



TECHNISCHE
UNIVERSITÄT
DARMSTADT

On the Understandability of Rule Learning

Johannes Fürnkranz

TU Darmstadt

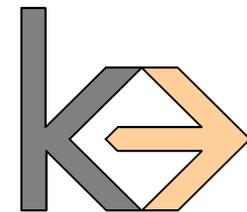
Knowledge Engineering Group

Hochschulstrasse 10

D-64289 Darmstadt

06151/166238

`juffi@ke.tu-darmstadt.de`



Joint Work with H. Paulheim, T. Kliegr, F. Janssen and J. Stecher

Data Mining

statistics



background
knowledge



utility theory
cost models



Data Mining is the non-trivial
process of identifying

▪ valid

▪ novel

▪ potentially useful

▪ ultimately understandable

patterns in data.

(Fayyad et al. 1996)



Understandability → Rules?

- Data:
 - Fertility and Family Survey 1995/96 for Italians and Austrians
 - Features based on general descriptors and variables that describes whether (quantum), at which age (timing) and in what order (sequencing) typical life course events have occurred.
- Objective:
 - Find rules that capture typical life courses for either country
- Examples:

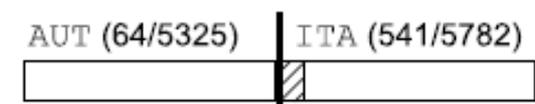
```
IF LeftHome < Marriage
THEN AUT
```



```
IF Union = Marriage
AND Education <= 14
THEN ITA
```

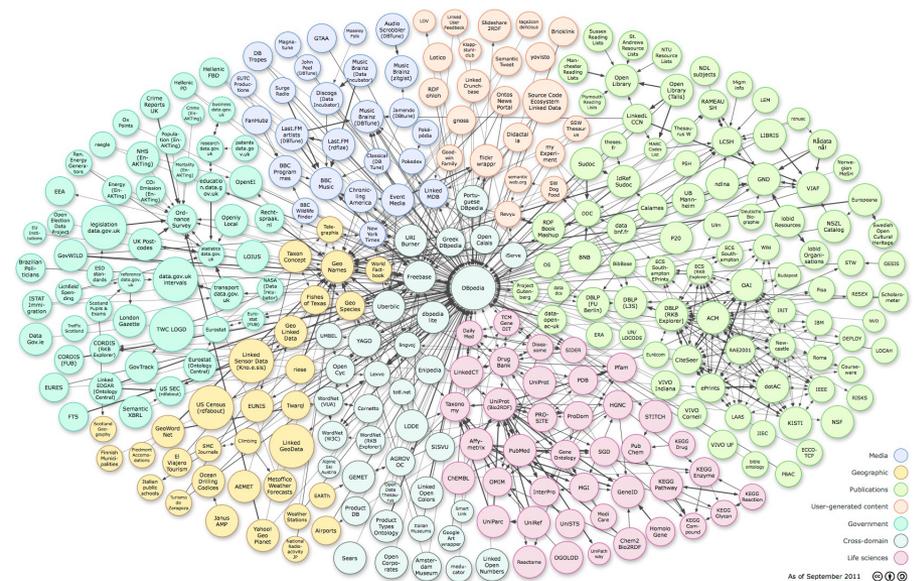


```
IF Union = Marriage
AND Education >= 22
THEN ITA
```



Why Rules?

- Rules provide a good (the best?) trade-off between
 - human understandability
 - machine executability
- Used in many applications which will gain importance in the near future
 - Security
 - Spam Mail Filters
 - Semantic Web
- But they are not a universal tool
 - e.g., learned rules sometimes lack in predictive accuracy
→ challenge to close or narrow this gap



Understandability – State of Affairs

Data Mining essentially assume

- Rules are inherently understandable
- Shorter rules are more understandable than longer rules
- Good explanations = Good fit to the data
- No additional criteria or algorithms are needed to address understandability

