

# Multi-Target Prediction: A Unifying View on Problems and Methods

Willem Waegeman

Research Unit Knowledge-based Systems (KERMIT)  
Department of Mathematical Modelling, Statistics and Bioinformatics  
Ghent University, Belgium

October 13 2017

# UEFA Nations League: Rode Duivels ontmoeten Zwitserland en IJsland



WO 24/01/2018 - 12:47

In Zwitserland is vandaag geloot voor de UEFA Nations League, een gloednieuw toernooi met alle Europese voetballanden. De Rode Duivels spelen tegen Zwitserland en IJsland. De vier groepswinnaars uit divisie A, de divisie van de Belgen, strijden in 2019 in een Final Four om toernooiwinst.

Door hun goede prestaties van de voorbije jaren zitten de Rode Duivels in divisie A van de gloednieuwe

## RODE DUIVELS



Martinez: "Moeten alle Rode Duivels op dezelfde tactische pagina krijgen"



Martinez: "Zou vergissing zijn om Radja te selecteren"



Vertonghen: "We weten dat er een kans is om te winnen en daar gaan we vol voor"



Batshuayi: "Ik wil topschutter van het WK worden"



Defour stopt bij de Duivels: "Wil carrière zo lang mogelijk uitoefenen"



## Rode Duivels

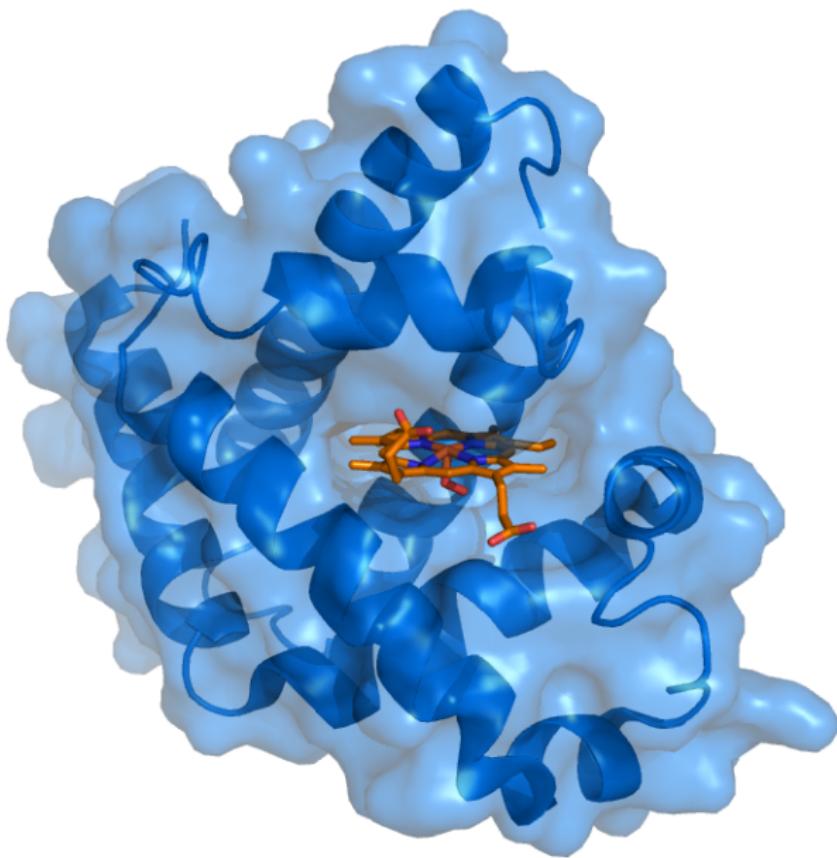
Een Twitter-lijst door @sporza

Rode Duivels



## Multi-label classification: the example of document categorization

		Tennis	Football	Biking	Movies	TV	Belgium
01101	Text1	0	1	0	0	1	1
00111	Text2	1	0	0	0	0	1
01110	Text3	0	0	0	1	1	0
10001	Text4	0	0	1	0	1	0
01011	Text5	1	0	0	1	0	0
11110	Text6	?	?	?	?	?	?



# Multivariate regression: the example of protein-ligand interaction prediction

	Mol1	Mol2	Mol3	Mol4	Mol5	Mol6	
01101		1,3	0,2	1,4	1,7	3,5	1,3
00111		2	1,7	1,5	7,5	8,2	7,6
01110		0,2	0	0,3	0,4	1,2	2,2
10001		3,1	1,1	1,3	1,1	1,7	5,2
01011		4,7	2,1	2,5	1,5	2,3	8,5
11110		?	?	?	?	?	



# Multi-task learning: the example of predicting student marks

	School1	School2	School3
01101		7	
00111		9	
01110			5
10001			8
01011			9
11110		?	?

There are a lot of multi-target prediction problems around...



# Upcoming article

## Multi-target prediction: A unifying view on problems and methods

Willem Waegeman · Krzysztof Dembczyński ·  
Eyke Hüllermeier

Received: date / Accepted: date

**Abstract** Multi-target prediction (MTP) is concerned with the simultaneous prediction of multiple target variables of diverse type. Due to its enormous application potential, it has developed into an active and rapidly expanding research field that combines several subfields of machine learning, including multivariate regression, multi-label classification, multi-task learning, dyadic prediction, zero-shot learning, network inference, and matrix completion. In this paper, we present a unifying view on MTP problems and methods. First, we formally discuss commonalities and differences between existing MTP problems. To this end, we introduce a general framework that covers the above subfields as special cases. As a second contribution, we provide a structured overview of MTP methods. This is accomplished by identifying a number of key properties, which distinguish such methods and determine their suitability for different types of problems. Finally, we also discuss a few challenges for future research.

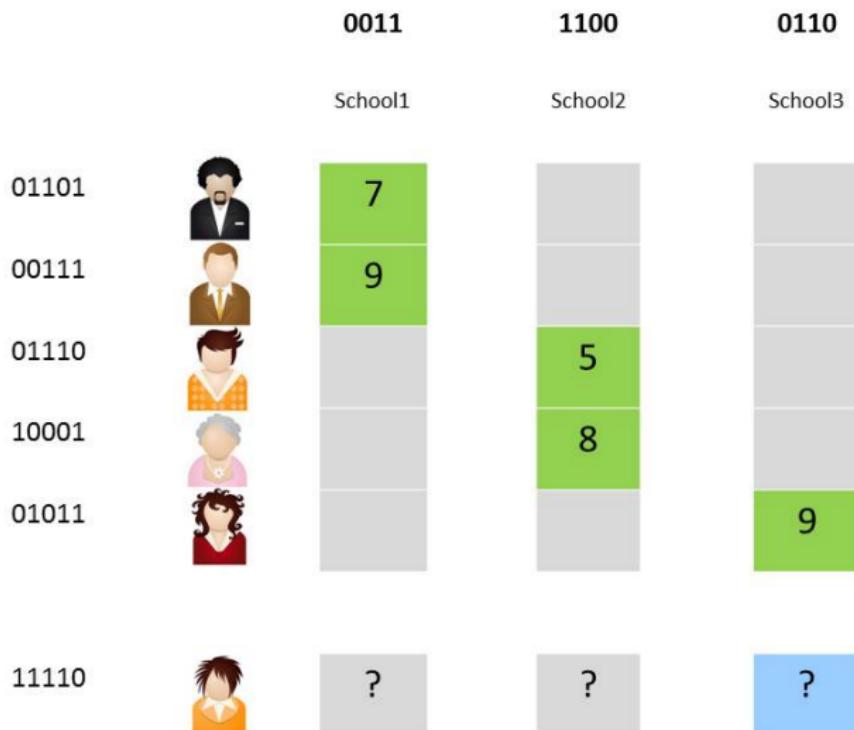
# Overview of this talk

- 1 Introduction
- 2 A unifying view on MTP problems
- 3 Loss functions in multi-target prediction
- 4 A unifying view on MTP methods
- 5 Conclusions

Let's assume a document hierarchy:  
How would you call this machine learning problem?

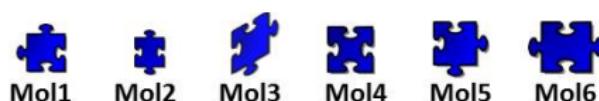
		Tags					
		Sports			Celebrities		Countries
		Tennis	Football	Biking	Movies	Tv	Belgium
01101	Text1	0	0	0	0	0	1
00111	Text2	0	0	1	0	1	1
01110	Text3	0	0	0	1	1	0
10001	Text4	0	0	1	0	1	0
01011	Text5	1	0	0	1	0	0
11110	Text6	?	?	?	?	?	?

Let's assume a target representation:  
How would you call this machine learning problem?



Let's assume a target representation:

How would you call this machine learning problem?



01101		1,3	0,2	1,4	1,7	3,5	1,3
00111		2	1,7	1,5	7,5	8,2	7,6
01110		0,2	0	0,3	0,4	1,2	2,2
10001		3,1	1,1	1,3	1,1	1,7	5,2
01011		4,7	2,1	2,5	1,5	2,3	8,5

11110		?	?	?	?	?	?
-------	--	---	---	---	---	---	---

## Generalizing to new targets

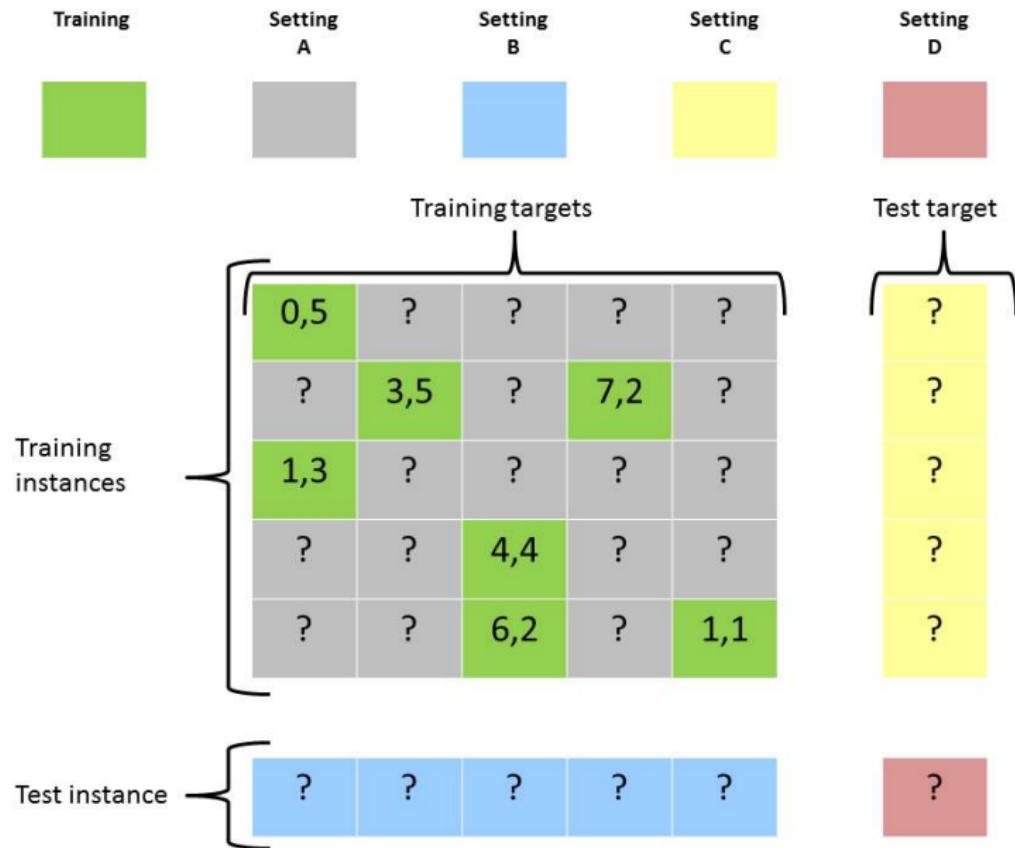
$g(., .)$  : target similarity



01101		1,3	0,2	1,4	1,7	3,5	1,3	?
00111		2	1,7	1,5	7,5	8,2	7,6	?
01110		0,2	0	0,3	0,4	1,2	2,2	?
10001		3,1	1,1	1,3	1,1	1,7	5,2	?
01011		4,7	2,1	2,5	1,5	2,3	8,5	?

11110		?	?	?	?	?	?	?
-------	--	---	---	---	---	---	---	---

# Important subdivision of different learning settings



## General framework

### Definition (Multi-target prediction)

A multi-target prediction setting is characterized by instances  $\mathbf{x} \in \mathcal{X}$  and targets  $\mathbf{t} \in \mathcal{T}$  with the following properties:

- P1. A training dataset  $\mathcal{D}$  consists of triplets  $(\mathbf{x}_i, t_j, y_{ij})$ , where  $y_{ij} \in \mathcal{Y}$  denotes a score that characterizes the relationship between the instance  $\mathbf{x}_i$  and the target  $t_j$ .
- P2. In total,  $n$  different instances and  $m$  different targets are observed during training, with  $n$  and  $m$  finite numbers. Thus, the scores  $y_{ij}$  of the training data can be arranged in an  $n \times m$  matrix  $Y$ , which is in general incomplete, i.e.,  $Y$  has missing values.
- P3. The score set  $\mathcal{Y}$  is one-dimensional. It consists of nominal, ordinal or real values.
- P4. The goal consists of predicting scores for any instance-target couple  $(\mathbf{x}, t) \in \mathcal{X} \times \mathcal{T}$ .

## Conventional MTP settings

- Side information for targets is normally not available.
- **Multivariate regression** (e.g., predicting whether a protein will bind to a set of experimentally developed small molecules).
- **Multi-label classification** (e.g., assigning appropriate category tags to documents).
- **Multi-task learning** (e.g., predicting student marks in the final exam for a typical high-school course).

## Conventional MTP settings

### Definition (Multivariate regression)

A multivariate regression problem is a specific instantiation of the general framework, which exhibits the following additional properties:

- P5. The cardinality of  $\mathcal{T}$  is  $m$ . This implies that all targets are observed during training.
- P6. No side information is available for targets. Without loss of generality, we can hence assign the numbers 1 to  $m$  as identifiers to targets, such that the target space is  $\mathcal{T} = \{1, \dots, m\}$ .
- P7. The score matrix  $Y$  has no missing values.
- P8. The score set is  $\mathcal{Y} = \mathbb{R}$ .

## Conventional MTP settings

### Definition (Multi-task learning)

A multi-task learning problem is a specific instantiation of the general framework, which exhibits the following additional properties:

- P5. The cardinality of  $\mathcal{T}$  is  $m$ ; this implies that all targets are observed during training.
- P6. No side information is available for targets. Again, the target space can hence be taken as  $\mathcal{T} = \{1, \dots, m\}$ .
- P8a. The score set is homogenous across columns of  $Y$ , e.g.,  $\mathcal{Y} = \{0, 1\}$  or  $\mathcal{Y} = \mathbb{R}$ .

# Conventional MTP settings

## Definition (Multi-label classification)

A multi-label classification problem is a specific instantiation of the general framework, which exhibits the following additional properties:

- P5. The cardinality of  $\mathcal{T}$  is  $m$ ; this implies that all targets are observed during training.
- P6. No side information is available for targets. Again, without loss of generality, we can hence identify targets with natural numbers, such that the target space is  $\mathcal{T} = \{1, \dots, m\}$ .
- P7. The score matrix  $Y$  has no missing values.
- P8b. The score set is  $\mathcal{Y} = \{0, 1\}$ .

## Conventional MTP settings

### Definition (Label ranking)

A multi-label classification problem is a specific instantiation of the general framework, which exhibits the following additional properties:

- P5. The cardinality of  $\mathcal{T}$  is  $m$ ; this implies that all targets are observed during training.
- P6. No side information is available for targets. Again, without loss of generality, we can hence identify targets with natural numbers, such that the target space is  $\mathcal{T} = \{1, \dots, m\}$ .
- P7. The score matrix  $Y$  has no missing values.
- P8c. The score set is  $\mathcal{Y} = \{1, \dots, m\}$ , and the scores (interpreted as ranks) are such that  $y_{ij} \neq y_{ik}$  for all  $1 \leq j, k \neq m$ .

## Conventional MTP settings

- In **label ranking**<sup>1</sup>, each instance is associated with a ranking (total order) of the targets.
- Thus, the score  $y_{ij} \in \{1, \dots, m\}$  for a pair  $(x_i, t_j)$  is the position of  $t_j$  in the ranking associated with  $x_i$ .
- Training data:  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , where  $y_i$  is a ranking (permutation) of a fixed number of labels/alternatives.
- Predict permutation  $(y_{\pi(1)}, y_{\pi(2)}, \dots, y_{\pi(m)})$  for a given  $x$ .

	$X_1$	$X_2$	$Y_1$	$Y_2$	$Y_3$
$x_1$	5.0	4.5	1	3	2
$x_2$	2.0	2.5	2	1	3
$\vdots$	$\vdots$	$\vdots$	$\vdots$		$\vdots$
$x_n$	3.0	3.5	3	1	2
$x$	4.0	2.5	?	?	?

<sup>1</sup> E.H., J. Fürnkranz, W. Cheng, K. Brinker. Label Ranking by Learning Pairwise Preferences, Artificial Intelligence, 172, 2008.

