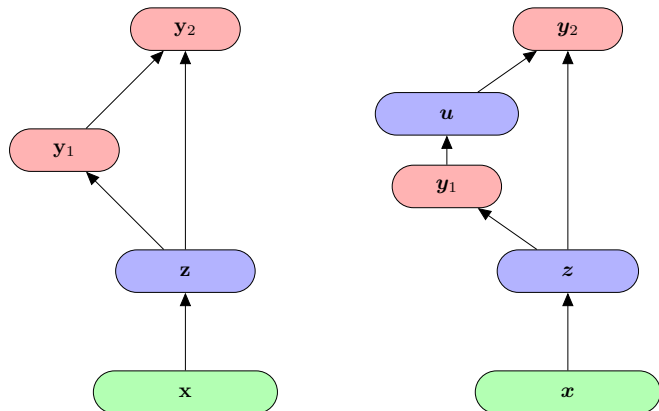
Basic idea: Transform multi-labels to single labels<sup>3</sup>;  
i.e., 'deep' version of meta labels<sup>4</sup>:



<sup>3</sup> Chenghua Li et al. "DeepBE: Learning deep binary encoding for multi-label classification". In: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 2016, pp. 39–46

<sup>4</sup> Tsoumakas, Katakis, and Vlahavas, "Random k-Labelsets for Multi-Label Classification", 2011; Read, Puurula, and Bifet, "Multi-label Classification with Meta Labels", 2014

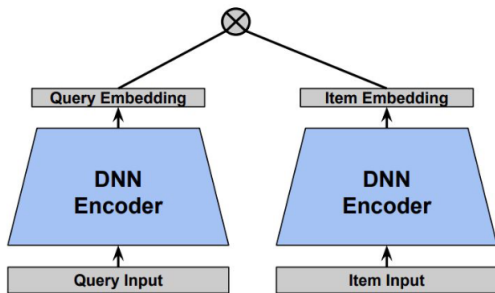
# Deep in the Label Space



Basic idea: Embeddings for the output space as well.

We can also import the ‘probabilistic chains’ into this context.

# Two-Tower Networks



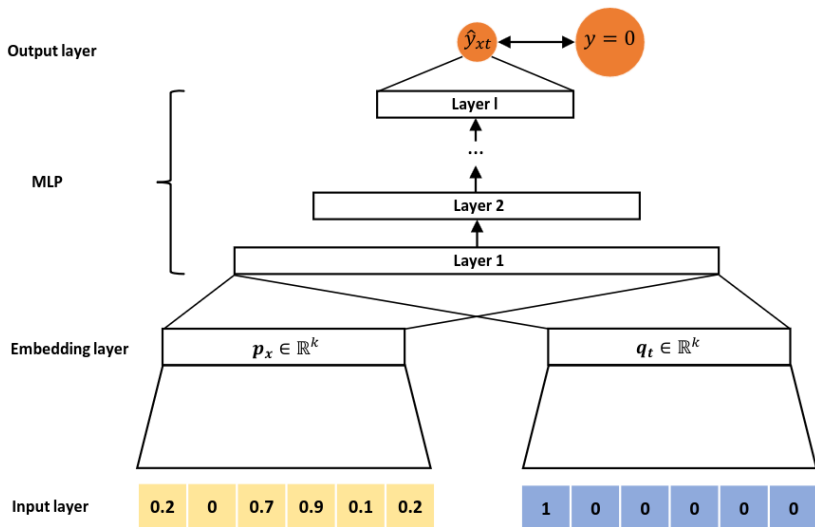
The two tower networks<sup>5</sup> have been generalized to multi-label learning<sup>6</sup>.

Basic idea: Embed the instance ( $x$ ; left); embed the item ( $j$ ; right); provide score  $y_j(x) \in \{0, 1\}$  (at the top).

---

<sup>5</sup>e.g., Yang et al., "Mixed negative sampling for learning two-tower neural networks in recommendations", 2020; He et al., "Neural collaborative filtering", 2017

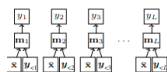
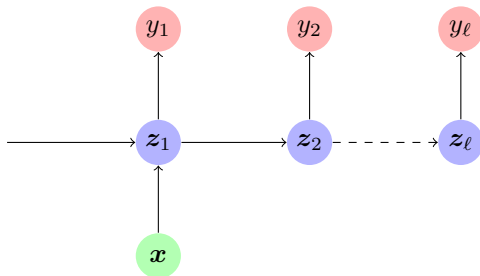
<sup>6</sup>Iliadis, De Baets, and Waegeman, "Multi-target prediction for dummies using two-branch neural networks", 2022 (in the general sense of multi-target prediction)



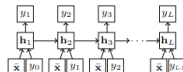
Related to the ‘independent models’ transformation, but more efficient, and flexible.