The point of this talk is to question these assumptions



Overview

- Motivation
 - Understandability has not received much attention
- Understandability
 - Conjunctive Fallacy
 - Gambler's Fallacy
 - Representativeness Heuristic
- Different Types of Rules
 - Discriminative vs. Characteristic Rules
 - Formal Concepts
 - Closed Itemsets
- Heuristic Rule Learning
 - Concept Learning
 - Coverage Spaces
 - Rule Learning Heuristics
- Inverted Heuristics
- Explain-A-LOD
 - Semantic Coherence
 - Representation Heuristics
- Algorithmic Enhancements
 - Structured theories
 - More complex problems
- Conclusions



Conjunctive Fallacy

(Tversky & Kahneman 1983)

Linda is 31 years old, single, outspoken, and very bright. She majored in philosophy. As a student, she was deeply concerned with issues of discrimination and social justice, and also participated in anti-nuclear demonstrations.

Which is more probable?

- A) Linda is a bank teller.
- B) Linda is a bank teller and is active in the feminist movement.

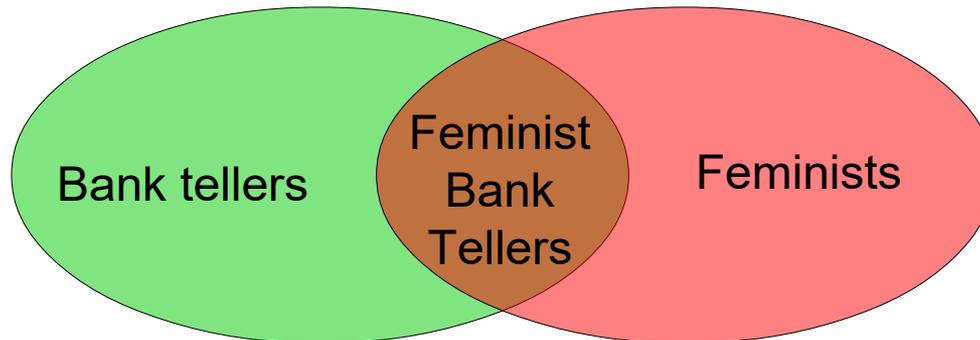


Conjunctive Fallacy

(Tversky & Kahneman 1983)

- The majority of people (85%) preferred B)
- However, B) is a specialization of A), so that A) cannot be less probable than B)

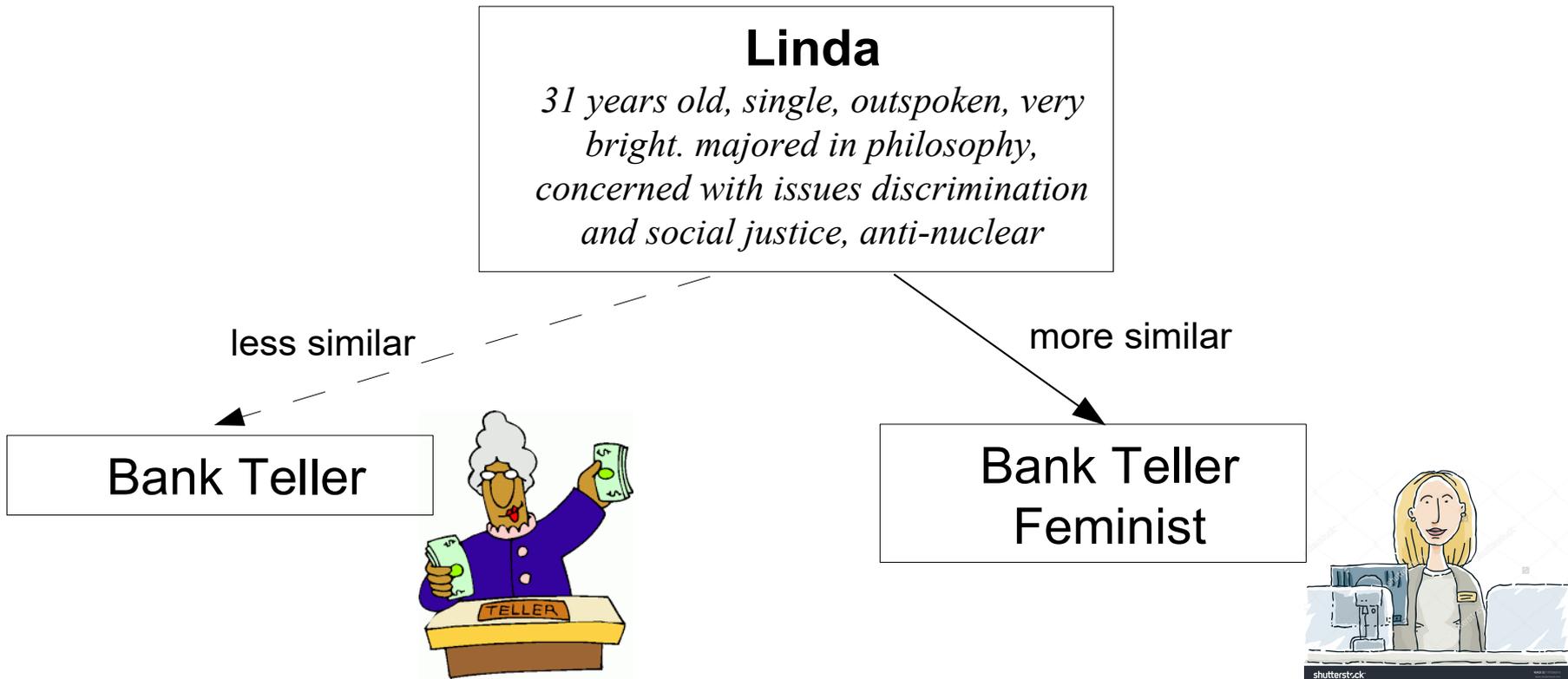
$$\Pr(\textit{bank} \wedge \textit{feminist}) = \Pr(\textit{feminist} | \textit{bank}) \cdot \Pr(\textit{bank}) \leq \Pr(\textit{bank})$$



