

Actively querying superset labels using indecision: the k-nn case

**Sebastien Destercke, Vu-Linh Nguyen, Mylene
Masson**

Université de Technologie de Compiègne

Learning with superset labels

Features			Labels
[0.1, 1.5]	...	0.6	a
0.3	...	0.2	$\{a, b\}$
0.3	...	[0.2, 0.5]	$\{a, b, c\}$

- ▶ Partial data can induce uncertainty in learning process.
- ▶ May happen in a number of situations :
 - Expert labelling,
 - Using easily accessible information to get label set (e.g., actor list to do facial recognition of TV series pictures),
 - Data collection with sensors of various qualities,
 - Data descriptions using different levels of details (coarsening)

Learning from partial data

Two view points

- ▶ adapting classical approaches to learn one optimal model. For instance, defining specific loss functions [T. Cour *et al*, 2011.]
 - Necessary assumptions on the missingness process
- ▶ learning (IP) models to get set of optimal models from all completions of partial data. For example, the paper of [E. Hullermeier, 2014].
 - No or few assumptions on the missingness process

This work

- ◇ We adopt the second view regarding data completions
- ◇ We wonder about which data to query to make better predictions

Some formalisation

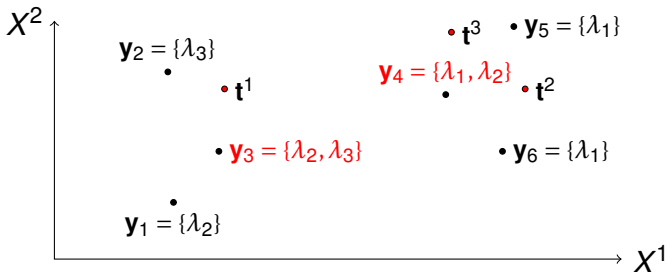
- A training set $D = \{x_i, \mathbf{y}_i\}$ with
 - $x_i \in R^d$ are precise values
 - $\mathbf{y}_i \subseteq \mathcal{Y} = \{\lambda_1, \dots, \lambda_M\}$ are partial, a.k.a. superset labels
- An evaluation set $T = \{t_i\}$ of input instances, $t_i \in R^d$
- Possibly a decision function $h: D \rightarrow \mathcal{Y}$ providing a precise prediction

Which labels \mathbf{y}_i should we query to improve our model accuracy/decisiveness ?

K-nn classifier for partially labelled data

A simple (maximax) way to take decision despite partial labels [E. Hullermeier & J. Beringer, 2006]

$$h(\mathbf{t}) = \arg \max_{\lambda \in \Omega} \sum_{\mathbf{y}_k \in \mathbf{N}_{\mathbf{t}}} w_k \mathbb{1}_{\lambda \in \mathbf{y}_k}. \quad (1)$$

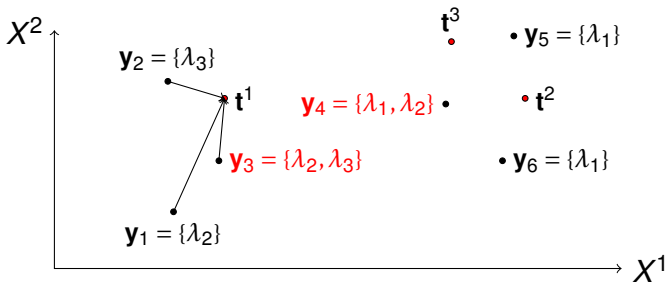


◇ Which partial label is the more informative ?

