

Goal:

Learn a classifier from positive and unlabeled relational data (relational PU learning)

Shortcoming of current solution [1]:

Cannot learn disjunctive concepts

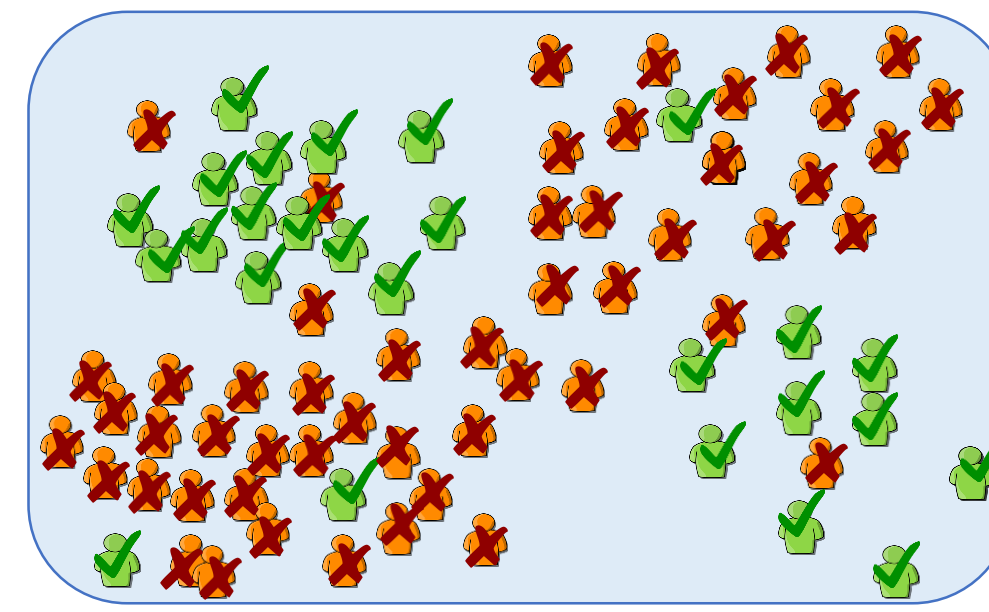
Solution:

Use the label frequency to adjust standard learning algorithms to the PU scenario. This is a common approach in propositional PU learning

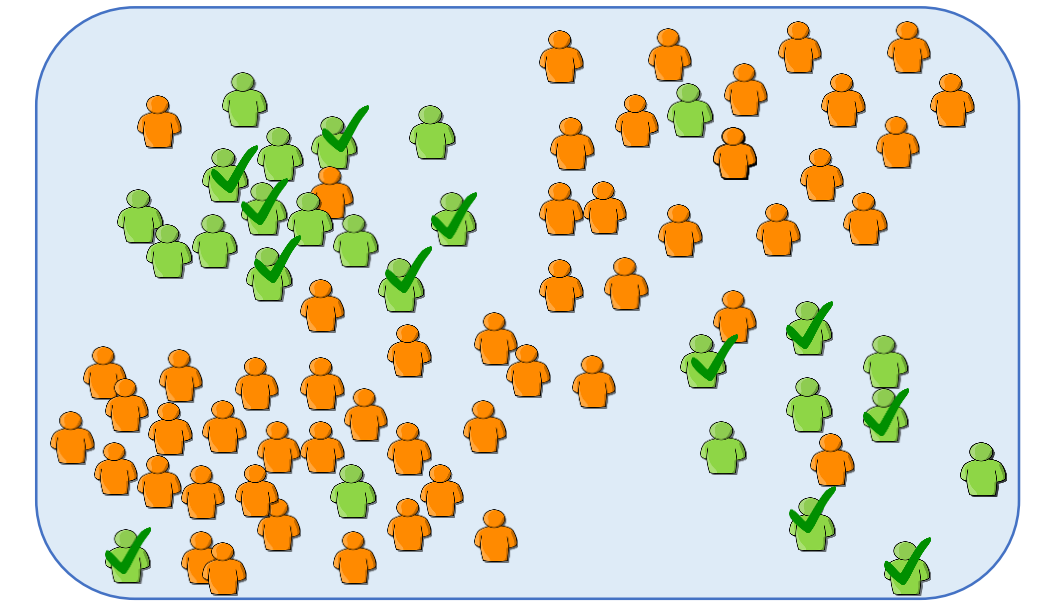
Challenge:

Estimate the label frequency in relational PU data

Positive and Unlabeled (PU) Learning



Supervised Data



Positive and Unlabeled Data

Common assumption:

Positive examples get labeled with constant probability c , the *label frequency*

$$c = P(\text{labeled} \mid \text{positive}, \text{facts}) \\ = P(\text{labeled} \mid \text{positive})$$

Estimate Label Frequency c from PU data

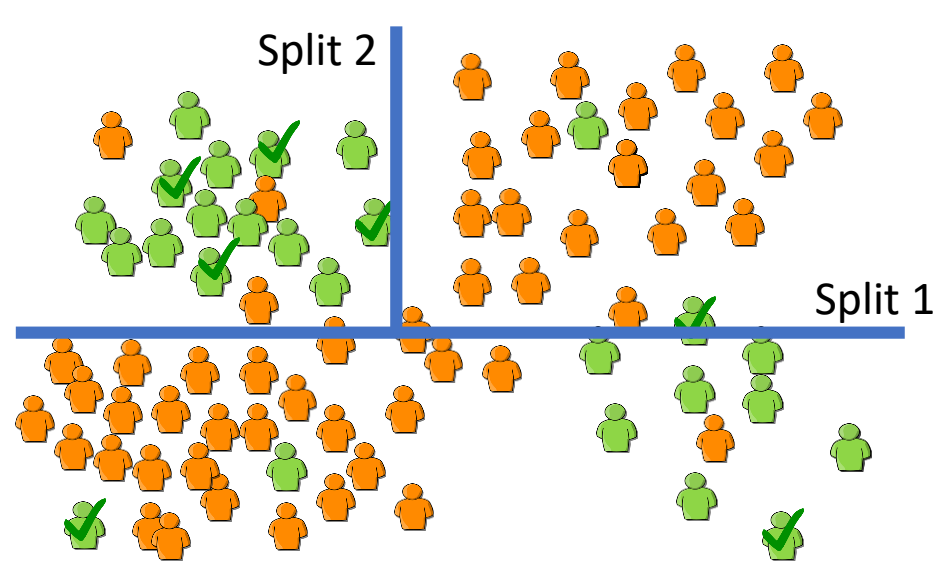
Insight 1: Data subset implies lower bound on c

$$P \leq T \Rightarrow c \geq \frac{L}{T} - \underbrace{\varepsilon(T)}_{\text{Error term from 1-sided Chebyshev inequality}}$$

Insight 2: Positive subsets give very tight bounds

Insight 3: Highly labeled subsets are likely positive

➔ Look for those through decision tree induction (Tilde)
Use subsets to tighten lower bound



$$\begin{aligned} \text{Init: } c &\geq \frac{7}{78} - \varepsilon(78) \\ &= 0.09 - \varepsilon(78) \\ \text{Split 1: } c &\geq \frac{5}{39} - \varepsilon(39) \\ &= 0.13 - \varepsilon(39) \\ \text{Split 2: } c &\geq \frac{4}{17} - \varepsilon(17) \\ &= 0.24 - \varepsilon(17) \end{aligned}$$

Simple PU Learning using the Label Frequency c [2]

Method 1: Probabilistic classifier that learns $P(\text{labeled} \mid \text{facts})$

Adjust output probabilities using the formula

E.g. Tilde: Probabilistic Relational Decision Trees

Method 2: Adjust learning algorithm using c :