## Learning with side information on targets

- Additional side information about the target space is available.
- Examples:
  - ▶ Representation for the target molecules in drug design application (**structured representation**).
  - ▶ Taxonomie on document categories (**hierarchy**).
  - ▶ Information about schools and courses (geographical location, qualifications of the teachers, reputation of the school, etc.) in student mark forecasting application (**feature representation**).
- Such problems are often referred to as dyadic prediction, link prediction, or network inference settings.

## Learning with side information on targets

- Generally speaking, such settings cover problems that obey the four properties listed in the MTP definition.
- Labels  $y_{ij}$  can be arranged in a matrix  $Y$ , which is often sparse.
- Thus, one may argue that **dyadic prediction** is nothing else than **multi-task learning with task features**.
- However, MTP terminology is rarely used in the dyadic prediction literature.

## Inductive versus transductive learning problems

- In the previous problems,
  - ▶ predictions need to be generated for novel instances,
  - ▶ whereas the set of targets is known beforehand and observed during the training phase.
- These problems are **inductive** w.r.t. instances and **transductive** w.r.t. targets.
- **Side information** is of crucial importance for generalizing to novel targets that are unobserved during the training phase.

# Inductive versus transductive learning problems

## Definition (Zero-shot learning)

A zero-shot learning problem is a specific instantiation of the general framework with the following additional property:

P5\*.  $m < m^* = |\mathcal{T}|$ . Some targets are hence not observed during training, but may nevertheless appear at prediction time.

- By substituting P5 with P5\*, one now tackles problems that are inductive instead of transductive w.r.t. targets.
- The same subdivision can be made for instances.
- In total, the four different settings referred to as A, B, C, D can be distinguished (in the presence of side information).
- Theoretically, settings B and C are identical/symmetric, though there are practical differences/asymmetries.

# Inductive versus transductive learning problems

## Definition (Matrix completion)

A matrix completion problem is a specific instantiation of the general framework with the following additional properties:

- P5. The cardinality of  $\mathcal{T}$  is  $m$ . This implies that all targets are observed during training.
- P6. No side information is available for targets. Without loss of generality, we can hence assign identifiers to targets from the set  $\{1, \dots, m\}$  such that the target space is  $\mathcal{T} = \{1, \dots, m\}$ .
- P9. The cardinality of  $\mathcal{X}$  is  $n$ . This implies that all instances are observed during training.
- P10. No side information is available for instances. Without loss of generality, we can hence assign identifiers to instances from the set  $\{1, \dots, n\}$ , such that the instance space is  $\mathcal{X} = \{1, \dots, n\}$ .

## What we don't cover under MTP

- Our formal framework is rather generic.
- In principle, every prediction problem with (original) output space  $\mathcal{Y}$  could be seen as a special case by taking  $\mathcal{T} = \mathcal{Y}$  and  $\{0, 1\}$  as a score set (or through any other binary reduction).
- Each candidate output is treated as a target, the task is to predict whether or not the sought output corresponds to that candidate.
- Consequently, a consistent prediction has to obey strong (deterministic) dependencies between the targets (a single 1, rest 0).
- Includes multi-class classification and structured output prediction (SOP) as special cases.
- Yet, we do not consider such problems as special cases of MTP.

## What we don't cover under MTP

- **Conceptually**, viewing each candidate prediction as a separate target appears artificial (actually, one is still interested in a single prediction, not multiple ones).
- To comply with the corresponding **consistency constraints**, a kind of post-processing (like the decoding step in ECOC) is normally required (separation over tasks is not even possible in principle).
- **Algorithmically**, the multi-target perspective is not typical of SOP. Instead, such methods are specifically tailored for output spaces that are often huge but equipped with a strong structure.
- Also excluded are prediction problems where the ground truth cannot be represented in a matrix format with optional side information, such as problems involving multi-instance learning representations or dyadic feature representations.

# Overview of this talk

- 1 Introduction
- 2 A unifying view on MTP problems
- 3 Loss functions in multi-target prediction
- 4 A unifying view on MTP methods
- 5 Conclusions

## Multi-target prediction

- For a feature vector  $\mathbf{x}$ , predict a vector of responses  $\mathbf{y} = (y_1, y_2, \dots, y_m)$  using a function/hypothesis  $\mathbf{h}(\mathbf{x})$ :

$$\mathbf{x} = (x_1, x_2, \dots, x_p) \quad \xrightarrow{\mathbf{h}(\mathbf{x})} \quad \hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_m)$$

- Compared to single-target prediction, a multitude of **multivariate loss functions**

$$\ell : \mathcal{Y}^m \times \mathcal{Y}^m \rightarrow \mathbb{R}$$

are conceivable.

- Problem:** Given a target loss  $\ell(\cdot)$ , find a (Bayes) predictor  $\mathbf{h}(\cdot)$  that minimizes expected loss with regard to  $\ell(\cdot)$ .
- Key question:** Can we achieve this goal through simple reduction, i.e., by training one model for each target independently? Or can we do better with more sophisticated methods?

## The individual target view

- How can we improve the predictive accuracy of a single label by exploiting information about other labels?
- Goal: predict a value of  $y_i$  using  $\mathbf{x}$  and any available information on other targets  $y_j$ .
- The problem is usually defined through univariate losses  $\ell_i(y_i, \hat{y}_i)$ .
- Domain of  $y_i$  is either continuous or nominal.
- Independent models vs. regularized (shrunken) models.

## The James-Stein paradox

## The joint target view

- The problem is defined through multivariate losses  $\ell(\mathbf{y}, \hat{\mathbf{y}})$ .
- Is reduction to single-target prediction (decomposition over targets) still possible, and even if so, can we improve over such strategies by using more expressive models?
- Important: **Structure of loss**  $\ell(\cdot)$ , possible **dependencies between targets**, multivariate distribution of  $\mathbf{y}$ .

## Multivariate loss functions

- **Decomposable** over examples: A loss  $L$  is decomposable over examples if it can be written in the form

$$L = \sum_{i=1}^n \ell(\mathbf{y}_i, \mathbf{h}(\mathbf{x}_i)),$$

i.e., as a sum of losses over all (test) examples.

- **Decomposable** over targets: A multivariate loss  $\ell$  is decomposable over targets if it can be written as

$$\ell(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \sum_{i=1}^m \ell_i(y_i, h_i(\mathbf{x}))$$

with suitable single-target losses  $\ell_i$ .

# Macro- and micro-averaging

- Macro-averaging

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{4} (L_1 + L_2 + L_3 + L_4)$$

# Macro- and micro-averaging

- Macro-averaging

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{4} (\textcolor{red}{L_1} + L_2 + L_3 + L_4)$$

# Macro- and micro-averaging

- Macro-averaging

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{4} (L_1 + L_2 + L_3 + L_4)$$

# Macro- and micro-averaging

- Macro-averaging

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{4} (L_1 + L_2 + L_3 + L_4)$$

# Macro- and micro-averaging

- Macro-averaging

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{4} (L_1 + L_2 + L_3 + L_4)$$

# Macro- and micro-averaging

- Micro-averaging

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \sum_{i,j} \ell(y_{ij}, \hat{y}_{ij})$$

# Macro- and micro-averaging

- Micro-averaging

True labels				Predicted labels			
$y_{11}$	$y_{12}$		$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$		$\hat{y}_{14}$
$y_{21}$		$y_{23}$	$y_{24}$	$\hat{y}_{21}$		$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$		$y_{43}$	$y_{44}$	$\hat{y}_{41}$		$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
	$y_{62}$	$y_{63}$			$\hat{y}_{62}$	$\hat{y}_{63}$	

$$L = \sum_{i,j} \ell(y_{ij}, \hat{y}_{ij})$$

Same weight of every predictions vs. same weight of every target.

## Instance-wise losses

- Averaging over instances

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{6} \left( \ell(\mathbf{y}_1, \hat{\mathbf{y}}_1) + \ell(\mathbf{y}_2, \hat{\mathbf{y}}_2) + \ell(\mathbf{y}_3, \hat{\mathbf{y}}_3) + \ell(\mathbf{y}_4, \hat{\mathbf{y}}_4) + \ell(\mathbf{y}_5, \hat{\mathbf{y}}_5) + \ell(\mathbf{y}_6, \hat{\mathbf{y}}_6) \right)$$

## Instance-wise losses

- Averaging over instances

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{6} \left( \ell(\mathbf{y}_1, \hat{\mathbf{y}}_1) + \ell(\mathbf{y}_2, \hat{\mathbf{y}}_2) + \ell(\mathbf{y}_3, \hat{\mathbf{y}}_3) + \ell(\mathbf{y}_4, \hat{\mathbf{y}}_4) + \ell(\mathbf{y}_5, \hat{\mathbf{y}}_5) + \ell(\mathbf{y}_6, \hat{\mathbf{y}}_6) \right)$$

## Instance-wise losses

- Averaging over instances

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{6} \left( \ell(\mathbf{y}_1, \hat{\mathbf{y}}_1) + \ell(\mathbf{y}_2, \hat{\mathbf{y}}_2) + \ell(\mathbf{y}_3, \hat{\mathbf{y}}_3) + \ell(\mathbf{y}_4, \hat{\mathbf{y}}_4) + \ell(\mathbf{y}_5, \hat{\mathbf{y}}_5) + \ell(\mathbf{y}_6, \hat{\mathbf{y}}_6) \right)$$

## Instance-wise losses

- Averaging over instances

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{6} \left( \ell(\mathbf{y}_1, \hat{\mathbf{y}}_1) + \ell(\mathbf{y}_2, \hat{\mathbf{y}}_2) + \ell(\mathbf{y}_3, \hat{\mathbf{y}}_3) + \ell(\mathbf{y}_4, \hat{\mathbf{y}}_4) + \ell(\mathbf{y}_5, \hat{\mathbf{y}}_5) + \ell(\mathbf{y}_6, \hat{\mathbf{y}}_6) \right)$$

## Instance-wise losses

- Averaging over instances

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{6} \left( \ell(\mathbf{y}_1, \hat{\mathbf{y}}_1) + \ell(\mathbf{y}_2, \hat{\mathbf{y}}_2) + \ell(\mathbf{y}_3, \hat{\mathbf{y}}_3) + \ell(\mathbf{y}_4, \hat{\mathbf{y}}_4) + \ell(\mathbf{y}_5, \hat{\mathbf{y}}_5) + \ell(\mathbf{y}_6, \hat{\mathbf{y}}_6) \right)$$

## Instance-wise losses

- Averaging over instances

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{6} \left( \ell(\mathbf{y}_1, \hat{\mathbf{y}}_1) + \ell(\mathbf{y}_2, \hat{\mathbf{y}}_2) + \ell(\mathbf{y}_3, \hat{\mathbf{y}}_3) + \ell(\mathbf{y}_4, \hat{\mathbf{y}}_4) + \ell(\mathbf{y}_5, \hat{\mathbf{y}}_5) + \ell(\mathbf{y}_6, \hat{\mathbf{y}}_6) \right)$$

## Instance-wise losses

- Averaging over instances

True labels				Predicted labels			
$y_{11}$	$y_{12}$	$y_{13}$	$y_{14}$	$\hat{y}_{11}$	$\hat{y}_{12}$	$\hat{y}_{13}$	$\hat{y}_{14}$
$y_{21}$	$y_{22}$	$y_{23}$	$y_{24}$	$\hat{y}_{21}$	$\hat{y}_{22}$	$\hat{y}_{23}$	$\hat{y}_{24}$
$y_{31}$	$y_{32}$	$y_{33}$	$y_{34}$	$\hat{y}_{31}$	$\hat{y}_{32}$	$\hat{y}_{33}$	$\hat{y}_{34}$
$y_{41}$	$y_{42}$	$y_{43}$	$y_{44}$	$\hat{y}_{41}$	$\hat{y}_{42}$	$\hat{y}_{43}$	$\hat{y}_{44}$
$y_{51}$	$y_{52}$	$y_{53}$	$y_{54}$	$\hat{y}_{51}$	$\hat{y}_{52}$	$\hat{y}_{53}$	$\hat{y}_{54}$
$y_{61}$	$y_{62}$	$y_{63}$	$y_{64}$	$\hat{y}_{61}$	$\hat{y}_{62}$	$\hat{y}_{63}$	$\hat{y}_{64}$

$$L = \frac{1}{6} \left( \ell(\mathbf{y}_1, \hat{\mathbf{y}}_1) + \ell(\mathbf{y}_2, \hat{\mathbf{y}}_2) + \ell(\mathbf{y}_3, \hat{\mathbf{y}}_3) + \ell(\mathbf{y}_4, \hat{\mathbf{y}}_4) + \ell(\mathbf{y}_5, \hat{\mathbf{y}}_5) + \ell(\mathbf{y}_6, \hat{\mathbf{y}}_6) \right)$$

## Multi-label losses

- The **Hamming loss** averages over mistakes on individual labels:

$$\ell_H(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{m} \sum_{i=1}^m \llbracket y_i \neq \hat{y}_i \rrbracket$$

- The **subset 0/1 loss** simply checks for entire correctness:

$$\begin{aligned}\ell_{0/1}(\mathbf{y}, \hat{\mathbf{y}}) &= \llbracket \mathbf{y} \neq \hat{\mathbf{y}} \rrbracket \\ &= \max_i \llbracket y_i \neq \hat{y}_i \rrbracket\end{aligned}$$

## Hamming vs. subset 0/1 loss

- What is the risk-minimizing (Bayes) prediction for the Hamming loss and the subset 0/1 loss, respectively, given the following conditional distribution  $P(Y_1, Y_2 | \mathbf{x})$ ?

$y_1$	$y_2$	$P(y_1, y_2   \mathbf{x})$
0	0	0.3
0	1	0.3
1	0	0.0
1	1	0.4

## Hamming vs. subset 0/1 loss

- What is the risk-minimizing (Bayes) prediction for the Hamming loss and the subset 0/1 loss, respectively, given the following conditional distribution  $P(Y_1, Y_2 | \mathbf{x})$ ?

$y_1$	$y_2$	$P(y_1, y_2   \mathbf{x})$
0	0	0.3
0	1	0.3
1	0	0.0
1	1	0.4

$$\begin{aligned}\mathbb{E}_{\mathbf{Y}} \ell_{0/1}(\mathbf{y}, \hat{\mathbf{y}}) &= 0.3 \ell_{0/1}((0, 0), \hat{\mathbf{y}}) + 0.3 \ell_{0/1}((0, 1), \hat{\mathbf{y}}) + \\ &\quad 0.0 \ell_{0/1}((1, 0), \hat{\mathbf{y}}) + 0.4 \ell_{0/1}((1, 1), \hat{\mathbf{y}}) \\ &= 1 - P(\hat{\mathbf{y}} | \mathbf{x})\end{aligned}$$

## Hamming vs. subset 0/1 loss

- What is the risk-minimizing (Bayes) prediction for the Hamming loss and the subset 0/1 loss, respectively, given the following conditional distribution  $P(Y_1, Y_2 | \mathbf{x})$ ?

$y_1$	$y_2$	$P(y_1, y_2)$
0	0	0.3
0	1	0.3
1	0	0.0
1	1	0.4

$$\begin{aligned}\mathbb{E}_{\mathbf{Y}} \ell_H(\mathbf{y}, \hat{\mathbf{y}}) &= 0.3 \ell_H((0, 0), \hat{\mathbf{y}}) + 0.3 \ell_H((0, 1), \hat{\mathbf{y}}) + \\ &\quad 0.0 \ell_H((1, 0), \hat{\mathbf{y}}) + 0.4 \ell_H((1, 1), \hat{\mathbf{y}}) \\ &= 0.3(\llbracket 0 \neq y_1 \rrbracket + \llbracket 0 \neq y_2 \rrbracket) + 0.3(\llbracket 0 \neq y_1 \rrbracket + \llbracket 1 \neq y_2 \rrbracket) + \\ &\quad 0.0(\llbracket 1 \neq y_1 \rrbracket + \llbracket 0 \neq y_2 \rrbracket) + 0.4(\llbracket 1 \neq y_1 \rrbracket + \llbracket 1 \neq y_2 \rrbracket) \\ &= 0.6\llbracket 0 \neq y_1 \rrbracket + 0.4\llbracket 1 \neq y_1 \rrbracket + 0.3\llbracket 0 \neq y_2 \rrbracket + 0.7\llbracket 0 \neq y_2 \rrbracket \\ &= \mathbb{E}_{Y_1} \ell_H(y_1, \hat{y}_1) + \mathbb{E}_{Y_2} \ell_H(y_2, \hat{y}_2)\end{aligned}$$

## Hamming vs. subset 0/1 loss

- The risk minimizer for the Hamming loss is the **marginal mode**:

$$h_i^*(\mathbf{x}) = \arg \max_{y_i \in \{0,1\}} P(y_i | \mathbf{x}), \quad i = 1, \dots, m,$$

while for the subset 0/1 loss it is the **joint mode**:

$$\mathbf{h}^*(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x}).$$

- Marginal mode vs. joint mode.