First idea : the more decision it is involved in, the more informative it is

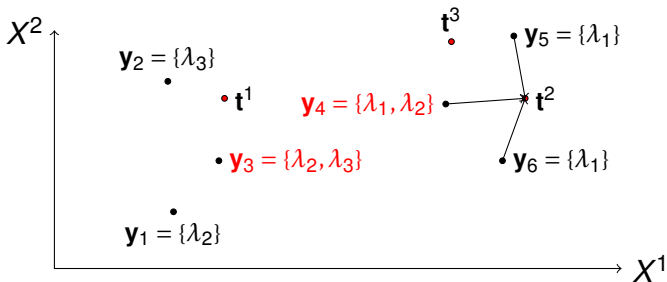
1. define a given number K of neighbours
2. for each x_i with partial label, count the number of items in T of which it is a neighbour
3. query the item y_j involved in most decisions

First idea : the more decision it is involved in, the more informative it is



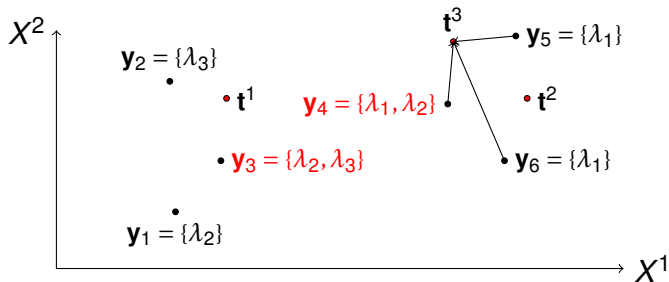
◇ y_3 involves in decision of $\{t^1\}$.

First idea : the more decision it is involved in, the more informative it is



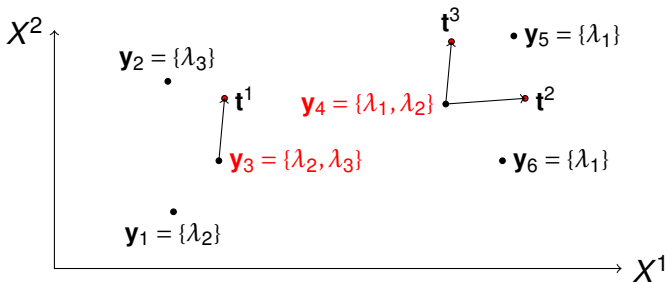
- ◇ y_3 involves in decision of $\{t^1\}$.
- ◇ y_4 involves in decisions of $\{t^2\}$.

First idea : the more decision it is involved in, the more informative it is



- ◇ y_3 involves in decision of $\{t^1\}$.
- ◇ y_4 involves in decisions of $\{t^2, t^3\}$.

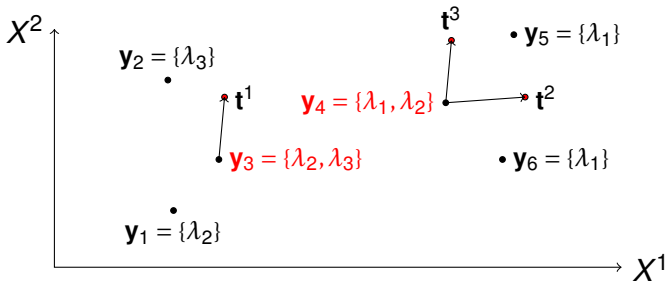
First idea : the more decision it is involved in, the more informative it is



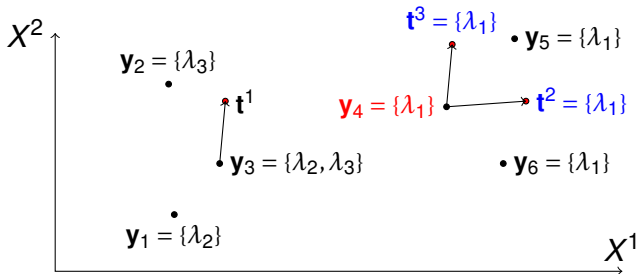
- ◇ y_3 involves in decision of $\{t^1\}$.
- ◇ y_4 involves in decisions of $\{t^2, t^3\}$.

▷ This strategy chooses y_4 as the optimal query.