Representativeness Heuristics

(Kahneman & Tversky 1972)

Humans tend to judge probability of a subgroup according to how similar it is to a prototype of the base group.



Gambler's Fallacy

Which sequence of outcomes on the roulette table is more likely?

- A) 0 0 0 0 0 0
- B) 27 18 4 23 8 17



People tend to think the 2nd sequence is more likely because it is *more representative of a random sequence*.



Gambler's Fallacy

Which sequence of outcomes on the roulette table is more likely?

- A) 4 8 17 18 23 27
- B) 27 18 4 23 8 17



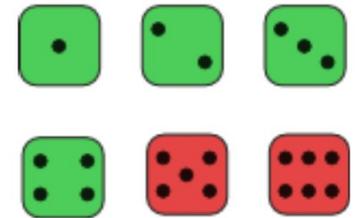
People tend to think the 2nd sequence is more likely because it is *more representative of a random sequence*.



Gambler's Fallacy

(Tversky & Kahneman 1983)

Consider a regular six-sided die with four green faces and two red faces. The die will be rolled 20 times and the sequence of greens (G) and reds (R) will be recorded. You are asked to select one sequence, from a set of three, and you will win \$25 if the sequence you choose appears on successive rolls of the die.



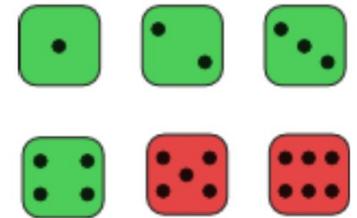
- A)
- B)
- C)

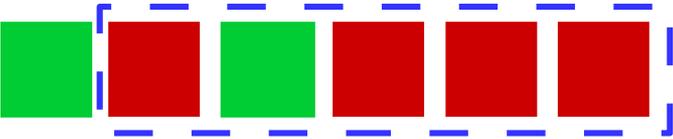


Gambler's Fallacy

(Tversky & Kahneman 1983)

Consider a regular six-sided die with four green faces and two red faces. The die will be rolled 20 times and the sequence of greens (G) and reds (R) will be recorded. You are asked to select one sequence, from a set of three, and you will win \$25 if the sequence you choose appears on successive rolls of the die.