$\mathbf{y}$	$P(\mathbf{y})$
0 0 0 0	0.30
0 1 1 1	0.17
1 0 1 1	0.18
1 1 0 1	0.17
1 1 1 0	0.18

Marginal mode: 1 1 1 1  
Joint mode: 0 0 0 0

## Hamming vs. subset 0/1 loss

- **Proposition<sup>2</sup>:** The following upper bound holds for  $m > 3$ :

$$\mathbb{E}_{\mathbf{Y}} \ell_H(\mathbf{Y}, \mathbf{h}_{0/1}^*(\mathbf{x})) - \mathbb{E}_{\mathbf{Y}} \ell_H(\mathbf{Y}, \mathbf{h}_H^*(\mathbf{x})) < \frac{m-2}{m+2}$$

Moreover, this bound is tight, i.e.

$$\sup_P (\mathbb{E}_{\mathbf{Y}} \ell_H(\mathbf{Y}, \mathbf{h}_{0/1}^*(\mathbf{x})) - \mathbb{E}_{\mathbf{Y}} \ell_H(\mathbf{Y}, \mathbf{h}_H^*(\mathbf{x}))) = \frac{m-2}{m+2},$$

where the supremum is taken over all probability distributions on  $\mathcal{Y}$ .

---

<sup>2</sup> K.D., W.W., W. Cheng, E.H. On Label Dependence and Loss Minimization in Multi-Label Classification. Machine Learning, 88, 2012.

## Hamming vs. subset 0/1 loss

Under specific conditions, the risk minimizers for  $\ell_H$  and  $\ell_{0/1}$  are **equivalent**,

$$\mathbf{h}_H^*(\mathbf{x}) = \mathbf{h}_{0/1}^*(\mathbf{x}),$$

for example, when

- the probability of the joint mode satisfies

$$P(\mathbf{h}_{0/1}^*(\mathbf{x})|\mathbf{x}) > 0.5,$$

- or the targets  $Y_1, \dots, Y_m$  are **conditionally independent**.

# Target dependence

We distinguish between conditional and unconditional (in)dependence of targets<sup>3</sup>.

- **Unconditional/marginal dependence:**

$$P(\mathbf{Y}) \neq \prod_{i=1}^m P(Y_i)$$

Often due to model similarities, i.e.,  $f_i(\mathbf{x}) = g_i(\mathbf{x}) + \epsilon_i$  for  $i = 1, \dots, m$ , with similarities in the structural parts  $g_i(\cdot)$ .

- **Conditional dependence:**

$$P(\mathbf{Y} \mid \mathbf{x}) \neq \prod_{i=1}^m P(Y_i \mid \mathbf{x})$$

---

<sup>3</sup> K.D., W.W., W. Cheng, E.H. On Label Dependence and Loss Minimization in Multi-Label Classification. Machine Learning, 88, 2012.

## Target dependence

- marginal (in)dependence  $\not\Rightarrow$  conditional (in)dependence
- Example:

$x_1$	$y_1$	$y_2$	$P$	$x_1$	$y_1$	$y_2$	$P$
0	0	0	0.25	1	0	0	0
0	0	1	0	1	0	1	0.25
0	1	0	0	1	1	0	0.25
0	1	1	0.25	1	1	1	0

- Strong conditional dependence, for example  
 $P(Y_1 = 0|x_1 = 1)P(Y_2 = 0|x_1 = 1) = 0.5 \times 0.5 = 0.25 \neq 0.$
- Yet, labels are marginally independent: Joint probability is the product of the marginals  $P(y_1) = P(y_2) = 0.5.$

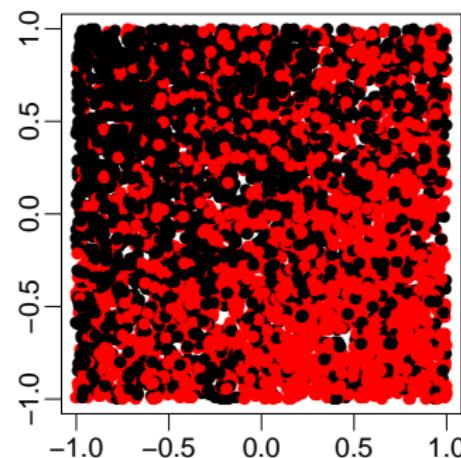
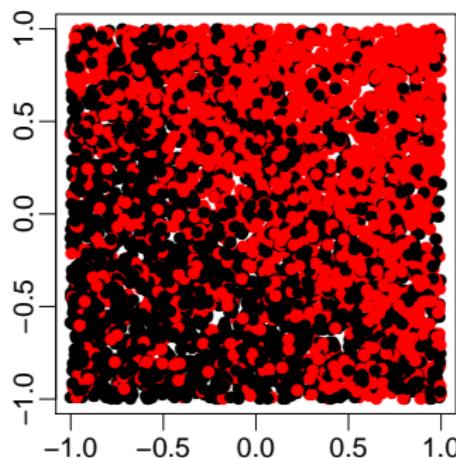
## Synthetic data

- Two conditionally independent models:

$$f_1(\mathbf{x}) = \frac{1}{2}x_1 + \frac{1}{2}x_2, \quad f_2(\mathbf{x}) = \frac{1}{2}x_1 - \frac{1}{2}x_2$$

- Logistic model to get labels:

$$P(y_i = 1) = \frac{1}{1 + \exp(-2f_i)}$$



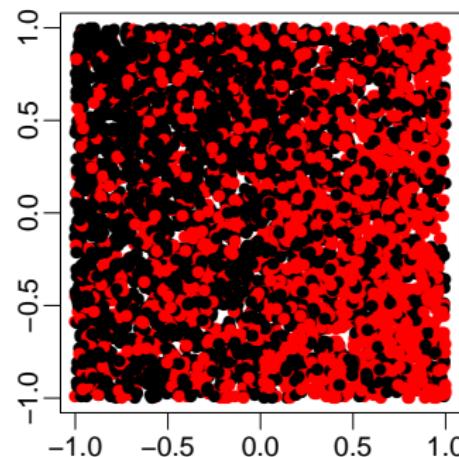
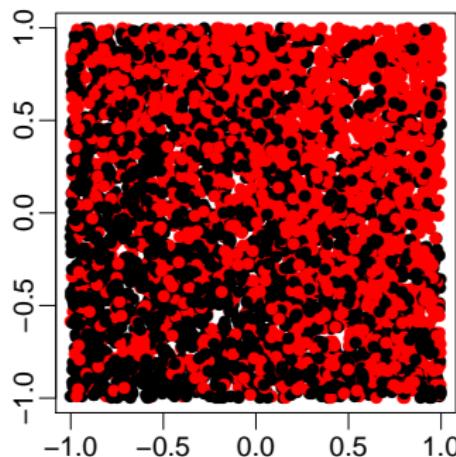
## Synthetic data

- Two dependent models:

$$f_1(\mathbf{x}) = \frac{1}{2}x_1 + \frac{1}{2}x_2 \quad f_2(y_1, \mathbf{x}) = y_1 + \frac{1}{2}x_1 - \frac{1}{2}x_2 - \frac{2}{3}$$

- Logistic model to get labels:

$$P(y_i = 1) = \frac{1}{1 + \exp(-2f_i)}$$



## Hamming vs. subset 0/1 loss

- Binary relevance (BR): Train two binary classifiers for targets  $y_1$  and  $y_2$  independently.
- Label powerset (LP): Train a 4-class classifier on meta-classes  $c_1 = (0, 0)$ ,  $c_2 = (0, 1)$ ,  $c_3 = (1, 0)$ ,  $c_4 = (1, 1)$ .

CONDITIONAL INDEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR	0.4232	0.6723
LP LR	0.4232	0.6725

CONDITIONAL DEPENDENCE		
CLASSIFIER	HAMMING LOSS	SUBSET 0/1 LOSS
BR LR	0.3470	0.5499
LP LR	0.3610	0.5146

## Rank loss

- The rank loss compares binary targets with a predicted ranking:

$$L_r(\mathbf{y}, \mathbf{h}(\mathbf{x})) = \sum_{(i,j): y_i > y_j} \left( [\![h_i(\mathbf{x}) < h_j(\mathbf{x})]\!] + \frac{1}{2} [\![h_i(\mathbf{x}) = h_j(\mathbf{x})]\!] \right)$$

- To minimize this loss, it is enough to sort the targets by their probability of relevance.
- **Theorem:** A ranking function that sorts the labels according to their probability of relevance, i.e., using the scoring function  $\mathbf{h}(\cdot)$  with

$$h_i(\mathbf{x}) = P(Y_i = 1 \mid \mathbf{x}),$$

minimizes the expected rank loss.

## Loss functions and target dependence

- Optimal (pointwise) prediction requires information about  $P(\mathbf{y} | \mathbf{x})$ .
- **Independence** simplifies learning a lot, since learning marginals is much easier than learning a joint distribution ( $\rightarrow$  graphical models).
- **Structure of the loss function** has an important influence, too, due to the “interaction” of  $\ell$  and  $P$ :

$$\mathbb{E}_{\mathbf{Y}} \ell(\mathbf{Y}, \hat{\mathbf{y}}) = \sum_{\mathbf{y}} \ell(\mathbf{y}, \hat{\mathbf{y}}) P(\mathbf{y} | \mathbf{x}).$$

In some cases, such as F-measure optimization, knowledge of properties of  $P$  instead of complete distribution is therefore enough.

- **Conditional independence** of  $P$  and **decomposability** of  $\ell$  are sufficient conditions for Bayes-optimality of target-wise Bayes predictor.

# Overview of this talk

- 1 Introduction
- 2 A unifying view on MTP problems
- 3 Loss functions in multi-target prediction
- 4 A unifying view on MTP methods
- 5 Conclusions

# A unifying view on MTP methods



Group of methods	Applicable setting
<b>Independent models</b>	B
Similarity-enforcing methods	B
Relation-exploiting methods	B, C and D
Relation-constructing methods	B and C
Representation-exploiting methods	B, C and D
Representation-constructing methods	A, B and C

A baseline method:  
learning a model for each target independently

	Mol1	Mol2	Mol3	Mol4	Mol5	Mol6	
01101		1,3	0,2	1,4	1,7	3,5	1,3
00111		2	1,7	1,5	7,5	8,2	7,6
01110		0,2	0	0,3	0,4	1,2	2,2
10001		3,1	1,1	1,3	1,1	1,7	5,2
01011		4,7	2,1	2,5	1,5	2,3	8,5
11110		?	?	?	?	?	

## A baseline method: learning a model for each target independently

	Mol1	Mol2	Mol3	Mol4	Mol5	Mol6	
01101		1,3	0,2	1,4	1,7	3,5	1,3
00111		2	1,7	1,5	7,5	8,2	7,6
01110		0,2	0	0,3	0,4	1,2	2,2
10001		3,1	1,1	1,3	1,1	1,7	5,2
01011		4,7	2,1	2,5	1,5	2,3	8,5
11110		?	?	?	?	?	?

## A baseline: Independent Models

	Mol1	Mol2	Mol3	Mol4	Mol5	Mol6	
01101		1,3	0,2	1,4	1,7	3,5	1,3
00111		2	1,7	1,5	7,5	8,2	7,6
01110		0,2	0	0,3	0,4	1,2	2,2
10001		3,1	1,1	1,3	1,1	1,7	5,2
01011		4,7	2,1	2,5	1,5	2,3	8,5
11110		?	?	?	?	?	?

## A baseline: Independent Models

Linear basis function model for  $i$ -th target:

$$f_i(\mathbf{x}) = \mathbf{a}_i^\top \phi(\mathbf{x}),$$

Solving as a joint optimization problem:

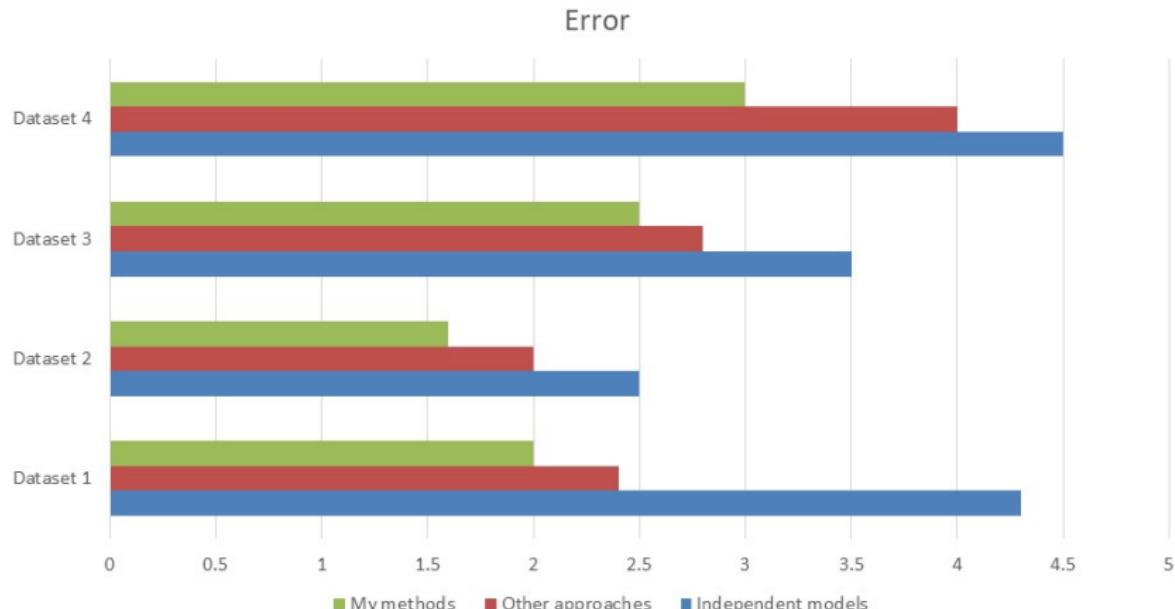
$$\min_A \|Y - XA\|_F^2 + \sum_{i=1}^m \lambda_i \|\mathbf{a}_i\|^2,$$

$$Y : (n \times m) \quad X : (n \times p) \quad A : (p \times m)$$

With the following notations:

$$X = \begin{bmatrix} \phi(\mathbf{x}_1)^T \\ \vdots \\ \phi(\mathbf{x}_n)^T \end{bmatrix} \quad A = [\mathbf{a}_1 \quad \cdots \quad \mathbf{a}_m].$$

## The results section of a typical MTP paper...



Independent models a.k.a. binary relevance, models that do not exploit target dependencies, one-versus-all, etc.

Learning a model for each target independently is still state-of-the-art in extreme multi-label classification<sup>4</sup>:

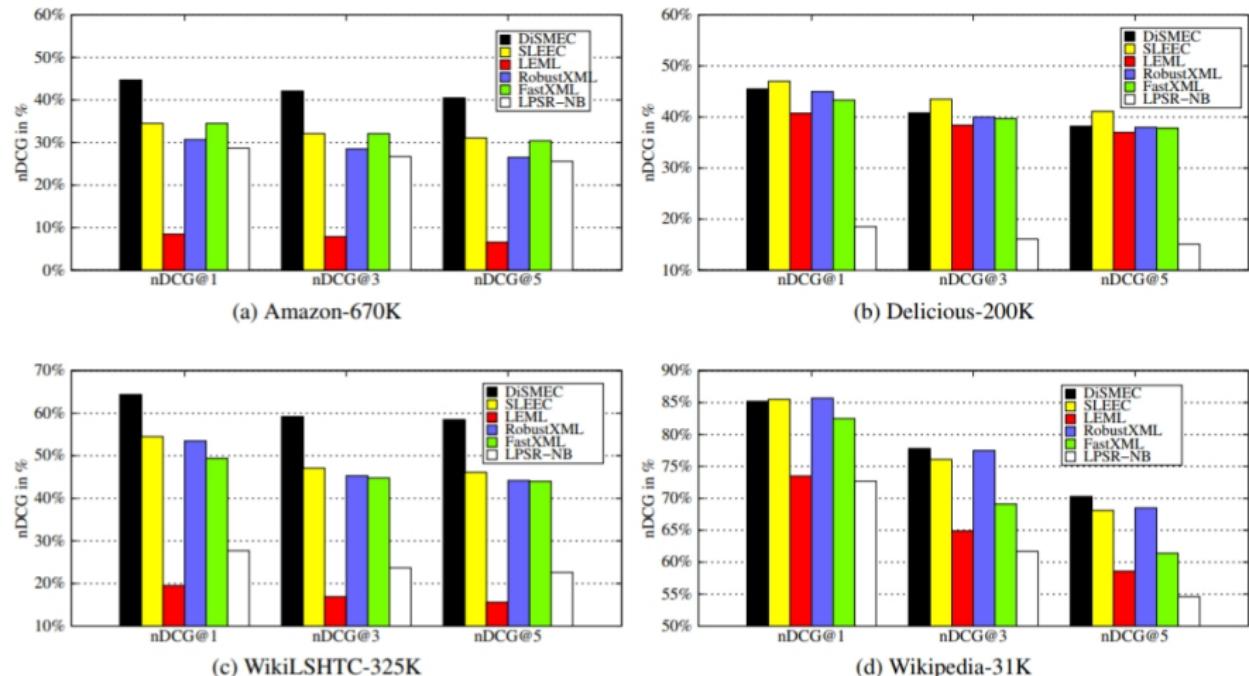


Figure 3: nDCG@k for k=1, 3 and 5

<sup>4</sup> Babbar and Schölkopf, DISMEC: Distributed Sparse Machines for Extreme Multi-label classification, WSDM 2017

# DiSMEC - Distributed Sparse Machines for XMC

**Require:** Training data  $\mathcal{T} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots (\mathbf{x}_n, \mathbf{y}_n)\}$ , input dimensionality  $p$ , label set  $\{1 \dots m\}$ ,  $B = \lfloor \frac{m}{1000} \rfloor + 1$  and pruning threshold  $\Delta$