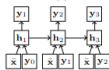
# Recurrent Neural Networks



(a) PCC



(b) RNN<sup>b</sup>



(c) RNN<sup>m</sup>



(d) EncDec

Main idea: only predict positive labels

$y_1, y_2, \dots, y_l \subset \{1, 2, \dots, m\}$ ; more efficient use of architecture.

# Modern Applications, Trends, and Open Areas (20 mins)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)
- 3 Adaptation of Classic ML Methods (5 mins)
- 4 Model-Agnostic Methods and Graphical Models (20 mins)
- 5 Deep Multi-label Learning (20 mins)
- 6 Modern Applications, Trends, and Open Areas (20 mins)**
- 7 Summary and Questions (5 mins)

# Missing Value Imputation

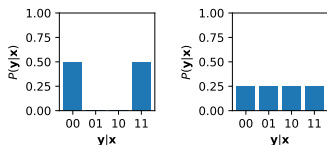
- Missing values – what to do?
- Connection to recommender systems, multi-label learning
- Missing inputs  $\approx$  noisy labels
- Where to start: Two-Tower Networks, Denoising Auto-Encoders, Expectation Maximization

$X_1$	$X_2$	$X_3$	$X_4$	Y
0	1	1	0	0
1	?	0	?	2
0	1	0	0	2
1	?	0	1	1
0	0	?	?	2

 $\Rightarrow$ 

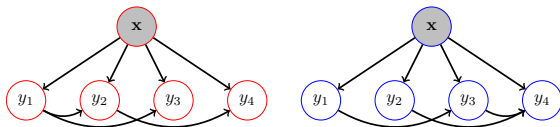
$X_1$	$X_2$	$X_3$	$X_4$	Y
0	1	1	0	0
1	1	0	0	2
0	1	0	0	2
1	0	0	1	1
0	0	0	0	2

Should keep information about uncertainty.

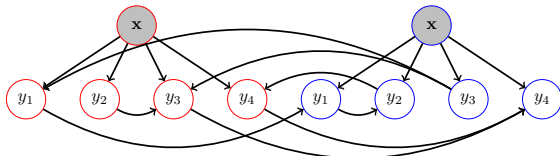


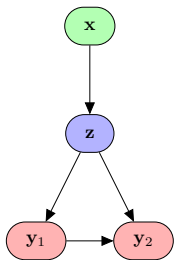
# Multi-Task and Transfer Learning

Common opinion: **label dependence** is fundamental. So, if we take two **totally unrelated datasets**; and stick them together; search for inherent structure, we should find something like this,

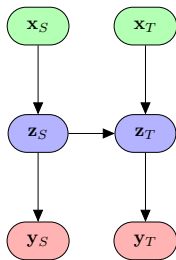


In reality: we can find something like this,

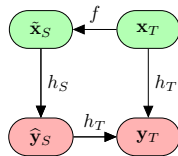




Multi-label (Chain)



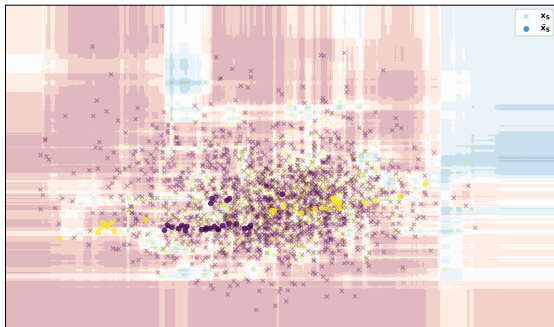
vs Deep Transfer



vs Chain Transfer.

A case study (toy example):

Source dataset: function of yeast genomes; Target dataset: insomnia diagnosis among human patients.



The predictions of genome functionality are useful features for insomnia prediction (+2% accuracy).

A hint towards Foundation Models without back-propagation; reduce, re-use (modularize), recycle models.

# Partial and Weak Labels


We get **partial labels** from noisy annotators<sup>7</sup>:



The set of candidate labels

building window  
sky street  
people car  
tree

And **weak labels**<sup>8</sup> from lazy annotators (unknown missing labels):

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

Other scenarios: **ambiguous/imprecise labels** (multiple annotators).