- A) 
- B) 
- C) 

65% bet on B) even though A) is a subsequence of B) and will thus appear more frequently



Understandability vs. Rule Length

Conventional Rule learning algorithms tend to learn short rules

- They favor to add conditions that exclude many negative examples

Typical intuition: Short rules are better

- long rules are less understandable, therefore short rules are preferable
- short rules are more general, therefore (statistically) more reliable and would have been easier to falsify on the training data

Claim: Shorter rules are not always better

- **Predictive Performance:** Longer rules often cover the same number of examples than shorter rules so that (statistically) there is no preference for choosing one over the other
- **Understandability:** In many cases, longer rules may be much more intuitive than shorter rules

→ *we need to understand understandability!*



Overview

- Motivation
 - Understandability has not received much attention
- Understandability
 - Conjunctive Fallacy
 - Gambler's Fallacy
 - Representativeness Heuristic
- **Different Types of Rules**
 - Discriminative vs. Characteristic Rules
 - Formal Concepts
 - Closed Itemsets
- Heuristic Rule Learning
 - Concept Learning
 - Coverage Spaces
 - Rule Learning Heuristics
- Inverted Heuristics
- Explain-A-LOD
 - Semantic Coherence
 - Representation Heuristics
- Algorithmic Enhancements
 - Structured theories
 - More complex problems
- Conclusions



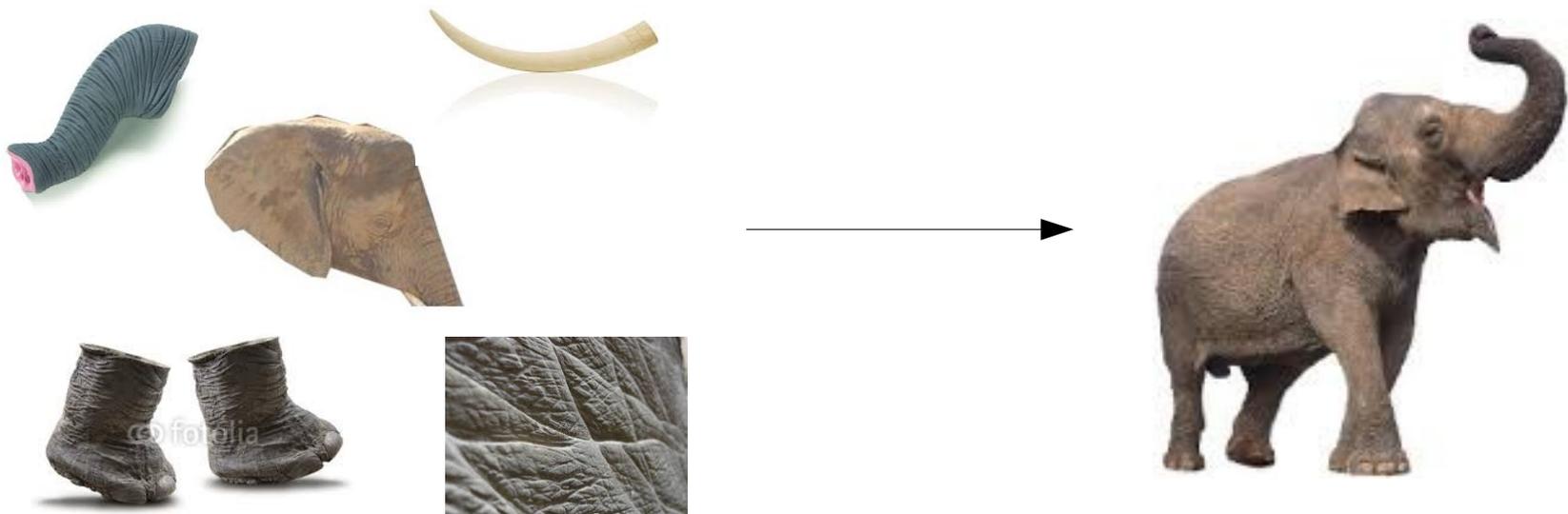
Discriminative Rules

- Allow to quickly **discriminate an object** of one category from objects of other categories
- Typically a few properties suffice
- Example:



Characteristic Rules

- Allow to characterize an object of a category
- Focus is on all properties that are **representative** for objects of that category
- Example:



Discriminative Rules vs. Characteristic Rules

(Michalski 1983)

Michalski (1983) discerns two kinds of classification rules:

- **Discriminative Rules:**

- A way to distinguish the given class from other classes

Features → **Class**

- Most interesting are *minimal discriminative rules*.