**Ensure:** Learnt  $p \times m$  matrix  $A$  in sparse format

# DiSMEC - Distributed Sparse Machines for XMC

**Require:** Training data  $\mathcal{T} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots (\mathbf{x}_n, \mathbf{y}_n)\}$ , input dimensionality  $p$ , label set  $\{1 \dots m\}$ ,  $B = \lfloor \frac{m}{1000} \rfloor + 1$  and pruning threshold  $\Delta$

**Ensure:** Learnt  $p \times m$  matrix  $A$  in sparse format

- 1: Load single copy of input vectors  $X = \{\mathbf{x}_1 \dots \mathbf{x}_n\}$  in the main memory
- 2: Load binary sign vectors  $s_j = \{+1, -1\}_{i=1}^n$  separately for each label

# DiSMEC - Distributed Sparse Machines for XMC

**Require:** Training data  $\mathcal{T} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots (\mathbf{x}_n, \mathbf{y}_n)\}$ , input dimensionality  $p$ , label set  $\{1 \dots m\}$ ,  $B = \lfloor \frac{m}{1000} \rfloor + 1$  and pruning threshold  $\Delta$

**Ensure:** Learnt  $p \times m$  matrix  $A$  in sparse format

- 1: Load single copy of input vectors  $X = \{\mathbf{x}_1 \dots \mathbf{x}_n\}$  in the main memory
- 2: Load binary sign vectors  $s_j = \{+1, -1\}_{i=1}^n$  separately for each label
- 3: **for**  $\{b = 0; b < B; b++\}$  **do** ▷ 1st parallelization
- 4:     **#pragma omp parallel for private(j)** ▷ 2nd parallelization
- 5:     **for**  $\{j = b \times 1000; j \leq (b + 1) \times 1000; j++\}$  **do**
- 6:         Using  $(X, s_j)$ , train weight vector  $a_j$  on a single core
- 7:         **Prune ambiguous weights** in  $a_j$
- 8:     **return**  $p \times 1000$  matrix  $A$
- 9: **return**  $p \times m$  weight matrix  $A$

# DiSMEC - Distributed Sparse Machines for XMC

**Require:** Training data  $\mathcal{T} = \{(\mathbf{x}_1, \mathbf{y}_1) \dots (\mathbf{x}_n, \mathbf{y}_n)\}$ , input dimensionality  $p$ , label set  $\{1 \dots m\}$ ,  $B = \lfloor \frac{m}{1000} \rfloor + 1$  and pruning threshold  $\Delta$

**Ensure:** Learnt  $p \times m$  matrix  $A$  in sparse format

- 1: Load single copy of input vectors  $X = \{\mathbf{x}_1 \dots \mathbf{x}_n\}$  in the main memory
- 2: Load binary sign vectors  $s_j = \{+1, -1\}_{i=1}^n$  separately for each label
- 3: **for**  $\{b = 0; b < B; b++\}$  **do** ▷ 1st parallelization
- 4:     **#pragma omp parallel for private(j)** ▷ 2nd parallelization
- 5:     **for**  $\{j = b \times 1000; j \leq (b + 1) \times 1000; j++\}$  **do**
- 6:         Using  $(X, s_j)$ , train weight vector  $a_j$  on a single core
- 7:         **Prune ambiguous weights** in  $a_j$
- 8:     **return**  $p \times 1000$  matrix  $A$
- 9: **return**  $p \times m$  weight matrix  $A$

- Learns model for LSHTCWiki-325K in 6 hours on 400 cores
- Model size is 3GB due to pruning step

# Distribution of Learnt Weights

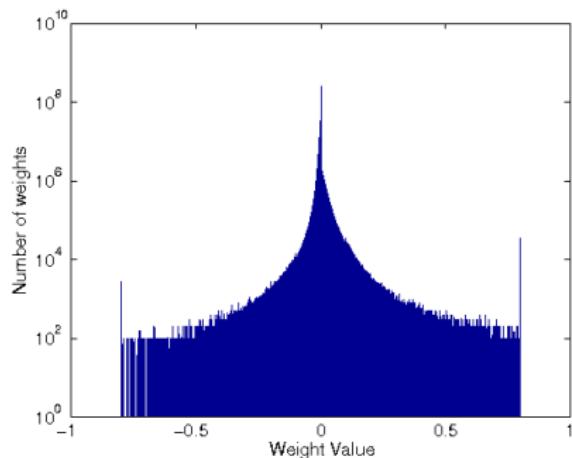


Figure: Before pruning

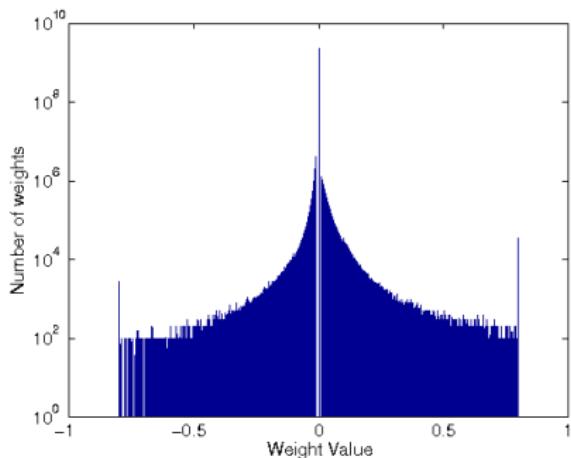


Figure: After pruning small weights

- Of the 3 Billion weights, 97% are s.t.,  $A_{ij} \leq |0.01|$ , and hence non-discriminative
- Storing them leads to large model sizes but no benefit in classification

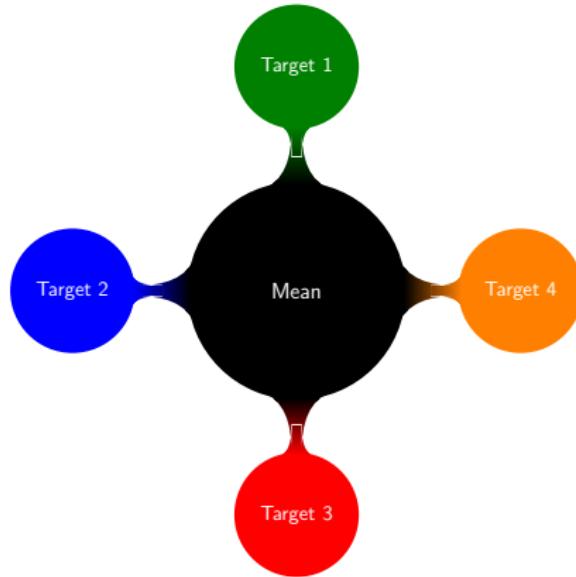
# A unifying view on MTP methods



Group of methods	Applicable setting
Independent models	B
<b>Similarity-enforcing methods</b>	B
Relation-exploiting methods	B, C and D
Relation-constructing methods	B and C
Representation-exploiting methods	B, C and D
Representation-constructing methods	A, B and C

# Mean-regularized multi-task learning<sup>5</sup>

- **Simple assumption:** models for different targets are related to each other.
- **Simple solution:** the parameters of these models should have similar values.
- **Approach:** bias the parameter vectors towards their mean vector.



$$\min_A \|Y - XA\|_F^2 + \lambda \sum_{i=1}^m \left\| \mathbf{a}_i - \frac{1}{m} \sum_{j=1}^m \mathbf{a}_j \right\|^2,$$

<sup>5</sup> Evgeniou and Pontil, Regularized multi-task learning, KDD 2004.

## Joint feature selection

- Enforce that the same features are selected for different targets<sup>6</sup>:

$$\min_A ||Y - XA||_F^2 + \lambda \sum_{j=1}^p ||\mathbf{a}_j||^2$$

---

<sup>6</sup> Obozinski et al. Joint covariate selection and joint subspace selection for multiple classification problems. Statistics and Computing 2010

## Joint feature selection

- Enforce that the same features are selected for different targets<sup>6</sup>:

$$\min_A \|Y - XA\|_F^2 + \lambda \sum_{j=1}^p \|\mathbf{a}_j\|^2$$

- The vectors  $\mathbf{a}_j$  now represent the columns of matrix  $A^T$ :

L1-norm per target:

$$A^T = \begin{bmatrix} a_{11} & \cdots & a_{1p} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mp} \end{bmatrix} \in \diamond$$

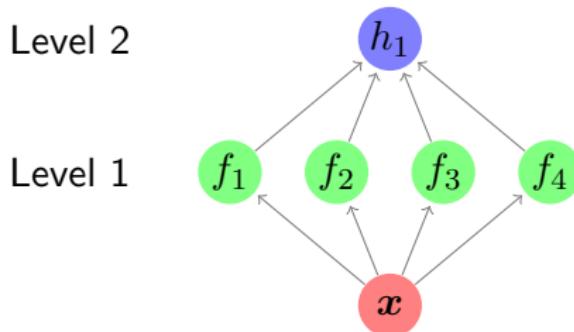
Mix L1/L2-norm per target:

$$A^T = \begin{bmatrix} a_{11} & \cdots & a_{1p} \\ \vdots & & \vdots \\ a_{m1} & \cdots & a_{mp} \end{bmatrix} \in \diamond$$
$$(\mathbf{a}_1, \dots, \mathbf{a}_p) \in \diamond$$

<sup>6</sup> Obozinski et al. Joint covariate selection and joint subspace selection for multiple classification problems. Statistics and Computing 2010

## Stacking (Stacked generalization)

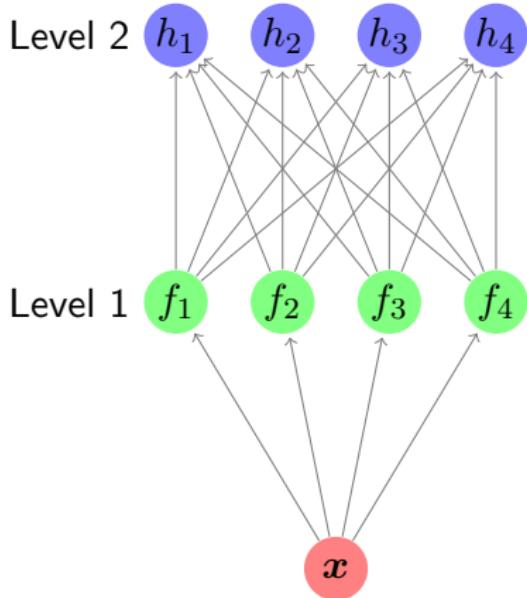
- Originally introduced as a general ensemble learning or blending technique.<sup>7</sup>
- Level 1 classifiers: apply a series of ML methods on the same dataset (or, one ML method on bootstrap samples of the dataset)
- Level 2 classifier: apply an ML method to a new dataset consisting of the predictions obtaining at Level 1



<sup>7</sup> Wolpert, Stacked generalization. Neural Networks 1992

# Stacking applied to multi-target prediction<sup>8</sup>

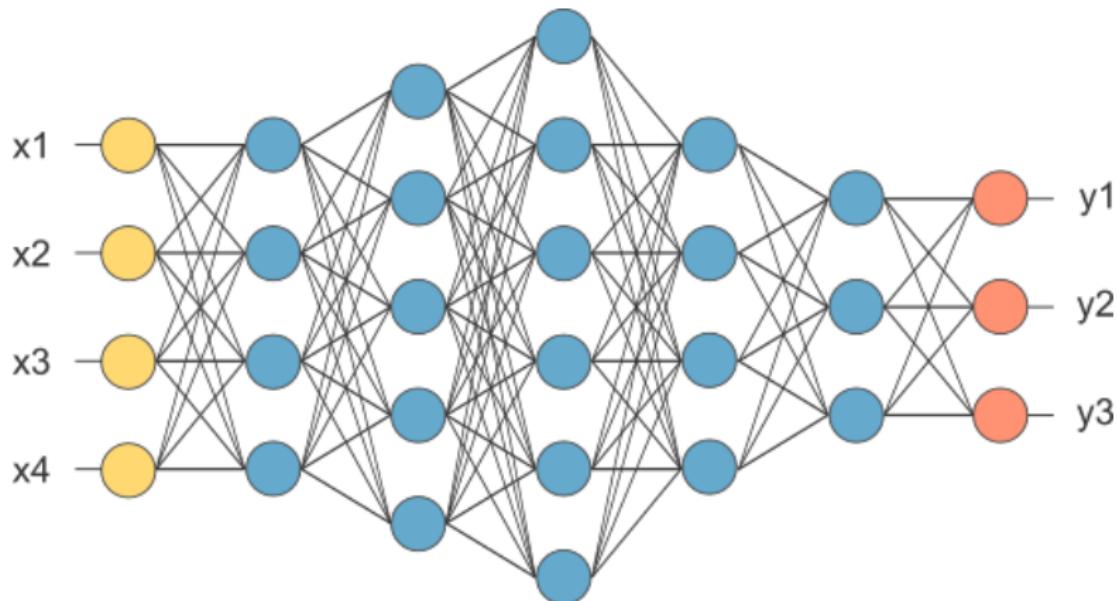
- Level 1 classifiers: learn a model for every target independently
- Level 2 classifier: learn again a model for every target independently, using the predictions of the first step as features



<sup>8</sup> Cheng and Hüllermeier, Combining Instance-based learning and Logistic Regression for Multi-Label classification, Machine Learning, 2009

# Enforcing similarity in (Deep) Neural Networks

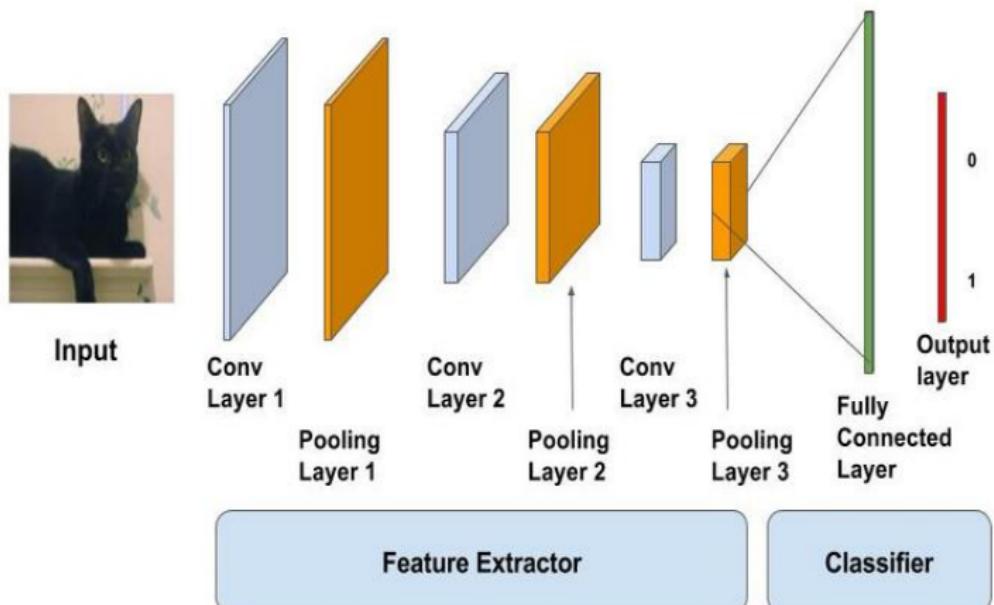
Commonly-used architecture: weight sharing among targets<sup>9</sup>



<sup>9</sup> Caruana, Multitask learning: A knowledge-based source of inductive bias. Machine Learning 1997

# Re-using Pretrained Models in (Deep) Neural Networks

Commonly-used training method: first train on targets that have a lot of observations, only train some parameters for targets that have few observations<sup>10</sup>



<sup>10</sup> Keras Tutorial: Transfer Learning using pre-trained models



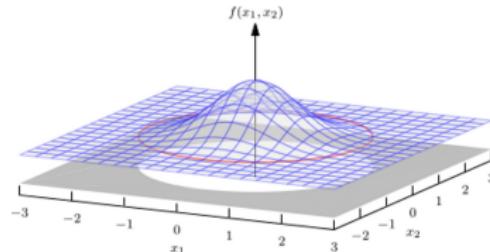
## Question

In which situations are similarity-enforcing models capable of outperforming independent models w.r.t. predictive performance?

- Always
- When  $p$  is sufficiently large
- When  $m$  is sufficiently large
- When the targets are sufficiently correlated

## An intuitive explanation: James-Stein estimation

- Consider a sample of a multivariate normal distribution  $\mathbf{y} \sim N(\boldsymbol{\theta}, \sigma^2 \mathbf{I})$ .



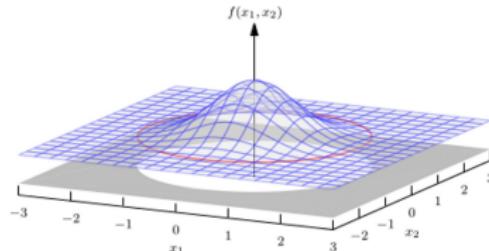
- What is the best estimator of the mean vector  $\boldsymbol{\theta}$ ?
- Evaluation w.r.t. MSE:  $\mathbb{E}[(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^2]$

---

<sup>11</sup>W. James and C. Stein. Estimation with quadratic loss. In Proc. Fourth Berkeley Symp. Math. Statist. Prob. 1, pages 361-379, 1961

## An intuitive explanation: James-Stein estimation

- Consider a sample of a multivariate normal distribution  $\mathbf{y} \sim N(\boldsymbol{\theta}, \sigma^2 \mathbf{I})$ .