<sup>7</sup> e.g., Xie and Huang, "Partial multi-label learning", 2018

<sup>8</sup> e.g., Sun, Zhang, and Zhou, "Multi-label learning with weak label", 2010

# Multiple Problems in MLL: A Case Study

- Multi-labelled ECG signals (heart multi-diagnostic)
- A pre-trained deep neural network works well, but
- **poor domain transfer** (multiple collections); and
- different label sets; **missing labels** when combined.

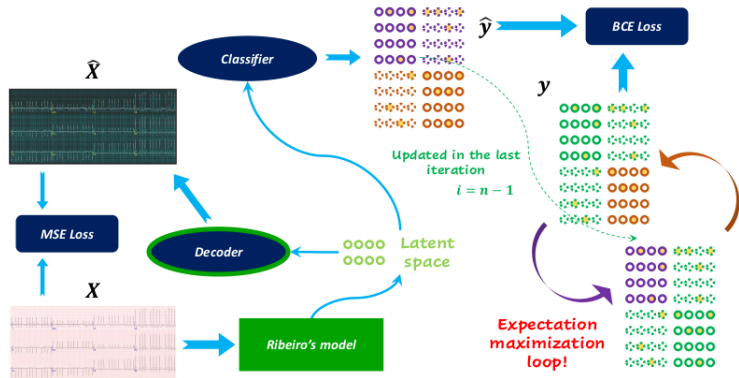


Image credit: Eran Zvuloni

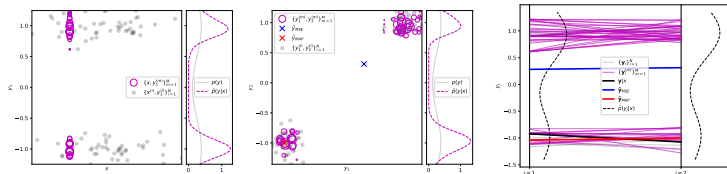


# Multi-Target Regression

So far our focus was **multi-label classification**. But modelling **continuous targets is essential** for many tasks: e.g., forecasting, structured output.

Key points:

- Several methods (e.g., greedy chains, decision trees, neural networks) can be applied off the shelf as in multi-label classification.
- Expect relatively less improvement from modelling labels together (Why? Think: **loss metric**; **non-linearities**)
- Difficulty to model  $p(\mathbf{y} | \mathbf{x})$ : tree search not possible (unless **discretization**; **Monte Carlo** tree search).



## Other Issues and Open Questions

- Data streams and concept drift (in the label space)
- Dynamic structures
- Interpretation and explainability: which graph/structure makes sense?



## Summary and Questions (5 mins)

- 1 Introduction and Motivation (10 mins)
- 2 Formalization: Loss Metrics and Label Dependence (10 mins)
- 3 Adaptation of Classic ML Methods (5 mins)
- 4 Model-Agnostic Methods and Graphical Models (20 mins)
- 5 Deep Multi-label Learning (20 mins)
- 6 Modern Applications, Trends, and Open Areas (20 mins)
- 7 Summary and Questions (5 mins)

Questions? Comments?




# References I

-  Bogatinovski, Jasmin et al. “Comprehensive comparative study of multi-label classification methods”. In: *Expert Systems with Applications* 203 (2022), p. 117215.
-  Cisse, Moustapha, Maruan Al-Shedivat, and Samy Bengio. “ADIOS: Architectures Deep In Output Space”. In: *Proceedings of The 33rd International Conference on Machine Learning*. Vol. 48. New York, New York, USA: PMLR, 2016, pp. 2770–2779.
-  Dembczyński, Krzysztof, Weiwei Cheng, and Eyke Hüllermeier. “Bayes optimal multilabel classification via probabilistic classifier chains”. In: *ICML '10: 27th International Conference on Machine Learning*. Haifa, Israel: Omnipress, 2010, pp. 279–286.