- **Characteristic Rules:**

- A conjunction of all properties that are common to all objects in the class

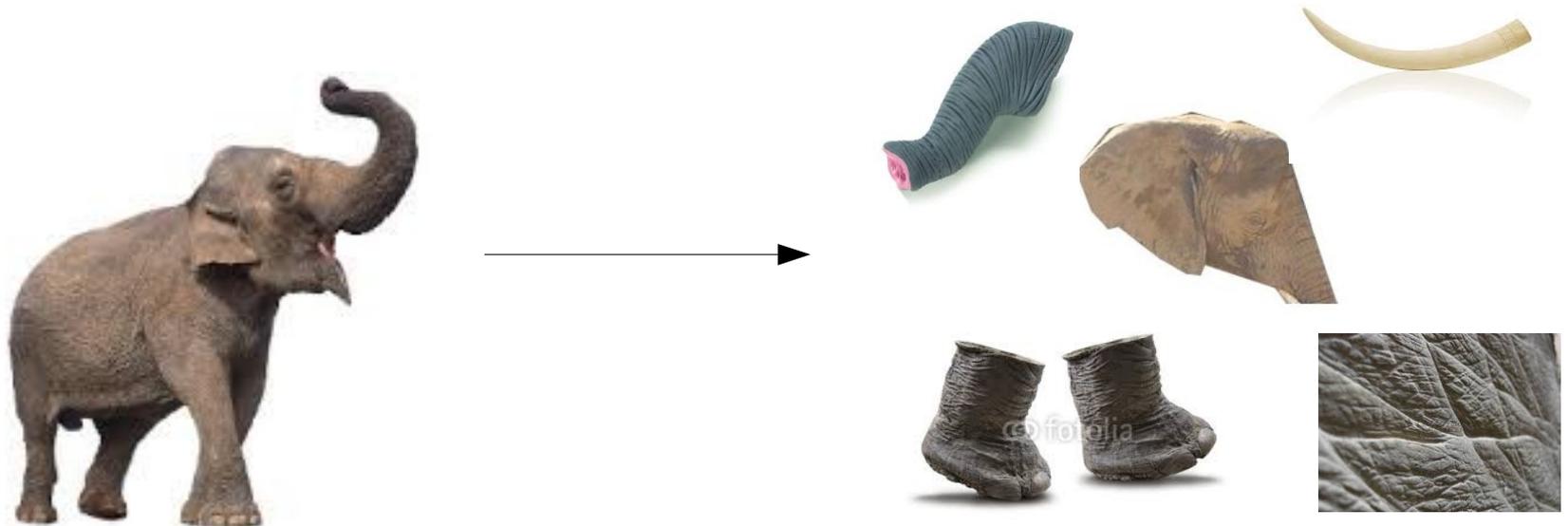
Class → **Features**

- Most interesting are *maximal characteristic rules*.



Characteristic Rules

- An alternative view of characteristic rules is to invert the implication sign
- All properties that are implied by the category
- Example:



(Informal) Formal Concept Analysis

(Wille 1982)

Intent of a Concept (Rule)

- Conjunction of Features

Extent of a Concept (Coverage)

- All objects (examples) that are covered by a rule

Formal Concept:

- A rule that cannot be further extended without losing coverage of one of its covered examples (*maximal intent*)
- Along with *all* covered examples (*maximal extent*)
- Essentially, a formal concept is a maximal discriminative / characteristic rule