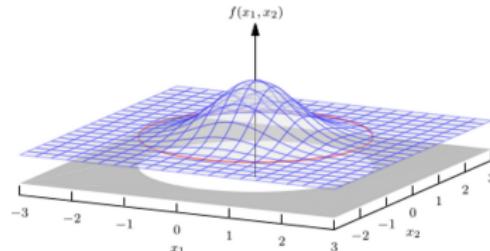
- What is the best estimator of the mean vector  $\boldsymbol{\theta}$ ?
- Evaluation w.r.t. MSE:  $\mathbb{E}[(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^2]$
- Single-observation maximum likelihood estimator:  $\hat{\boldsymbol{\theta}}^{\text{ML}} = \mathbf{y}$

---

<sup>11</sup> W. James and C. Stein. Estimation with quadratic loss. In Proc. Fourth Berkeley Symp. Math. Statist. Prob. 1, pages 361-379, 1961

## An intuitive explanation: James-Stein estimation

- Consider a sample of a multivariate normal distribution  $\mathbf{y} \sim N(\boldsymbol{\theta}, \sigma^2 \mathbf{I})$ .

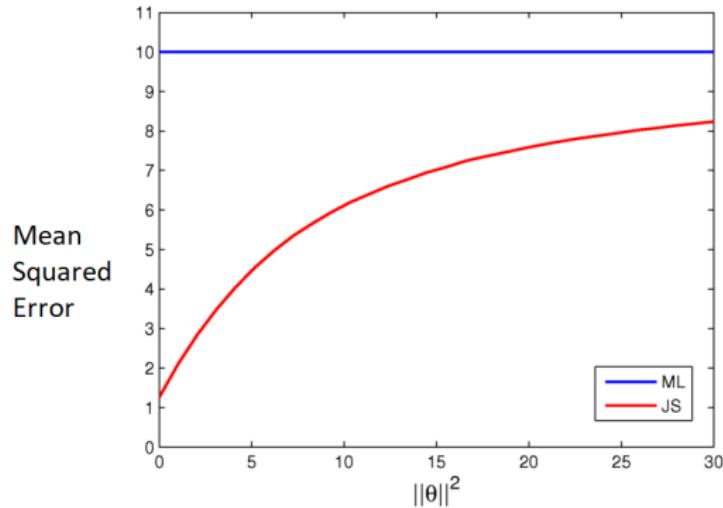


- What is the best estimator of the mean vector  $\boldsymbol{\theta}$ ?
- Evaluation w.r.t. MSE:  $\mathbb{E}[(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}})^2]$
- Single-observation maximum likelihood estimator:  $\hat{\boldsymbol{\theta}}^{\text{ML}} = \mathbf{y}$
- James-Stein estimator<sup>11</sup>:

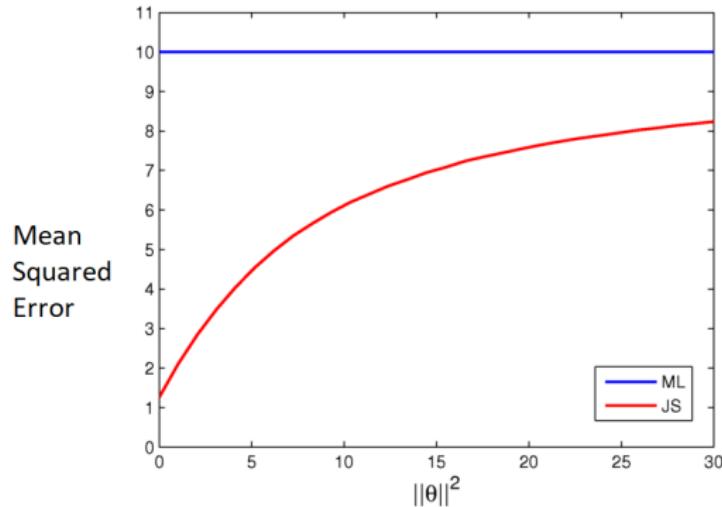
$$\hat{\boldsymbol{\theta}}^{\text{JS}} = \left(1 - \frac{(m-2)\sigma^2}{\|\mathbf{y}\|^2}\right) \mathbf{y}$$

<sup>11</sup>W. James and C. Stein. Estimation with quadratic loss. In Proc. Fourth Berkeley Symp. Math. Statist. Prob. 1, pages 361-379, 1961

- Works best when the norm of the mean vector is close to zero:



- Works best when the norm of the mean vector is close to zero:



- Regularization towards other directions is also possible:

$$\hat{\theta}^{\text{JS+}} = \left(1 - \frac{(m-2)\sigma^2}{\|\mathbf{y} - \mathbf{v}\|^2}\right) (\mathbf{y} - \mathbf{v}) + \mathbf{v}$$

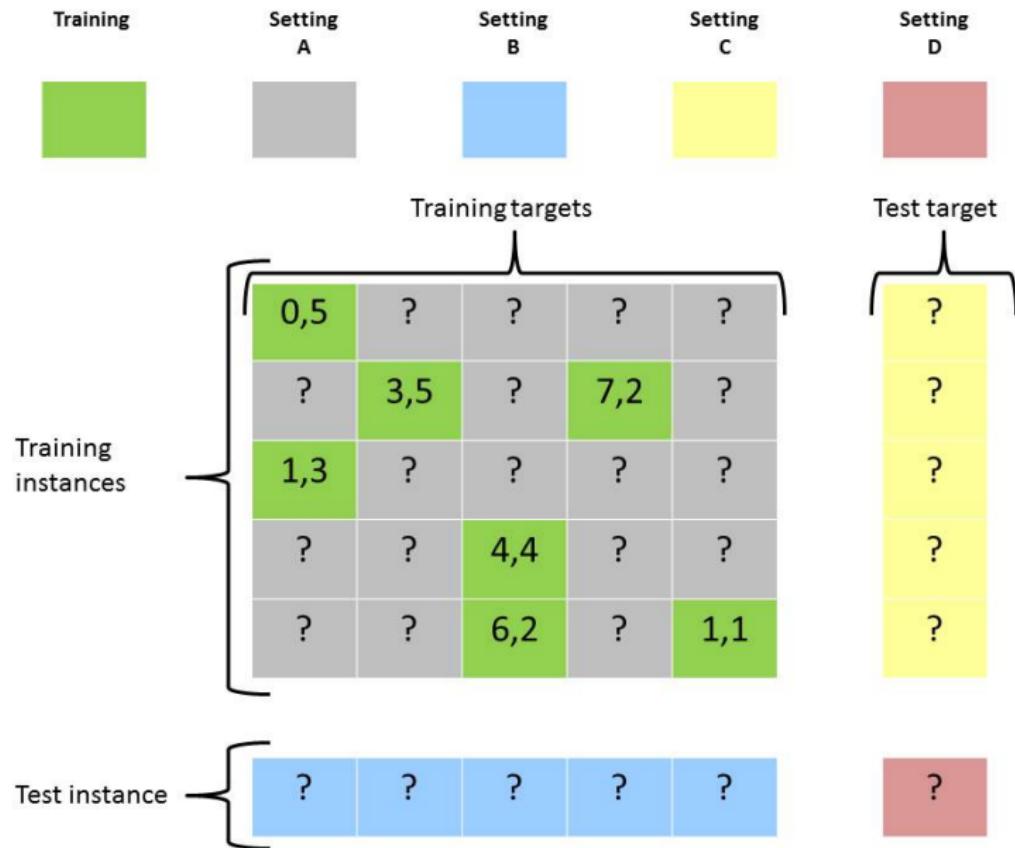
- Only outperforms the maximum likelihood estimator w.r.t. the sum of squared errors over all components.

# A unifying view on MTP methods



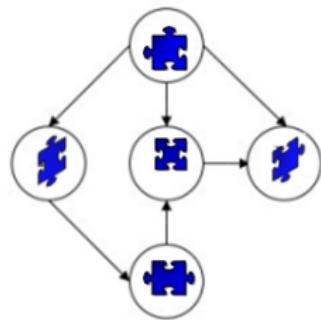
Group of methods	Applicable setting
Independent models	B
Similarity-enforcing methods	B
<b>Relation-exploiting methods</b>	B, C and D
Relation-constructing methods	B and C
Representation-exploiting methods	B, C and D
Representation-constructing methods	A, B and C

# Different learning settings revisited

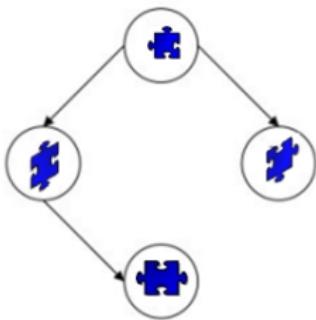


## Exploiting relations in regularization terms

Graph



Tree



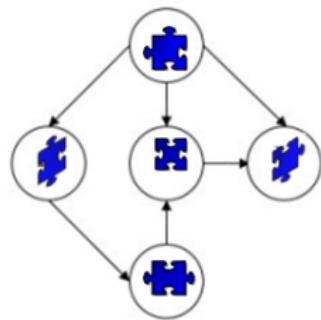
Similarity

	1	2	3	4
1	1	0.26	0.26	0.04
2	0.26	1	0.7	0.57
3	0.26	0.7	1	0.44
4	0.04	0.57	0.44	1

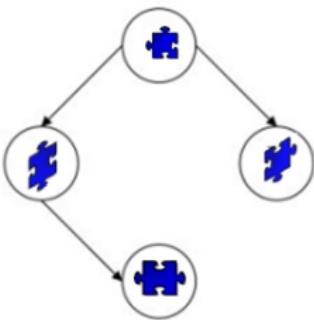
<sup>12</sup> Gopal and Yang, Recursive regularization for large-scale classification with hierarchical and graphical dependencies, KDD 2013

## Exploiting relations in regularization terms

Graph



Tree



Similarity

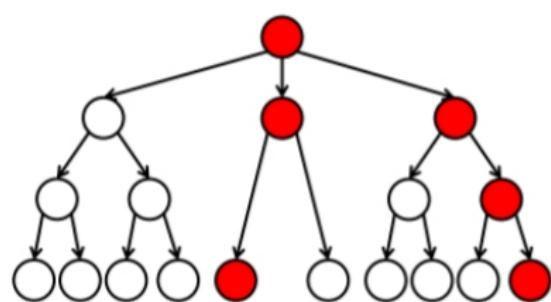
	1	2	3	4
1	1	0.26	0.26	0.04
2	0.26	1	0.7	0.57
3	0.26	0.7	1	0.44
4	0.04	0.57	0.44	1

Graph-based regularization is an approach that can be applied to the three types of relations<sup>12</sup>:

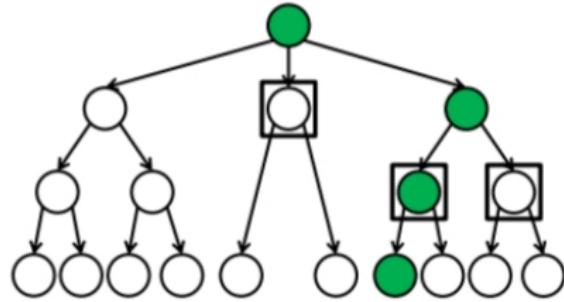
$$\min_A ||Y - XA||_F^2 + \lambda \sum_{i=1}^m \sum_{j \in \mathcal{N}(i)} ||\mathbf{a}_i - \mathbf{a}_j||^2$$

<sup>12</sup> Gopal and Yang, Recursive regularization for large-scale classification with hierarchical and graphical dependencies, KDD 2013

## Hierarchical multi-label classification



(a) Ground truth.



(b) Prediction A.

In addition to performance gains in general, hierarchies can also be used to define specific loss functions, such as the H-loss<sup>13</sup>:

$$L_H(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{j:y_j \neq \hat{y}_j} c_j I(\text{anc}(y_j) = \text{anc}(\hat{y}_j))$$

$c_i$  depends on the depth of node  $i$

<sup>13</sup> Bi and Kwok, Bayes-optimal hierarchical multi-label classification, IEEE Transactions on Knowledge and Data Engineering, 2014

## Exploiting similarity measures among targets

	1	0.26	0.26	0.04	
	0.26	1	0.7	0.57	
	0.26	0.7	1	0.44	
	0.04	0.57	0.44	1	

Can be done within the framework of vector-valued kernel functions<sup>14</sup>:

$$f(\mathbf{x}, \mathbf{t}) = \mathbf{w}^T \Psi(\mathbf{x}, \mathbf{t}) = \sum_{(\bar{\mathbf{x}}, \bar{\mathbf{t}}) \in \mathcal{D}} \alpha_{(\bar{\mathbf{x}}, \bar{\mathbf{t}})} \Gamma((\mathbf{x}, \mathbf{t}), (\bar{\mathbf{x}}, \bar{\mathbf{t}}))$$

Model the joint kernel as a product of an instance kernel  $k(\cdot, \cdot)$  and a target kernel  $g(\cdot, \cdot)$ :

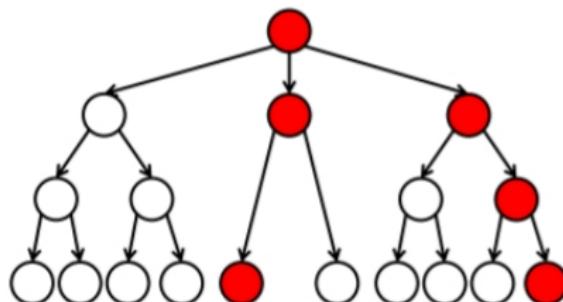
$$\Gamma((\mathbf{x}, \mathbf{t}), (\bar{\mathbf{x}}, \bar{\mathbf{t}})) = k(\mathbf{x}, \bar{\mathbf{x}}) \cdot g(\mathbf{t}, \bar{\mathbf{t}})$$

---

<sup>14</sup> Alvarez et al., Kernels for vector-valued functions: a review, Foundation and Trends in Machine Learning, 2012

# Converting graphs to similarities or target representations

- **Similarities:** use graph structure to express target similarities
  - e.g. the shortest-path kernel between two nodes
- **Representations:** often characteristics of a specific vertex or edge
  - e.g. the number of positive labels that are siblings of a vertex<sup>15</sup>



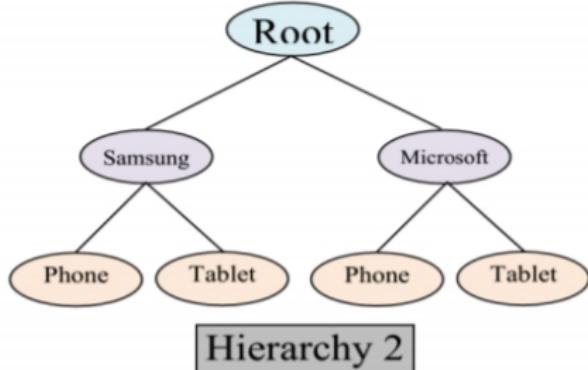
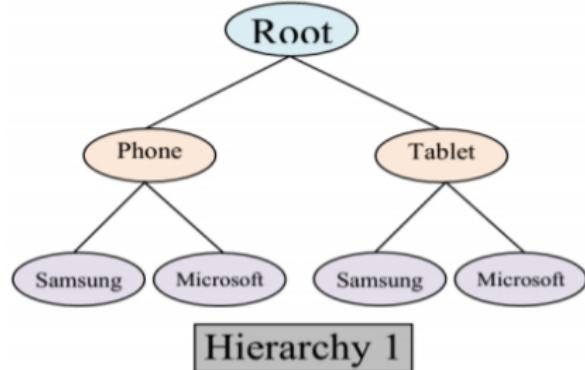
<sup>15</sup> Rousu et al., Kernel-based learning of hierarchical multilabel classification models, JMLR 2006

# A unifying view on MTP methods



Group of methods	Applicable setting
Independent models	B
Similarity-enforcing methods	B
Relation-exploiting methods	B, C and D
<b>Relation-constructing methods</b>	B and C
Representation-exploiting methods	B, C and D
Representation-constructing methods	A, B and C

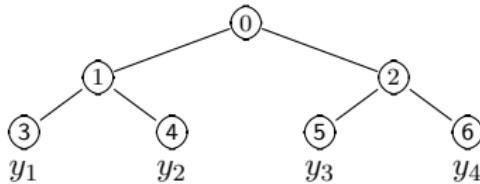
# Constructing target hierarchies



- It might be difficult for a human expert to define a hierarchy<sup>16</sup>
- Perhaps one can try to learn the hierarchy from data?
- Algorithms: level flattening, node removal, hierarchy modification, hierarchy generation, etc.

<sup>16</sup> Rangwala and Naik, Tutorial on Large-Scale Hierarchical Classification, KDD 2017.