## References II

-  Dembczyński, Krzysztof et al. “On Label Dependence and Loss Minimization in Multi-label Classification”. In: *Mach. Learn.* 88.1-2 (July 2012), pp. 5–45. ISSN: 0885-6125. DOI: 10.1007/s10994-012-5285-8.
-  He, Xiangnan et al. “Neural collaborative filtering”. In: *Proceedings of the 26th international conference on world wide web.* 2017, pp. 173–182.
-  Iliadis, Dimitrios, Bernard De Baets, and Willem Waegeman. “Multi-target prediction for dummies using two-branch neural networks”. In: *Machine Learning (2022)*, pp. 1–34.
-  Jasinska-Kobus, Kalina et al. “Probabilistic Label Trees for Extreme Multi-label Classification”. In: (2020). URL: <https://arxiv.org/pdf/2009.11218v1.pdf>.
-  Li, Chenghua et al. “DeepBE: Learning deep binary encoding for multi-label classification”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition workshops.* 2016, pp. 39–46.

## References III

-  Mena, Deiner et al. “An Overview of Inference Methods in Probabilistic Classifier Chains for Multilabel Classification”. In: *Wiley Int. Rev. Data Min. and Knowl. Disc.* 6.6 (Nov. 2016). <https://digibuo.uniovi.es/dspace/bitstream/handle/10651/39325/surveyppcc-gracc.pdf>, pp. 215–230. ISSN: 1942-4787. DOI: 10.1002/widm.1185. URL: <https://doi.org/10.1002/widm.1185>.
-  Nam, Jinseok et al. “Large-Scale Multi-label Text Classification - Revisiting Neural Networks”. In: *ECML-PKDD '14: 25th European Conference on Machine Learning and Knowledge Discovery in Databases*. 2014, pp. 437–452.
-  Nam, Jinseok et al. “Maximizing Subset Accuracy with Recurrent Neural Networks in Multi-label Classification”. In: *Advances in Neural Information Processing Systems 30*. 2017, pp. 5413–5423. URL: <http://papers.nips.cc/paper/7125-maximizing-subset-accuracy-with-recurrent-neural-networks-in-multi-label-classification.pdf>.

## References IV





-  Read, Jesse. *From Multi-label Learning to Cross-Domain Transfer: A Model-Agnostic Approach*. Tech. rep. 2207.11742. ArXiv. ArXiv.org, 2023. URL: <http://arxiv.org/abs/2207.11742>.
-  Read, Jesse and Jaakko Hollmén. *Multi-label Classification using Labels as Hidden Nodes*. Tech. rep. 1503.09022v3. ArXiv. ArXiv.org, 2017. URL: <http://arxiv.org/abs/1503.09022v3>.
-  Read, Jesse and Fernando Perez-Cruz. *Deep Learning for Multi-label Classification*. Tech. rep. ArXiv. ArXiv, 2013. URL: <http://arxiv.org/abs/1502.05988>.
-  Read, Jesse, Antti Puurula, and Albert Bifet. “Multi-label Classification with Meta Labels”. eng. In: *ICDM'14: IEEE International Conference on Data Mining*. Shenzhen, China: IEEE, 2014, pp. 941–946.



## References V

-  Read, Jesse et al. “Classifier Chains: A Review and Perspectives”. In: *Journal of Artificial Intelligence Research (JAIR)* 70 (2021). <https://jair.org/index.php/jair/article/view/12376/26658>, pp. 683–718. URL: <https://jair.org/index.php/jair/article/view/12376>.
-  — .“Classifier Chains for Multi-label Classification”. In: *ECML-PKDD 2009: 20th European Conference on Machine Learning*. Bled, Slovenia: Springer, 2009, pp. 254–269. URL: [http://link.springer.com/chapter/10.1007%2F978-3-642-04174-7\\_17](http://link.springer.com/chapter/10.1007%2F978-3-642-04174-7_17).
-  Sun, Yu-Yin, Yin Zhang, and Zhi-Hua Zhou. “Multi-label learning with weak label”. In: *Proceedings of the twenty-fourth AAAI conference on artificial intelligence*. 2010, pp. 593–598.

## References VI

-  Tsoumakas, Grigorios, Ioannis Katakis, and Ioannis Vlahavas. “Random k-Labelsets for Multi-Label Classification”. In: *IEEE Transactions on Knowledge and Data Engineering* 23.7 (2011). Ed. by IEEE, pp. 1079–1089.
-  Wang, Jiang et al. “CNN-RNN: A unified framework for multi-label image classification”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 2285–2294.
-  Xie, Ming-Kun and Sheng-Jun Huang. “Partial multi-label learning”. In: *Proceedings of the AAAI conference on artificial intelligence*. Vol. 32. 1. 2018.
-  Yang, Ji et al. “Mixed negative sampling for learning two-tower neural networks in recommendations”. In: *Companion Proceedings of the Web Conference 2020*. 2020, pp. 441–447.