Features ↔ Class

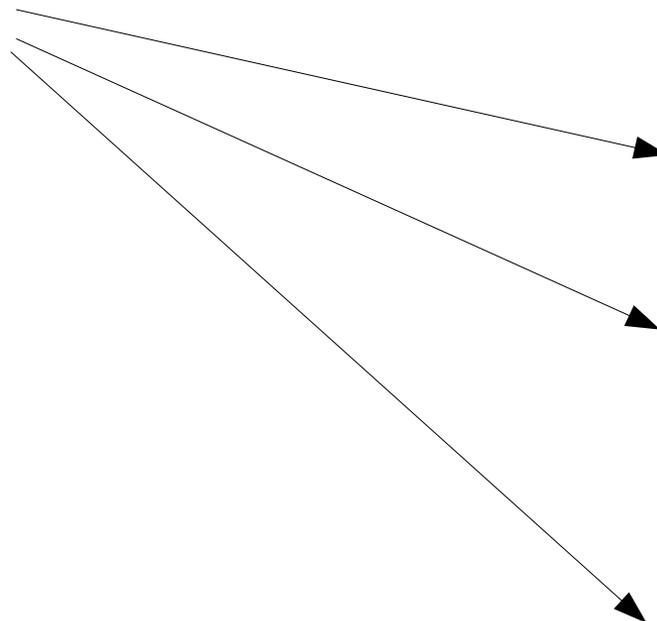


FCA Example

Concept “education = university”

- Maximal extent:

Education	Marital S.	Income
Primary	Single	Low
Primary	Single	Low
Primary	Married	Low
University	Divorced	High
University	Married	High
Secondary	Single	Low
University	Single	High
Secondary	Divorced	High
Secondary	Single	High
Secondary	Married	Low
Secondary	Divorced	Low
University	Divorced	High
Secondary	Divorced	Low



FCA Example

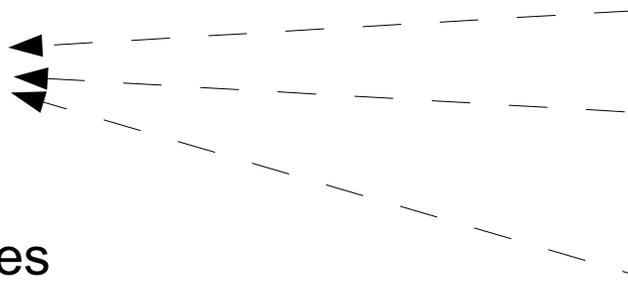
Concept “education = university”

- Maximal extent
 - All covered examples
- Maximal intent
 - All conditions common to the covered examples

→ Formal Concept

“Education = university
AND Income = high”

Education	Marital S.	Income
Primary	Single	Low
Primary	Single	Low
Primary	Married	Low
University	Divorced	High
University	Married	High
Secondary	Single	Low
University	Single	High
Secondary	Divorced	High
Secondary	Single	High
Secondary	Married	Low
Secondary	Divorced	Low
University	Divorced	High
Secondary	Divorced	Low



Closed Itemsets

In association rule discovery, formal concepts are called **closed itemsets**

- Although there is **no statistical difference** between an itemset and its closure (except for #items), their **interestingness may change**

Shopping Basket of a young family:



Itemset



Closed Itemset

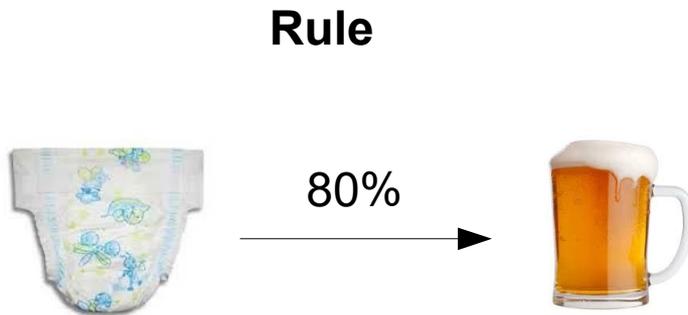


Rule Pruning

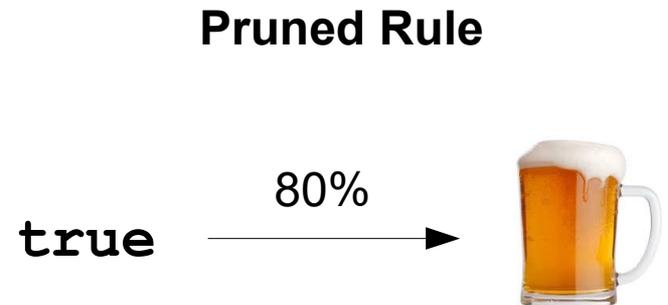
Rules are often pruned in order to get the shortest rule

- Remove conditions from the rule as long as the evaluation measure does not significantly change