## Label trees ( $\neq$ decision trees)



- Organize classifiers in a tree structure (one leaf  $\Leftrightarrow$  one label)
- Mainly used in multi-class and multi-label classification
- Goal is fast prediction: almost logarithmic in the number of labels
- Algorithms: Label embedding trees<sup>17</sup>, Nested dichotomies<sup>18</sup>, Conditional probability trees<sup>19</sup>, Hierarchical softmax<sup>20</sup>, FastText<sup>21</sup>, Probabilistic classifier chains<sup>22</sup>

---

<sup>17</sup> Bengio et al., Label embedding trees for large multi-class tasks, NIPS 2010

<sup>18</sup> Frank and Kramer, Ensembles of nested dichotomies for multi-class problems, ICML 2004

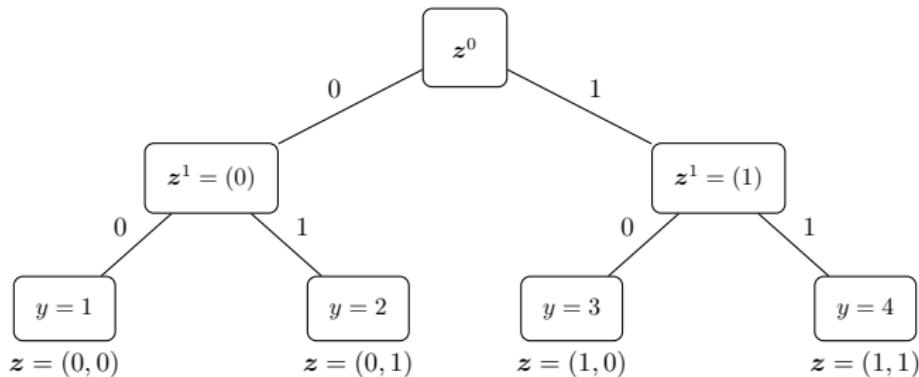
<sup>19</sup> Beygelzimer et al., Conditional probability tree estimation analysis and algorithms. UAI 2009

<sup>20</sup> Morin and Bengio, Hierarchical probabilistic neural network language model, AISTATS 2005

<sup>21</sup> Joulin et al., Bag of tricks for efficient text classification. CoRR, abs/1607.01759, 2016

<sup>22</sup> Dembczynski et al., Bayes optimal multilabel classification via probabilistic classifier chains, ICML 2010

# Hierarchical softmax / Probabilistic classifier trees



- Encode the targets by a **prefix code** ( $\Rightarrow$  tree structure)<sup>23</sup>
- Multi-class classification: each label  $y$  **coded** by  $z = (z_1, \dots, z_l) \in \mathcal{C}$
- Multi-label classification: a label vector  $y = (y_1, \dots, y_m)$  is a prefix code.

<sup>23</sup> Dembczynski et al., Consistency of probabilistic classifier trees. ECMLPKDD 2016

## Probabilistic classifier chains

- Estimate the joint conditional distribution  $P(\mathbf{Y} \mid \mathbf{x})$ .
- For optimizing the subset 0/1 loss:

$$\ell_{0/1}(\mathbf{y}, \hat{\mathbf{y}}) = \llbracket \mathbf{y} \neq \hat{\mathbf{y}} \rrbracket$$

- Repeatedly apply the **product rule of probability**:

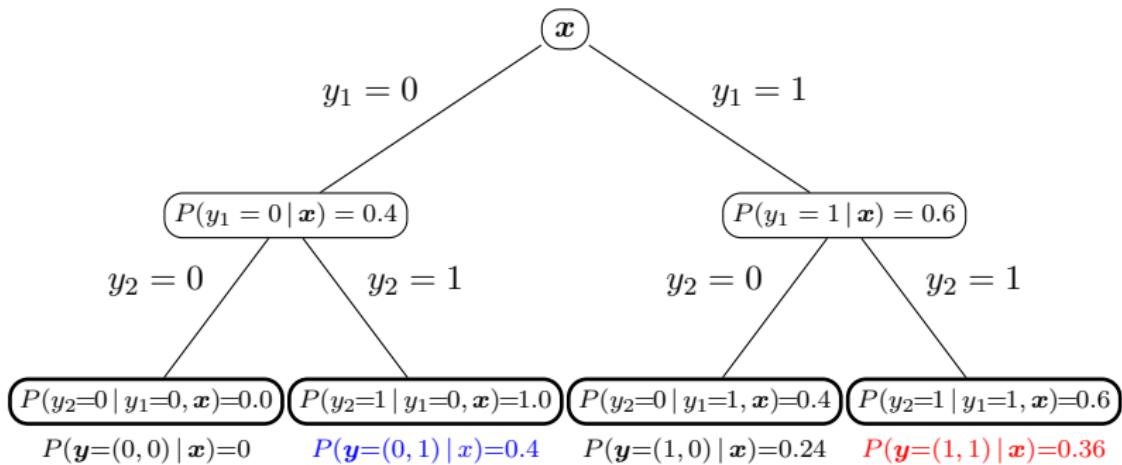
$$P(\mathbf{Y} = \mathbf{y} \mid \mathbf{x}) = \prod_{i=1}^m P(Y_i = y_i \mid \mathbf{x}, y_1, \dots, y_{i-1}).$$

- Learning relies on constructing probabilistic classifiers for estimating

$$P(Y_i = y_i \mid \mathbf{x}, y_1, \dots, y_{i-1}),$$

independently for each  $i = 1, \dots, m$ .

- Inference relies on exploiting a probability tree:



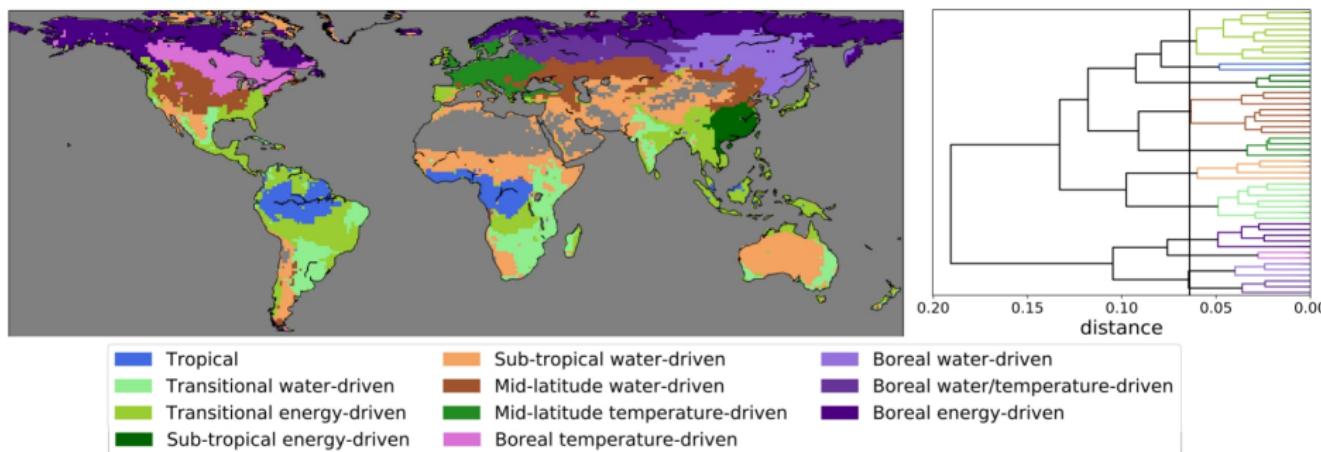
- For subset 0/1 loss one needs to find  $\mathbf{h}(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} P(\mathbf{y} | \mathbf{x})$ .
- Greedy and approximate search techniques with guarantees exist.<sup>24</sup>
- Other losses: compute the prediction on a sample from  $P(\mathbf{Y} | \mathbf{x})$ .<sup>25</sup>

<sup>24</sup> Kumar et al., Beam search algorithms for multilabel learning, Machine Learning 2013

<sup>25</sup> Dembczynski et al., An analysis of chaining in multi-label classification, ECAI 2012

# Constructing hierarchies to obtain additional insight

- Application in climate science
- Result of learning 20000 tasks simultaneously with a multi-task learning method
- Followed by hierarchical clustering of the learned weight vectors<sup>26</sup>:

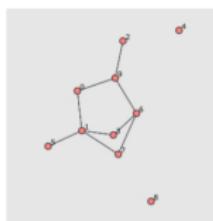


<sup>26</sup> Papagiannopoulou et al. Global hydro-climatic biomes identified with multi-task learning, Geoscientific Model Development Discussions 2018

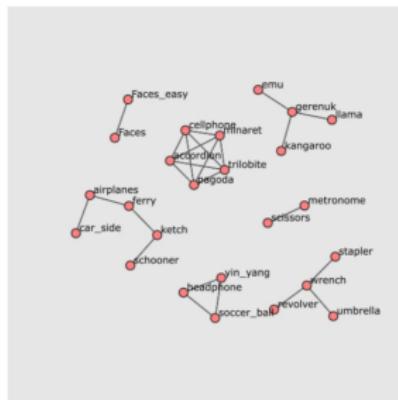
# Constructing target similarities by output kernel learning

- Consider models  $f : \mathcal{X} \rightarrow \mathbb{R}^m$
- Training dataset  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$
- Learnable correlation matrix  $\Gamma$  between targets
- Learn output kernel and model parameters jointly<sup>27</sup>:

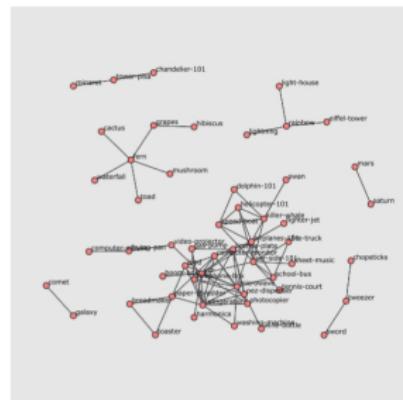
$$\min_{\Gamma \in \mathbb{R}^{m \times m}} \left[ \min_{f \in \mathcal{F}} \sum_{i=1}^n \frac{\|f(\mathbf{x}_i) - \mathbf{y}_i\|_2^2}{2\lambda} + \frac{\|f\|_{\mathcal{F}}^2}{2} + \frac{\|\Gamma\|_{\mathcal{F}}^2}{2} \right]$$



(a) USPS digits



(b) Caltech 101



(c) Caltech 256

<sup>27</sup> Dinuzzo et al., Learning Output Kernels with Block Coordinate Descent, ICML 2011

# Constructing decision rules among targets

**Table 2** Extended multi-label WEATHER dataset

Outlook	Temperature	Humidity	Windy	Play	Icecream	Tea	Lemonade	Dontplay
Rainy	65	70	Yes	0	0	1	0	1
Rainy	71	91	Yes	0	0	1	0	1
Sunny	85	85	No	0	1	0	1	1
Sunny	80	90	Yes	0	1	0	1	1
Sunny	72	95	No	0	1	0	1	1
Sunny	69	70	No	1	0	0	1	0
Sunny	75	70	Yes	1	0	0	1	0
Overcast	83	86	No	1	0	0	1	0
Overcast	64	65	Yes	1	0	0	1	0
Overcast	72	90	Yes	1	0	0	1	0
Overcast	81	75	No	1	0	0	1	0
Rainy	70	96	No	1	0	0	1	0
Rainy	68	80	No	1	0	0	1	0
Rainy	75	80	No	1	0	0	1	0

<sup>28</sup> Loza-Mencia and Janssen, Learning rules for multi-label classification: a stacking and separate-and-conquer approach

## Constructing decision rules among targets

**Table 2** Extended multi-label WEATHER dataset

Outlook	Temperature	Humidity	Windy	Play	Icecream	Tea	Lemonade	Dontplay
Rainy	65	70	Yes	0	0	1	0	1
Rainy	71	91	Yes	0	0	1	0	1
Sunny	85	85	No	0	1	0	1	1
Sunny	80	90	Yes	0	1	0	1	1
Sunny	72	95	No	0	1	0	1	1
Sunny	69	70	No	1	0	0	1	0
Sunny	75	70	Yes	1	0	0	1	0
Overcast	83	86	No	1	0	0	1	0
Overcast	64	65	Yes	1	0	0	1	0
Overcast	72	90	Yes	1	0	0	1	0
Overcast	81	75	No	1	0	0	1	0
Rainy	70	96	No	1	0	0	1	0
Rainy	68	80	No	1	0	0	1	0
Rainy	75	80	No	1	0	0	1	0

Potential inferred rule<sup>28</sup>: **TEA → NOT LEMONADE**

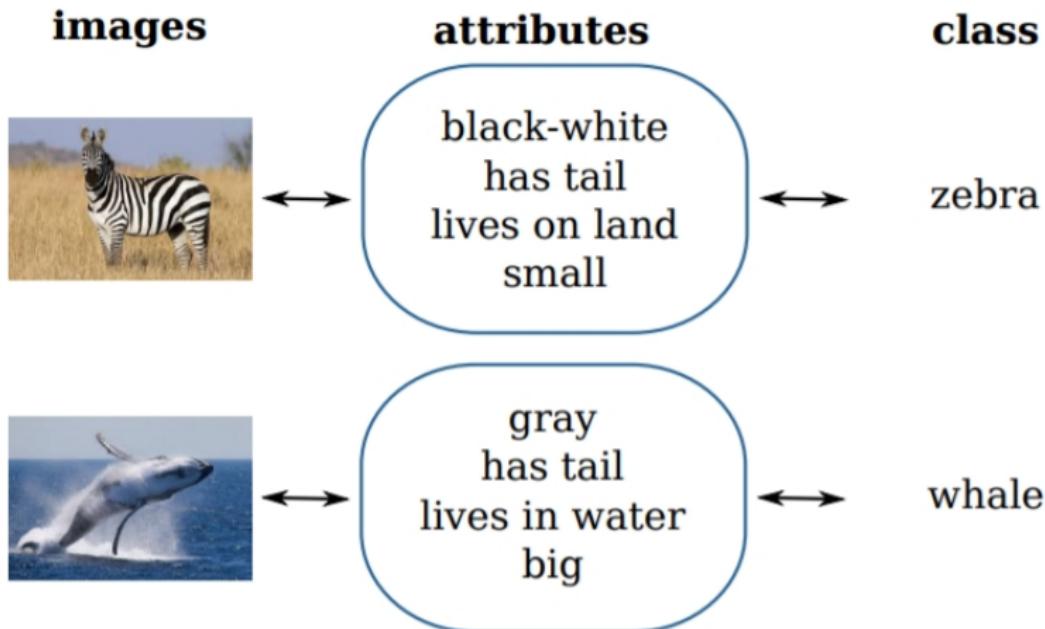
<sup>28</sup> Loza-Mencia and Janssen, Learning rules for multi-label classification: a stacking and separate-and-conquer approach

# A unifying view on MTP methods



Group of methods	Applicable setting
Independent models	B
Similarity-enforcing methods	B
Relation-exploiting methods	B, C and D
Relation-constructing methods	B and C
<b>Representation-exploiting methods</b>	B, C and D
Representation-constructing methods	A, B and C

# A target representation in computer vision



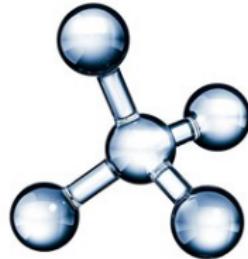
Target representations are the key element of zero-shot learning methods<sup>29</sup>

<sup>29</sup> Examples taken from the CVPR 2016 Tutorial on Zero-shot learning for Computer Vision