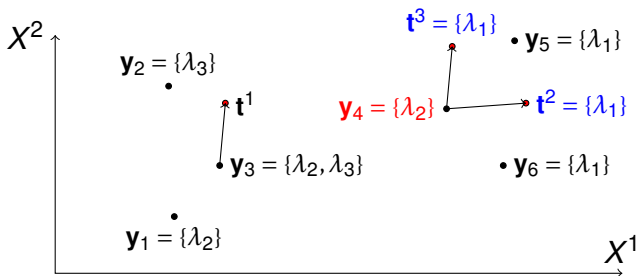
First idea : The more decision it is involved, the more informative it is



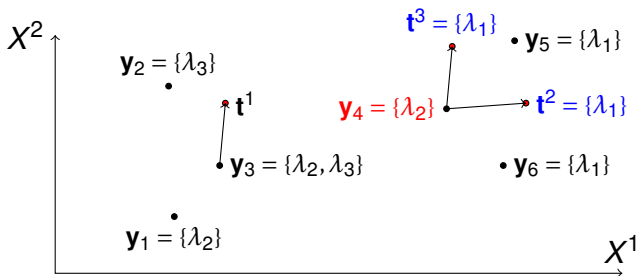
- ◇ It is simple, but is it a good idea ?
 - ▷ What do we gain if we query y_4 ?



- If $y_4 = \lambda_1$, then decisions are λ_1 .



- If $y_4 = \lambda_1$, then decisions are λ_1 .
- If $y_4 = \lambda_2$, then decisions are λ_1 .



- If $y_4 = \lambda_1$, then decisions are λ_1 .
 - If $y_4 = \lambda_2$, then decisions are λ_1 .
- ▷ Querying y_4 does not change predictions.

Second idea : The more uncertainty the label introduces, the more informative it is

- Assume $\mathbf{y}_1, \dots, \mathbf{y}_K$ are neighbours
- Set $\mathbf{L}_t = \{(l_1^t, \dots, l_K^t) \mid l_k^t \in \mathbf{y}_k^t\}$ is the selection of partial labels
- Set of possible predicted labels

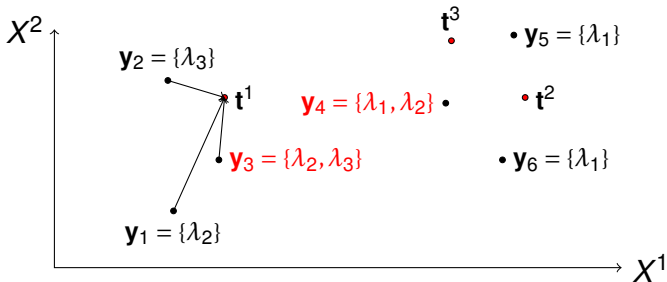
$$\mathbf{PL}_t = \{\lambda \in \Omega \mid \exists \mathbf{l}^t \in \mathbf{L}_t \text{ s.t. } \lambda \in \hat{\lambda}_{\mathbf{l}^t}\}$$

- Set of necessary predicted labels

$$\mathbf{NL}_t = \{\lambda \in \Omega \mid \forall \mathbf{l}^t \in \mathbf{L}_t, \lambda \in \hat{\lambda}_{\mathbf{l}^t}\}$$

An instance t is said to be ambiguous if $\mathbf{PL}_t \neq \mathbf{NL}_t$

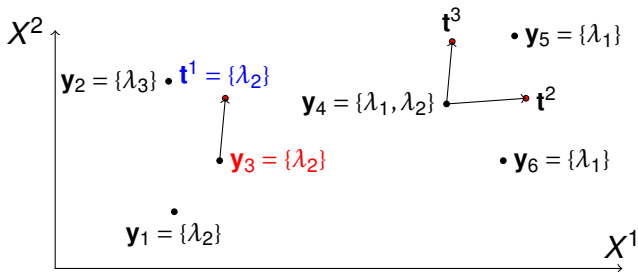
Second idea : The more uncertainty the label introduces, the more informative it is