This may also significantly change the semantics without changing the statistics



*80% of customers who buy
diapers also buy beer*



80% of all customers buy beer

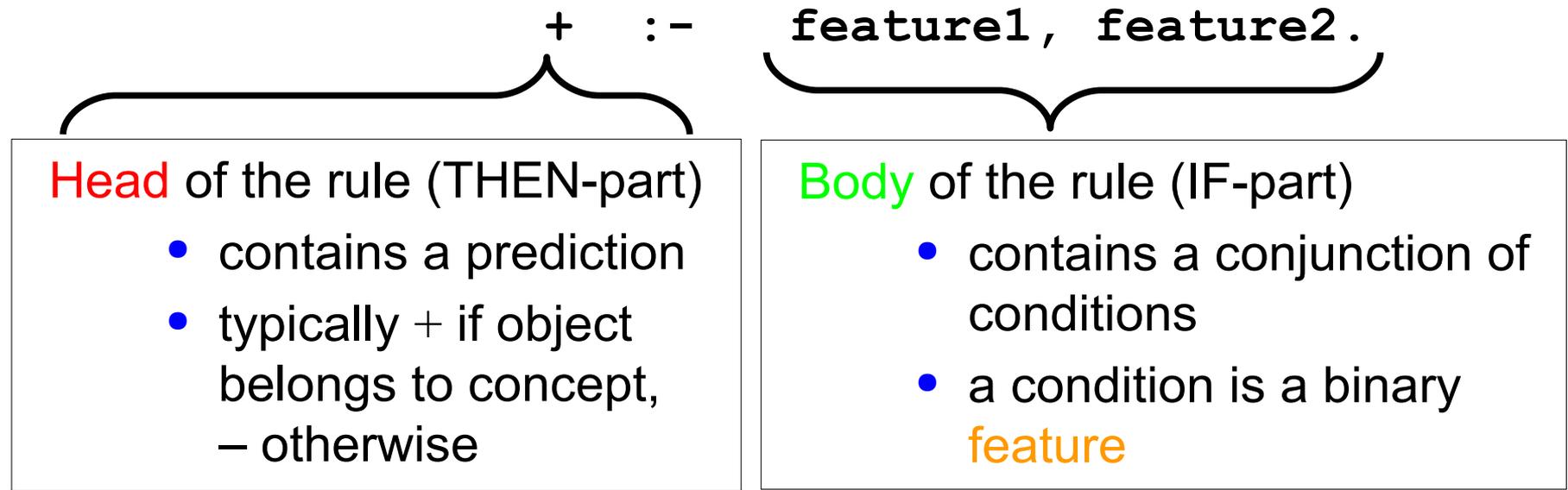


Overview

- Motivation
 - Understandability has not received much attention
- Understandability
 - Conjunctive Fallacy
 - Gambler's Fallacy
 - Representativeness Heuristic
- Different Types of Rules
 - Discriminative vs. Characteristic Rules
 - Formal Concepts
 - Closed Itemsets
- Heuristic Rule Learning
 - Concept Learning
 - Coverage Spaces
 - Rule Learning Heuristics
- Inverted Heuristics
- Explain-A-LOD
 - Semantic Coherence
 - Representation Heuristics
- Algorithmic Enhancements
 - Structured theories
 - More complex problems
- Conclusions



Conjunctive Rule



- Coverage
 - A rule is said to **cover** an example if the example satisfies the conditions of the rule.
- Prediction
 - If a rule covers an example, the rule's head is predicted for this example.



A Sample Database

No.	Education	Marital S.	Income	Children ?	Approved?
1	Primary	Single	Low	N	-
2	Primary	Single	Low	Y	-
3	Primary	Married	Low	N	+
4	University	Divorced	High	N	+
5	University	Married	High	Y	+
6	Secondary	Single	Low	N	-
7	University	Single	High	N	+
8	Secondary	Divorced	High	N	+
9	Secondary	Single	High	Y	+
10	Secondary	Married	Low	Y	+
11	Primary	Married	High	N	+
12	Secondary	Divorced	Low	Y	-
13	University	Divorced	High	Y	-
14	Secondary	Divorced	Low	N	+