# Target representations can take many forms



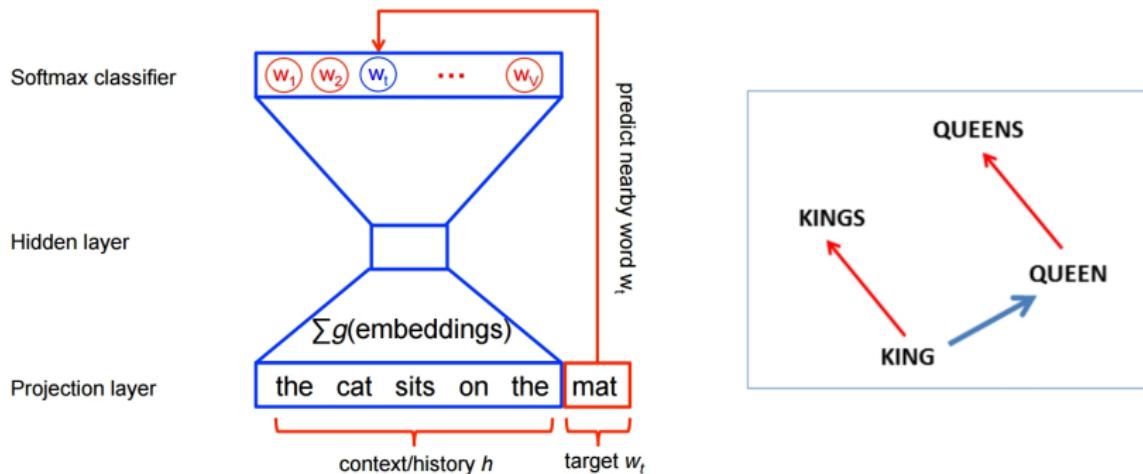
A T G G C T A C  
T A C G G A T A



# Learning target embeddings from text: Word2Vec

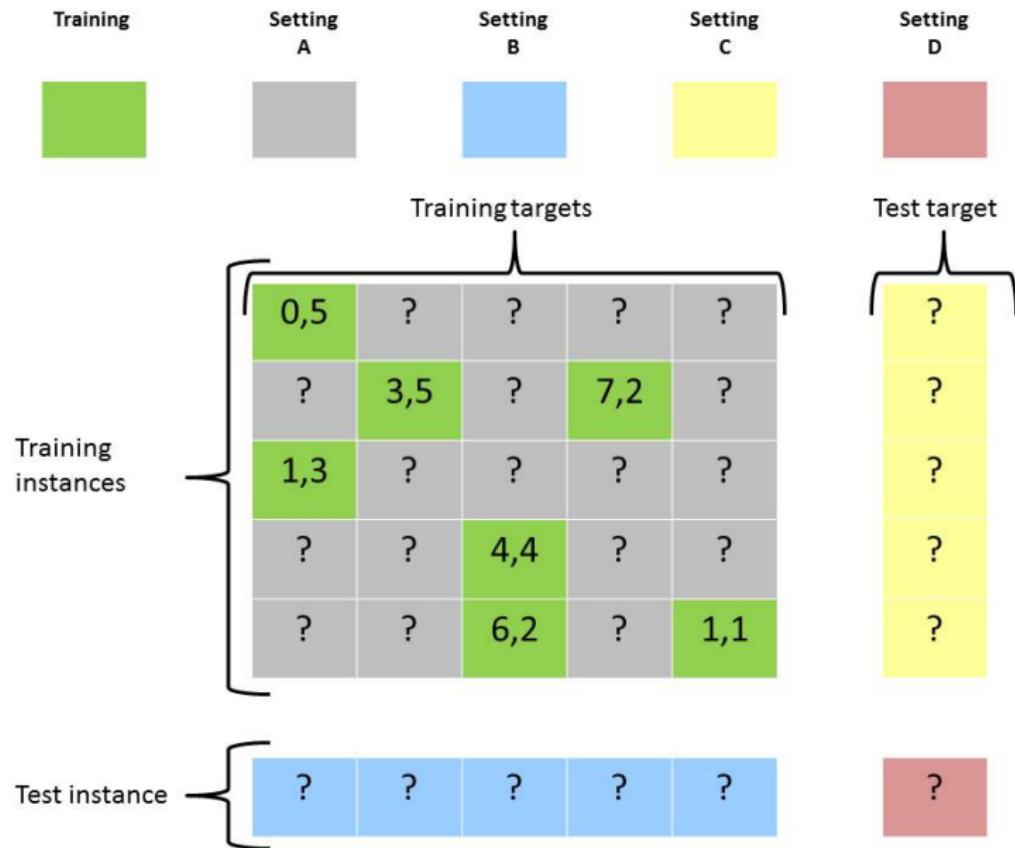
Predict the probability of the next word  $w_t$  given the previous words  $h^{30}$ :

$$P(w_t | h) = \frac{\exp(f(w_t, h))}{\sum_{\text{allwords}} \exp(f(w_t, h))}$$



<sup>30</sup> Mikolov et al., Efficient Estimation of Word Representations in Vector Space, Arxiv 2013

# Different learning settings revisited



# Kronecker kernel ridge regression

Pairwise model representation in the primal:

$$f(\mathbf{x}, \mathbf{t}) = \mathbf{w}^T (\phi(\mathbf{x}) \otimes \psi(\mathbf{t}))$$

Kronecker product pairwise kernel in the dual<sup>31</sup>:

$$f(\mathbf{x}, \mathbf{t}) = \sum_{(\bar{\mathbf{x}}, \bar{\mathbf{t}}) \in \mathcal{D}} \alpha_{(\bar{\mathbf{x}}, \bar{\mathbf{t}})} k(\mathbf{x}, \bar{\mathbf{x}}) \cdot g(\mathbf{t}, \bar{\mathbf{t}}) = \sum_{(\bar{\mathbf{x}}, \bar{\mathbf{t}}) \in \mathcal{D}} \alpha_{(\bar{\mathbf{x}}, \bar{\mathbf{t}})} \Gamma((\mathbf{x}, \mathbf{t}), (\bar{\mathbf{x}}, \bar{\mathbf{t}}))$$

Least-squares minimization with  $\mathbf{z} = \text{vec}(Y)$ :

$$\min_{\boldsymbol{\alpha}} \|\boldsymbol{\Gamma}\boldsymbol{\alpha} - \mathbf{z}\|_2^2 + \lambda \boldsymbol{\alpha}^\top \boldsymbol{\Gamma} \boldsymbol{\alpha}$$

---

<sup>31</sup> Stock et al., A comparative study of pairwise learning methods based on kernel ridge regression, Neural Computation 2018

## Two-step zero-shot learning<sup>32 33</sup>

	Mol1	Mol2	Mol3	Mol4	Mol5	Mol6	Mol7	
01101	?	1,3	0,2	1,4	1,7	3,5	1,3	?
00111	?	2	1,7	1,5	7,5	8,2	7,6	?
01110	?	0,2	0	0,3	0,4	1,2	2,2	?
10001	?	3,1	1,1	1,3	1,1	1,7	5,2	?
01011	?	4,7	2,1	2,5	1,5	2,3	8,5	?
11110	?	?	?	?	?	?	?	

<sup>32</sup> Pahikkala et al. A two-step approach for solving full and almost full cold-start problems in dyadic prediction, ECML/PKDD 2014.

<sup>33</sup> Romero-Paredes and Torr, An embarrassingly simple approach to zero-shot learning, ICML 2015.

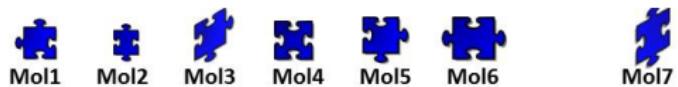
# Two-step zero-shot learning<sup>34 35</sup>

	Mol1	Mol2	Mol3	Mol4	Mol5	Mol6	
01101		1,3	0,2	1,4	1,7	3,5	1,3
00111		2	1,7	1,5	7,5	8,2	7,6
01110		0,2	0	0,3	0,4	1,2	2,2
10001		3,1	1,1	1,3	1,1	1,7	5,2
01011		4,7	2,1	2,5	1,5	2,3	8,5
11110		1,2	2,1	1,7	4,3	2,4	2,5

<sup>34</sup> Pahikkala et al. A two-step approach for solving full and almost full cold-start problems in dyadic prediction, ECML/PKDD 2014.

<sup>35</sup> Romero-Paredes and Torr, An embarrassingly simple approach to zero-shot learning, ICML 2015.

# Two-step zero-shot learning<sup>36 37</sup>



1,3	0,2	1,4	1,7	3,5	1,3	1,2
2	1,7	1,5	7,5	8,2	7,6	1,4
0,2	0	0,3	0,4	1,2	2,2	3,8
3,1	1,1	1,3	1,1	1,7	5,2	1,1
4,7	2,1	2,5	1,5	2,3	8,5	1,5

1,2	2,1	1,7	4,3	2,4	2,5	4,3
-----	-----	-----	-----	-----	-----	-----

<sup>36</sup> Pahikkala et al. A two-step approach for solving full and almost full cold-start problems in dyadic prediction, ECML/PKDD 2014.

<sup>37</sup> Romero-Paredes and Torr, An embarrassingly simple approach to zero-shot learning, ICML 2015.

## Two-step kernel ridge regression

- Kernel evaluations for new test instance:

$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))^T$$

$$\mathbf{g}(\mathbf{t}) = (g(\mathbf{t}, \mathbf{t}_1), \dots, g(\mathbf{t}, \mathbf{t}_q))^T$$

## Two-step kernel ridge regression

- Kernel evaluations for new test instance:

$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))^T$$

$$\mathbf{g}(\mathbf{t}) = (g(\mathbf{t}, \mathbf{t}_1), \dots, g(\mathbf{t}, \mathbf{t}_q))^T$$

- Step 1: prediction for  $\mathbf{x}$  on all the training targets

$$\mathbf{f}_T(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T A^{IT} = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda_d \mathbf{I})^{-1} Y$$

## Two-step kernel ridge regression

- Kernel evaluations for new test instance:

$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))^T$$

$$\mathbf{g}(\mathbf{t}) = (g(\mathbf{t}, \mathbf{t}_1), \dots, g(\mathbf{t}, \mathbf{t}_q))^T$$

- Step 1: prediction for  $\mathbf{x}$  on all the training targets

$$\mathbf{f}_T(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T A^{IT} = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda_d \mathbf{I})^{-1} Y$$

- Step 2: generalizing to new targets

$$f^{\text{TS}}(\mathbf{x}, \mathbf{t}) = \mathbf{g}(\mathbf{t})^T (\mathbf{G} + \lambda_t \mathbf{I})^{-1} \mathbf{f}_T(\mathbf{x})^T$$

## Two-step kernel ridge regression

- Kernel evaluations for new test instance:

$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))^T$$

$$\mathbf{g}(\mathbf{t}) = (g(\mathbf{t}, \mathbf{t}_1), \dots, g(\mathbf{t}, \mathbf{t}_q))^T$$

- Step 1: prediction for  $\mathbf{x}$  on all the training targets

$$\mathbf{f}_T(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T A^{IT} = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda_d \mathbf{I})^{-1} Y$$

- Step 2: generalizing to new targets

$$\begin{aligned} f^{\text{TS}}(\mathbf{x}, \mathbf{t}) &= \mathbf{g}(\mathbf{t})^T (\mathbf{G} + \lambda_t \mathbf{I})^{-1} \mathbf{f}_T(\mathbf{x})^T \\ &= \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda_d \mathbf{I})^{-1} Y (\mathbf{G} + \lambda_t \mathbf{I})^{-1} \mathbf{g}(\mathbf{t}) \end{aligned}$$

## Two-step kernel ridge regression

- Kernel evaluations for new test instance:

$$\mathbf{k}(\mathbf{x}) = (k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_n))^T$$

$$\mathbf{g}(\mathbf{t}) = (g(\mathbf{t}, \mathbf{t}_1), \dots, g(\mathbf{t}, \mathbf{t}_q))^T$$

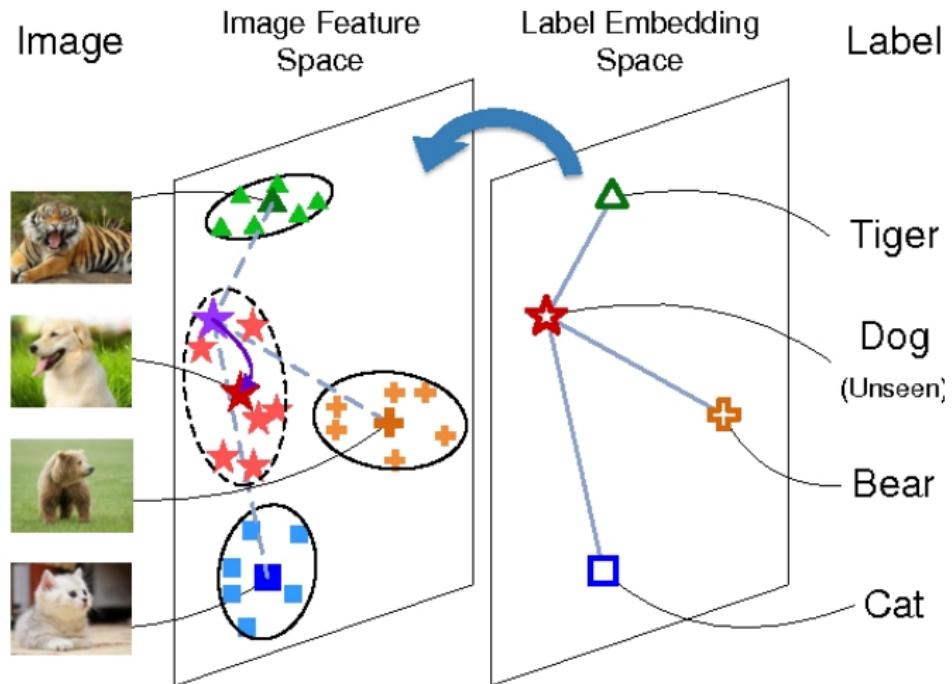
- Step 1: prediction for  $\mathbf{x}$  on all the training targets

$$\mathbf{f}_T(\mathbf{x}) = \mathbf{k}(\mathbf{x})^T A^{IT} = \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda_d \mathbf{I})^{-1} Y$$

- Step 2: generalizing to new targets

$$\begin{aligned} f^{TS}(\mathbf{x}, \mathbf{t}) &= \mathbf{g}(\mathbf{t})^T (\mathbf{G} + \lambda_t \mathbf{I})^{-1} \mathbf{f}_T(\mathbf{x})^T \\ &= \mathbf{k}(\mathbf{x})^T (\mathbf{K} + \lambda_d \mathbf{I})^{-1} Y (\mathbf{G} + \lambda_t \mathbf{I})^{-1} \mathbf{g}(\mathbf{t}) \\ &= \mathbf{k}(\mathbf{x})^T A^{TS} \mathbf{g}(\mathbf{t}) \\ &= \mathbf{w}^T (\phi(\mathbf{x}) \otimes \psi(\mathbf{t})) \end{aligned}$$

# Zero-shot learning in computer vision



# Zero-shot learning in computer vision

Pairwise model representation as before:

$$f(\mathbf{x}, \mathbf{t}) = \mathbf{w}^T (\phi(\mathbf{x}) \otimes \psi(\mathbf{t}))$$

Inference in a structured prediction fashion:

$$\hat{c}(\mathbf{x}) = \arg \max_{\mathbf{t} \in \mathcal{T}} f(\mathbf{x}, \mathbf{t})$$

Different optimization problems:

- Multi-class objective<sup>38</sup>
- Ranking objective<sup>39</sup>
- Regression objective<sup>40</sup>
- Canonical correlation analysis

Different model formulations:

- Linear embeddings
- Nonlinear embeddings

---

<sup>38</sup> Akata et al., Evaluation of Output Embeddings for Fine-Grained Image Classification, CVPR2015

<sup>39</sup> Frome et al., Devise: A deep visual-semantic embedding model, NIPS 2013

<sup>40</sup> Socher et al., g. Zero-shot learning through cross-modal transfer, NIPS 2013



## Question

In which situation(s) is it useful to exploit target relations and representations?

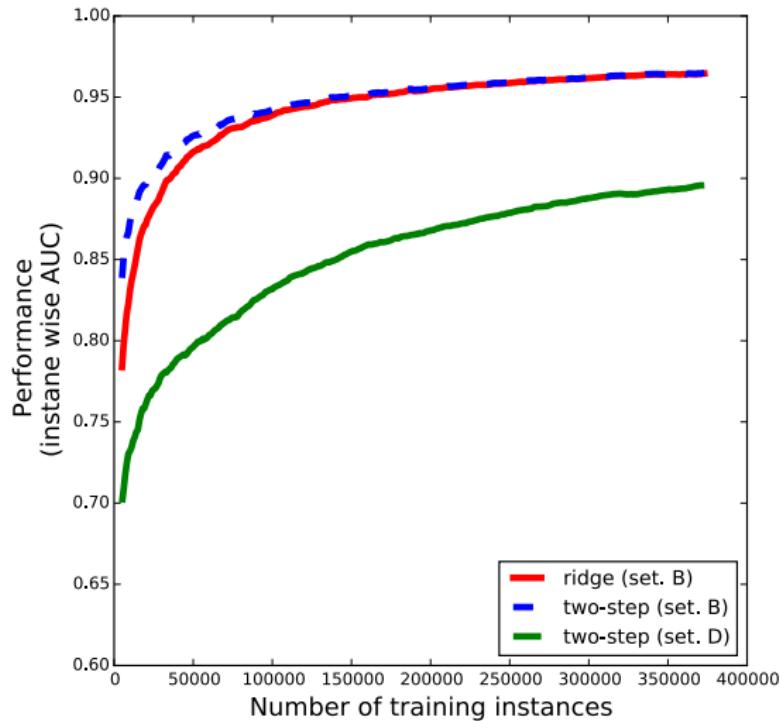
- In Setting B, when  $n$  is sufficiently large
- In Setting B, when  $n$  is sufficiently small
- In Setting D, when  $n$  is sufficiently large
- In Setting D, when  $n$  is sufficiently small

# A case study on the Wikipedia dataset

		Tags					
		Sports			Celebrities		Countries
		Tennis	Football	Biking	Movies	Tv	Belgium
01101	Text1	0	0	0	0	0	1
00111	Text2	0	0	1	0	1	1
01110	Text3	0	0	0	1	1	0
10001	Text4	0	0	1	0	1	0
01011	Text5	1	0	0	1	0	0

11110	Text6	?	?	?	?	?	?
-------	-------	---	---	---	---	---	---

# The answer



12,000 labels: from 5,000 to 350,000 instances<sup>41</sup>

<sup>41</sup> M. Stock, Exact and efficient algorithms for pairwise learning, PhD thesis, 2017

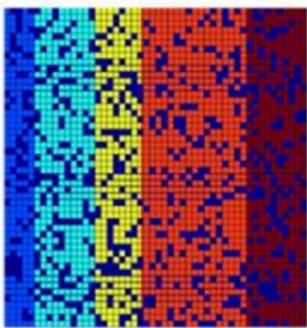
# A unifying view on MTP methods



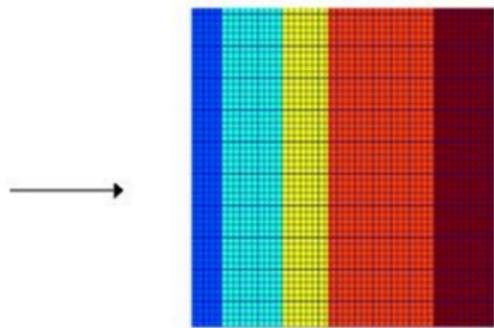
Group of methods	Applicable setting
Independent models	B
Similarity-enforcing methods	B
Relation-exploiting methods	B, C and D
Relation-constructing methods	B and C
Representation-exploiting methods	B, C and D
<b>Representation-constructing methods</b>	A, B and C

## Low-rank approximation in Settings B and C

High rank matrix



Low rank matrix



Typically perform a low-rank approximation of the parameter matrix<sup>42</sup>:

$$\min_A ||Y - XA||_F^2 + \lambda \text{rank}(A)$$

<sup>42</sup> Chen et al., A convex formulation for learning shared structures from multiple tasks, ICML 2009.