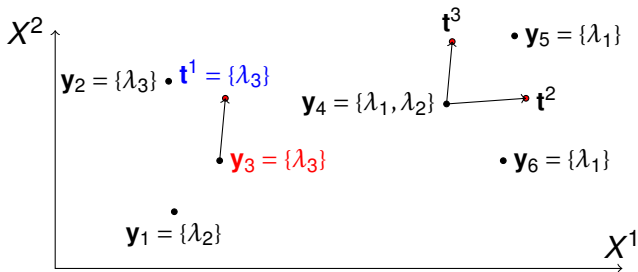
◇ Assuming a weighted 3-nn, $\mathbf{PL}_{t^1} = \{\lambda_2, \lambda_3\}$, $\mathbf{NL}_{t^1} = \{\}$

Second idea : The more uncertainty the label introduces, the more informative it is

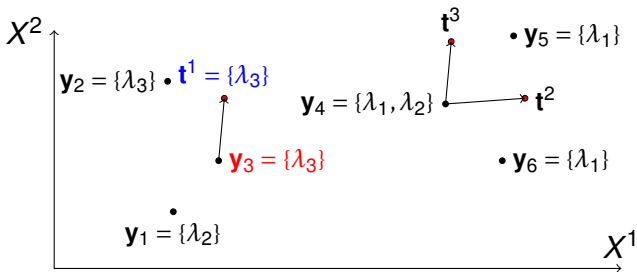
- A data x_i can reduce ambiguity on t if after knowing \mathbf{y}_i , \mathbf{NL}_t or \mathbf{PL}_t can potentially change
- For a given K and x_i , count the number of items in \mathbf{T} for which it can reduce ambiguity
- Query the \mathbf{y}_i with greatest potential ambiguity reductions



- If $y_3 = \lambda_2$, then decision is λ_2 .



- If $y_3 = \lambda_2$, then decision is λ_2 .
- If $y_3 = \lambda_3$, then decision is λ_3 .



- If $y_3 = \lambda_2$, then decision is λ_2 .
- If $y_3 = \lambda_3$, then decision is λ_3 .
- ◇ Querying y_3 can reduce ambiguity, while y_4 does not.
 - ▷ This strategy choose y_3 as the optimal query.

Computational considerations

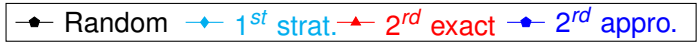
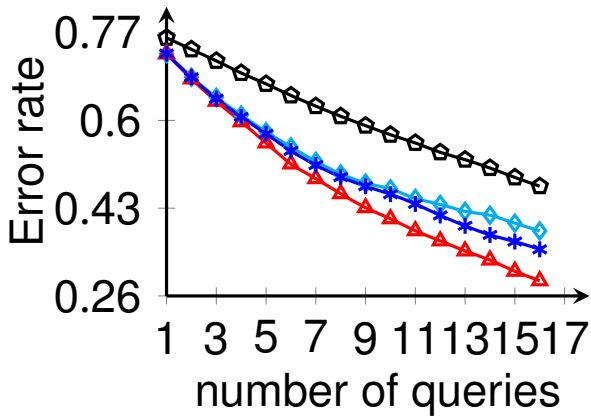
Link with computational social choice

- Computing **PL**, **NL** : equivalent to possible/necessary winner in plurality voting with partial voter preferences
- Querying a label y_j : eliciting precise preferences of a voter

Computational issues

- Computing **NL** : can be done in linear time
- Computing **PL** in unweighted case : can be done in cubic time (reducible to maximum flow problem)
- Computing **PL** in weighted case : seems NP-hard (reducible to 3-dimensional matching), but dynamical programming maybe possible, otherwise approximate

Experimental results illustrated on one data set



Other related questions

- Each query provides information about the "imprecisiation" process → can we use it to improve our results
- In practice, optimal model may be identifiable from partial information → under which conditions ?
- Can the idea be efficiently extended to other simple settings (learning NCC, Classification trees) or more complex ones (learning preferences, dynamical models)