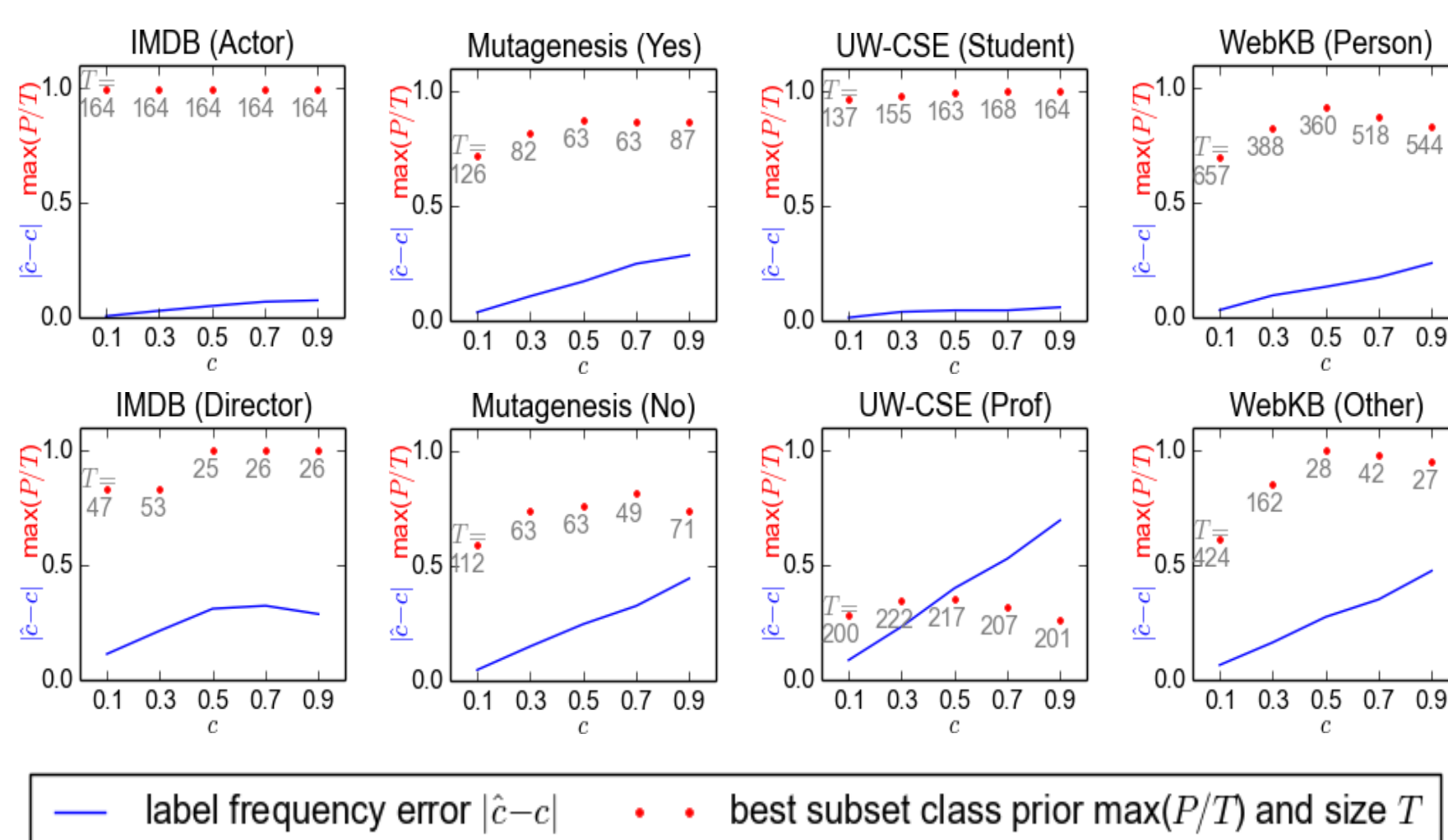
$$P = L/c \text{ and } N = T - P$$

E.g. Aleph: adjust score function

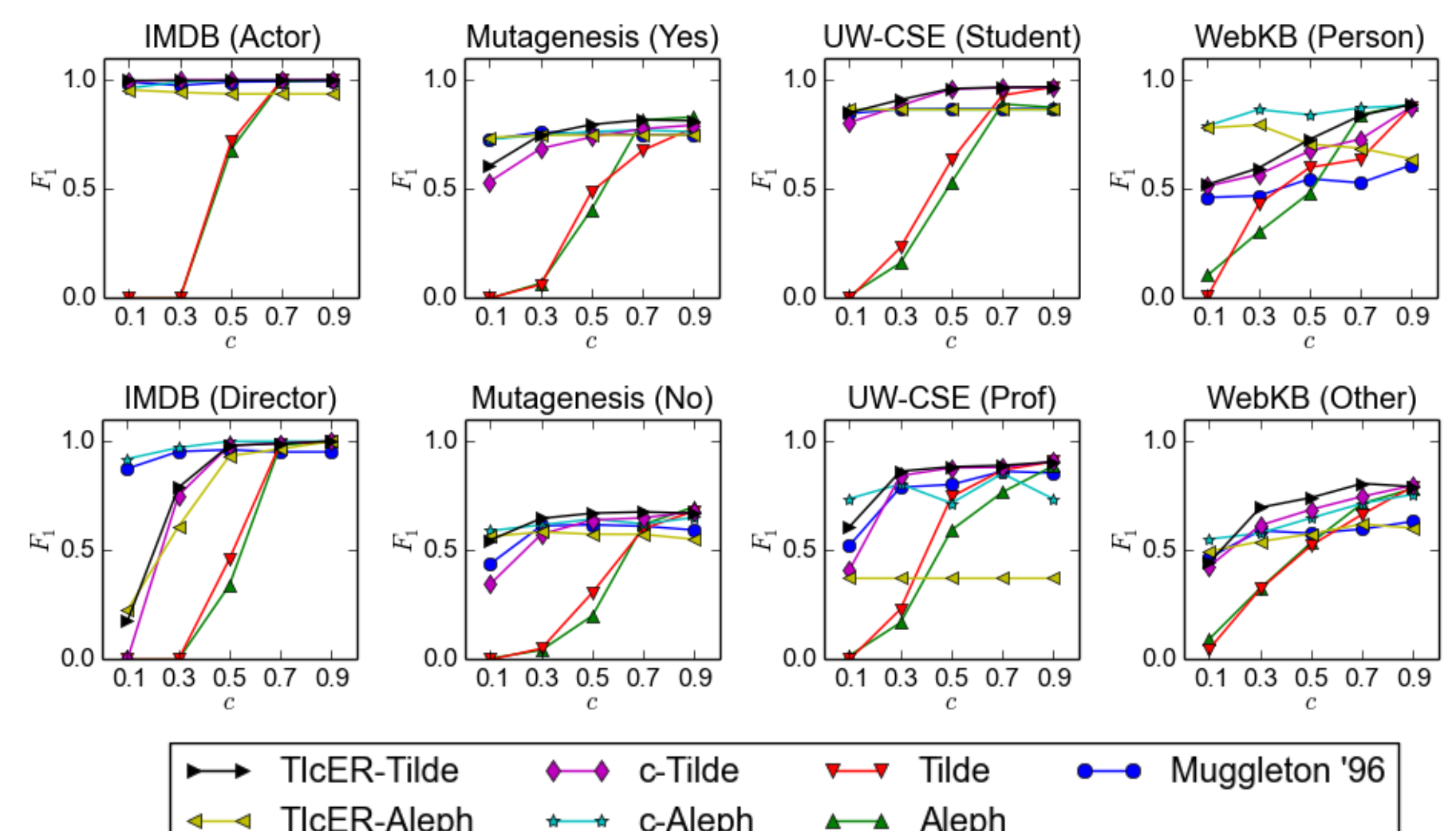
Supervised: Coverage = $P - N$

PU: Coverage = $L/c - (T - L/c) = 2L/c - T$

Label Frequency Estimation Results



Method Comparison



References

- [1] Muggleton, Stephen. Learning from positive data. ILP, 1996.
- [2] Elkan, Charles, and Noto, Keith. Learning classifiers from only positive and unlabeled data. KDD, 2008.