## Low-rank approximation in Settings B and C

- $A$ : parameter matrix of dimensionality  $p \times m$
- $p$ : the number of features
- $m$ : the number of targets
- Assume a low-rank structure of  $A$ :

$$U \quad \times \quad V \quad = \quad A$$

The diagram shows three matrices. On the left is a vertical matrix  $U$  with 3 rows and 3 columns. In the center is a horizontal matrix  $V$  with 3 rows and 5 columns. To the right is a large square matrix  $A$  with 5 rows and 5 columns. Between the first two matrices is a multiplication symbol ( $\times$ ). To the right of the second matrix is an approximation symbol ( $\approx$ ). This visualizes the equation  $U \times V \approx A$ .

- We can write  $A = VU$  and  $A\mathbf{x} = VU\mathbf{x}$
- $V$  is a  $p \times \hat{m}$  matrix
- $U$  is an  $\hat{m} \times m$  matrix
- $\hat{m}$  is the rank of  $A$

# Low-rank approximation in Settings B and C

## Overview of methods

- Popular for multi-output regression, multi-task learning and multi-label classification
- Linear as well as nonlinear methods
- Algorithms:
  - ▶ Principal component analysis<sup>43</sup>, Canonical correlation analysis<sup>44</sup>, Partial least squares
  - ▶ Singular value decomposition<sup>45</sup>, Alternating structure optimization<sup>46</sup>
  - ▶ Compressed sensing<sup>47</sup>, Output codes<sup>48</sup>, Landmark labels<sup>49</sup>, Bloom filters<sup>50</sup>, Auto-encoders<sup>51</sup>

---

<sup>43</sup> Weston et al., Kernel dependency estimation, NIPS 2002

<sup>44</sup> Multi-label prediction via sparse infinite CCA, NIPS 2009

<sup>45</sup> Tai and Lin, Multilabel classification with principal label space transformation, Neural Computation 2012

<sup>46</sup> Zhou et al., Clustered Multi-Task Learning Via Alternating Structure Optimization, NIPS 2011

<sup>47</sup> Hsu et al., Multi-label prediction via compressed sensing. NIPS 2009

<sup>48</sup> Zhang and Schneider, Multi-label Output Codes using Canonical Correlation Analysis, UAI 2011

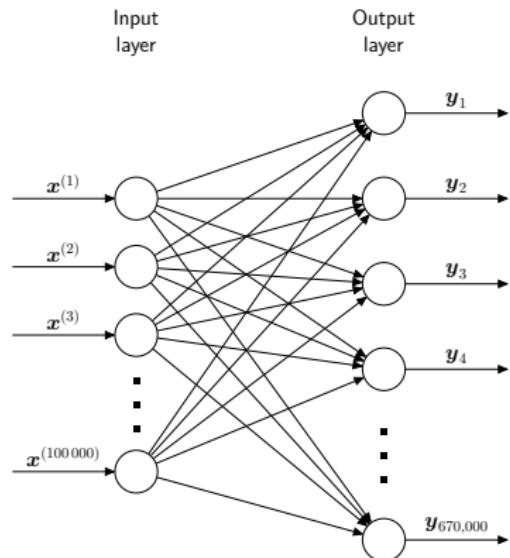
<sup>49</sup> Balasubramanian and Lebanon, The landmark selection method for multiple output prediction, ICML 2012

<sup>50</sup> Cissé et al., Robust bloom filters for large multilabel classification tasks, NIPS 2013

<sup>51</sup> Wicker et al., A nonlinear label compression and transformation method for multi-label classification using autoencoders, PAKDD 2016

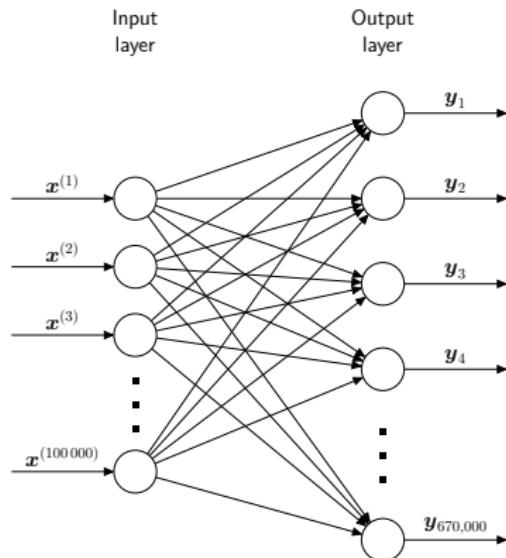
# Target embeddings in neural networks

- Shallow Networks - SVM
  - ▶ Direct mapping of input to output

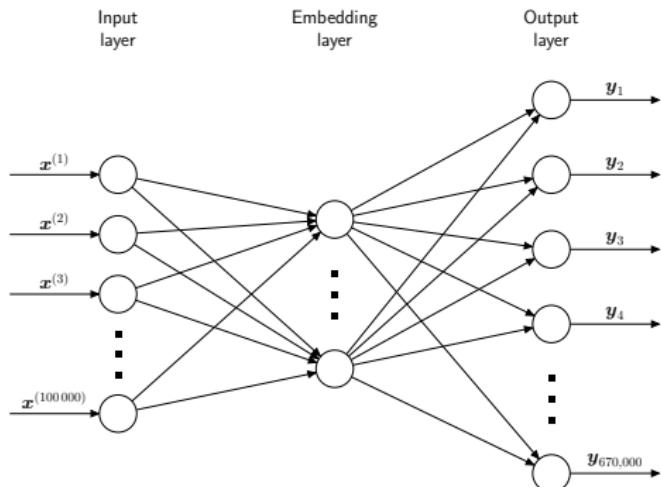


# Target embeddings in neural networks

- Shallow Networks - SVM
  - ▶ Direct mapping of input to output



- Label embedding
  - ▶ Mapping input to output via embedding layer



## Low-rank approximation in Setting A

Factorize the matrix  $Y$  instead of the parameter matrix  $A$ :

		Items									
		1	3	5		5	4				
		5		4				2	1	3	
		2	4	1	2	3	4	3	5		
		2	4	5		4			2		
			4	3	4	2			2	5	
		1	3	3		2			4		

Users

$\sim$

$$\begin{matrix} & \bullet & \end{matrix} \begin{matrix} \text{Users} \\ \begin{array}{|c|c|c|} \hline .1 & -.4 & 2 \\ \hline -.5 & .6 & 5 \\ \hline -.2 & .3 & .5 \\ \hline 1.1 & 2.1 & .3 \\ \hline -.7 & 2.1 & -2 \\ \hline -1 & .7 & .3 \\ \hline \end{array} \end{matrix} \bullet \begin{matrix} \text{Items} \\ \begin{array}{|c|c|c|c|c|c|c|c|c|c|c|c|} \hline 1.1 & -.2 & .3 & .5 & \textcolor{red}{-2} & -.5 & .8 & -.4 & .3 & 1.4 & 2.4 & -.9 \\ \hline -.8 & .7 & .5 & 1.4 & \textcolor{red}{-3} & -1 & 1.4 & 2.9 & -.7 & 1.2 & -1 & 1.3 \\ \hline 2.1 & -.4 & .6 & 1.7 & \textcolor{red}{2.4} & .9 & -.3 & .4 & .8 & .7 & -.6 & .1 \\ \hline \end{array} \end{matrix}$$

$$Y = U \times V$$

# Low-rank approximation in Setting A

## Overview of algorithms

- Nuclear norm minimization<sup>52</sup>
- Gaussian processes<sup>53</sup>
- Probabilistic methods<sup>54</sup>
- Spectral regularization<sup>55</sup>
- Non-negative matrix factorization<sup>56</sup>
- Alternating least-squares minimization<sup>57</sup>

---

<sup>52</sup> Candes and Recht, Exact low-rank matrix completion via convex optimization. Foundations of Computational Mathematics 2008

<sup>53</sup> Lawrence and Urtasun, Non-linear matrix factorization with Gaussian processes, ICML 2009

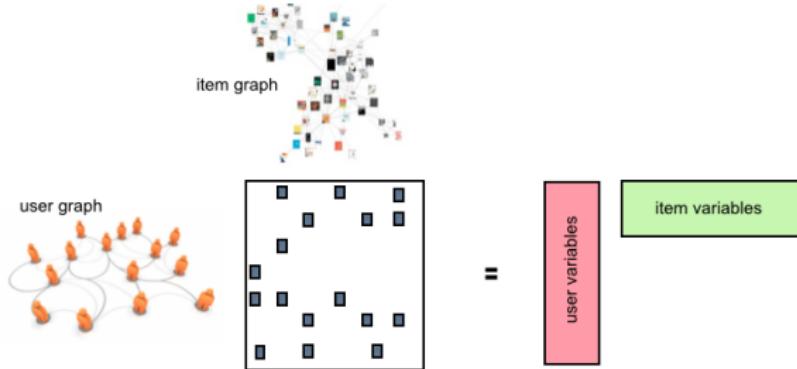
<sup>54</sup> Shan and Banerjee, Generalized probabilistic matrix factorizations for collaborative filtering, ICDM 2010

<sup>55</sup> Mazumder et al., Spectral regularization algorithms for learning large incomplete matrices., JMLR 2010

<sup>56</sup> Gaujoux and Seoighe, A flexible R package for nonnegative matrix factorization. BMC bioinformatics 2010

<sup>57</sup> Jain et al., Low-rank matrix completion using alternating minimization, ACM Symposium on Theory of Computing 2013

# Matrix factorization with side information for Setting A



- Construct **implicit** features  $(\mathbf{x}^I, \mathbf{t}^I)$  for users and items with matrix factorization methods
- Exploit **explicit** features  $(\mathbf{x}^E, \mathbf{t}^E)$  (a.k.a. side information)
- Concatenate:

$$\mathbf{x}^C = (\mathbf{x}^I, \mathbf{x}^E), \quad \mathbf{t}^C = (\mathbf{t}^I, \mathbf{t}^E)$$

- Apply methods that we have seen before<sup>5859</sup>:

$$f(\mathbf{x}^C, \mathbf{t}^C) = \mathbf{w}^T (\phi(\mathbf{x}^C) \otimes \psi(\mathbf{t}^C))$$

<sup>58</sup> Menon and Elkan, A log-linear model with latent features for dyadic prediction, ICDM 2010

<sup>59</sup> Volkovs and Zemel, Collaborative filtering with 17 parameters, NIPS 2012

## Hybrid matrix factorization for Setting A

Basilico and Hofmann, 2004; Abernethy et al, 2008; Adams et al, 2010;  
Fang and Si, 2011; Zhou et al, 2011a; Menon and Elkan, 2011; Zhou et al,  
2012b).

# When is it useful to construct target representations?

Does not work well in extreme multi-label classification<sup>60</sup>:

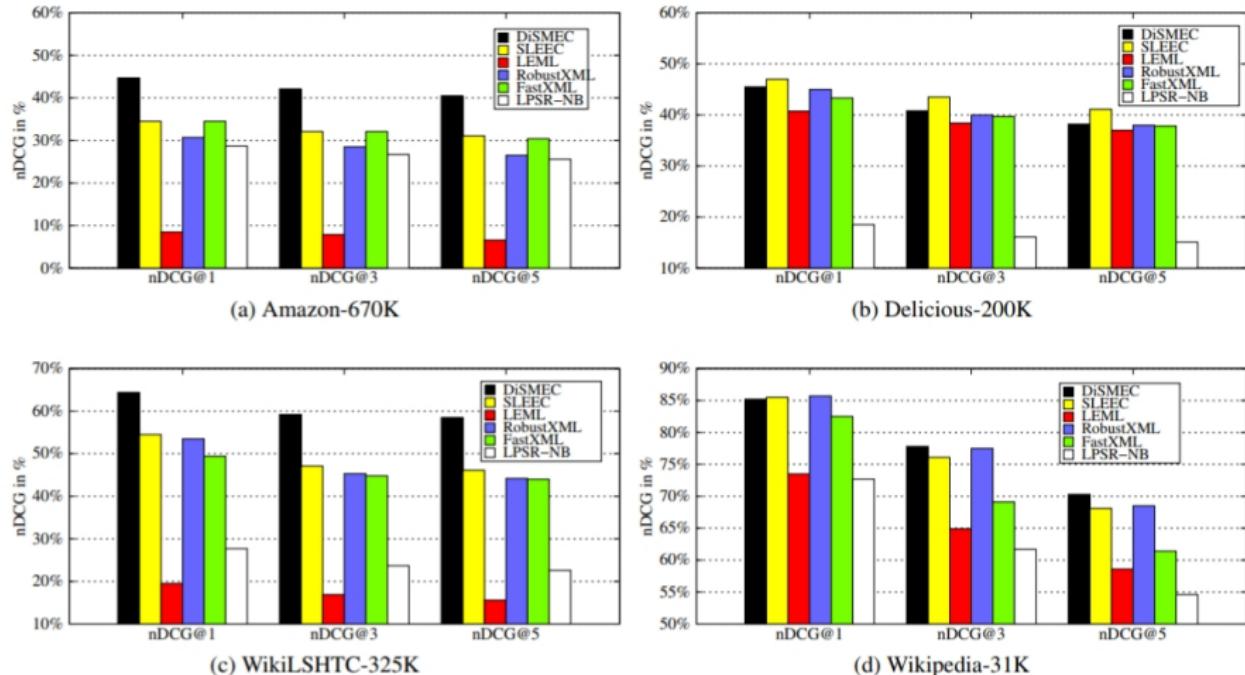


Figure 3: nDCG@k for k=1, 3 and 5

<sup>60</sup> Babbar and Schölkopf, DISMEC: Distributed Sparse Machines for Extreme Multi-label classification, WSDM 2017

# When is it useful to construct target representations?

## SVD interpretation

$$\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}^T$$

Diagram illustrating the Singular Value Decomposition (SVD) of matrix  $A$ :

- Representation of instance (or feature)**: Matrix  $\mathbf{U}$  is shown as a gray rectangle. A vertical line labeled  $j$  indicates a specific column, representing the representation of an instance or feature.
- Representation of target**: Matrix  $\mathbf{V}^T$  is shown as a gray rectangle. A vertical line labeled  $i$  indicates a specific row, representing the representation of a target.
- Matrix  $A$** : Matrix  $\mathbf{A}$  is shown as a gray rectangle.
- Sigma matrix  $\Sigma$** : Matrix  $\Sigma$  is shown as a gray rectangle with diagonal lines representing non-zero singular values  $\sigma_1, \sigma_2, \dots, \sigma_I$ .

Dimensions below the matrices:

- $\mathbf{A}$ :  $p \times m$
- $\mathbf{U}$ :  $p \times p$
- $\Sigma$ :  $p \times m$
- $\mathbf{V}^T$ :  $m \times m$

- $\sigma_1, \sigma_2, \dots$ : singular values of  $A$
- Rank of  $A$  = number of non-zero **singular values**
- **High rank** when a lot of singular values differ from zero
- **Low rank** when a lot of singular values are zero
- Singular values **give insight** in what can be gained

# Conclusions

- Multi-target prediction is an active field of research that connects different types of machine learning problems
- In the corresponding subfields of machine learning, problems have typically been solved in isolation, without establishing connections between methods
- Two-step zero-shot learning is a simple MTP method with a lot of interesting properties

**Upcoming paper:**

**Waegeman et al.**

**Multi-Target Prediction:**

**A Unifying View on Problems and Methods**

## Multi-target prediction papers at ECML/PKDD 2018

- **Djerrab et al.**, Output Fisher embedding regression, Tuesday, 15h20
- **Pikalos et al.**, Global multi-output decision trees for interaction prediction, Thursday 11h20
- **Masera and Blanzieri**, AWX: An integrated approach to hierarchical multi-label classification, Tuesday 14h40
- **Decubber et al.**, Deep F-measure maximization in multi-label classification: a comparative study, Tuesday 14h20
- **Park and Read**, A blended metric for multi-label optimization and evaluation, Thursday 11h40
- **Rafailidis and Crestani**, Deep collaborative filtering with multifaceted contextual information in location-based social networks, Thursday 14h20
- **Du et al.**, POLAR: Attention-based CNN for one-shot personalized article recommendation, Thursday 14h40
- **Lan et al.**, Personalized thread recommendation of MOOC discussion forums, Thursday 14h20
- **Marecek et al.**, Matrix completion under interval uncertainty, Thursday 16h30