Property of Interest
("class variable")



A Possible Solution

```
+ :- E=primary,      I=low,      M=married,  C=no.
+ :- E=university,  I=high,     M=divorced, C=no.
+ :- E=university,  I=high,     M=married,  C=no.
+ :- E=university,  I=high,     M=single,   C=no.
+ :- E=secondary,   I=high,     M=divorced, C=no.
+ :- E=secondary,   I=high,     M=single,   C=yes.
+ :- E=secondary,   I=low,      M=married,  C=yes.
+ :- E=primary,     I=high,     M=married,  C=no.
+ :- E=secondary,   I=low,      M=divorced, C=no.
```

The solution is

- a set of rules
- that is complete and consistent on the training examples
- but it does not generalize to new examples
- and is not easily understandable



A Better Solution

```
+ :- Marital = married.  
+ :- Marital = single,    Income = high.  
+ :- Marital = divorced, Children = no.
```

This solution is also

- a set of rules
- that is complete and consistent on the training examples
- but it does ~~not~~ generalize to new examples
- and is ~~not~~ easily understandable

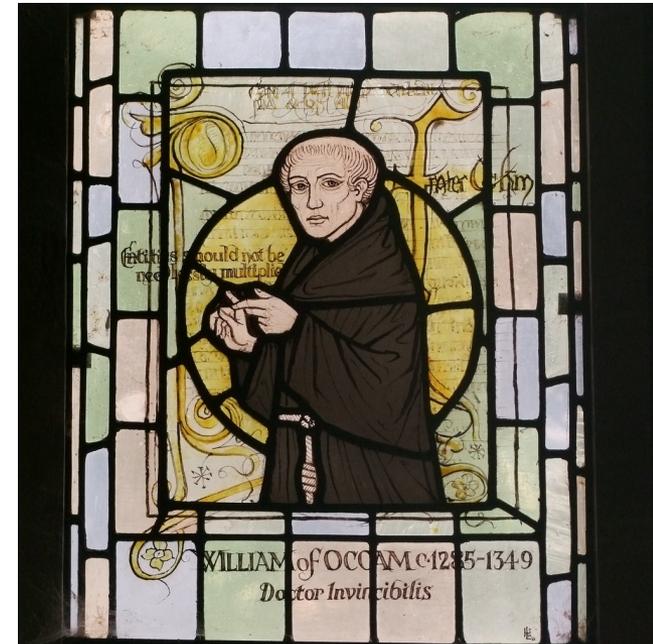


Occam's Razor

Entia non sunt multiplicanda sine necessitate.

William of Ockham (1285 - 1349)

- Machine Learning Interpretation:
 - Simple concepts are better
- (Debatable) Justifications:
 - There are more complex theories than simple theories, so that a simple theory is less likely to explain the data
 - Simpler theories are easier to falsify
- Empirically, we know that simpler theories perform better (overfitting)



Understandability and Rule length

- Humans sometimes prefer longer explanations
 - Bet on the longer sequence
- The reason seems to be that they do not operate on the logical level but construct a mental image of instances covered by a pattern
 - A prototypical object
 - They image a Linda, a bank teller, a feminist.
- While this has been observed in the context of fallible human reasoning, shouldn't we consider this when we talk about understandability?



Concept Learning

Given:

- **Positive Examples E^+**
 - examples for the concept to learn (e.g., people that approve an issue)
- **Negative Examples E^-**
 - counter-examples for the concept (e.g., people that don't approve)
- **Hypothesis Space H**
 - a (possibly infinite) set of candidate hypotheses
 - e.g., rules, rule sets, decision trees, linear functions, neural networks, ...

Find:

- Find the **target hypothesis $h \in H$**
- the target hypothesis is the concept that was used (or could have been used) to generate the training examples



Terminology

- training examples
 - P : total number of positive examples
 - N : total number of negative examples
- examples covered by the rule (predicted positive)
 - **true positives** p : positive examples covered by the rule
 - **false positives** n : negative examples covered by the rule
- examples not covered the rule (predicted negative)
 - **false negatives** $P-p$: positive examples not covered by the rule
 - **true negatives** $N-n$: negative examples not covered by the rule

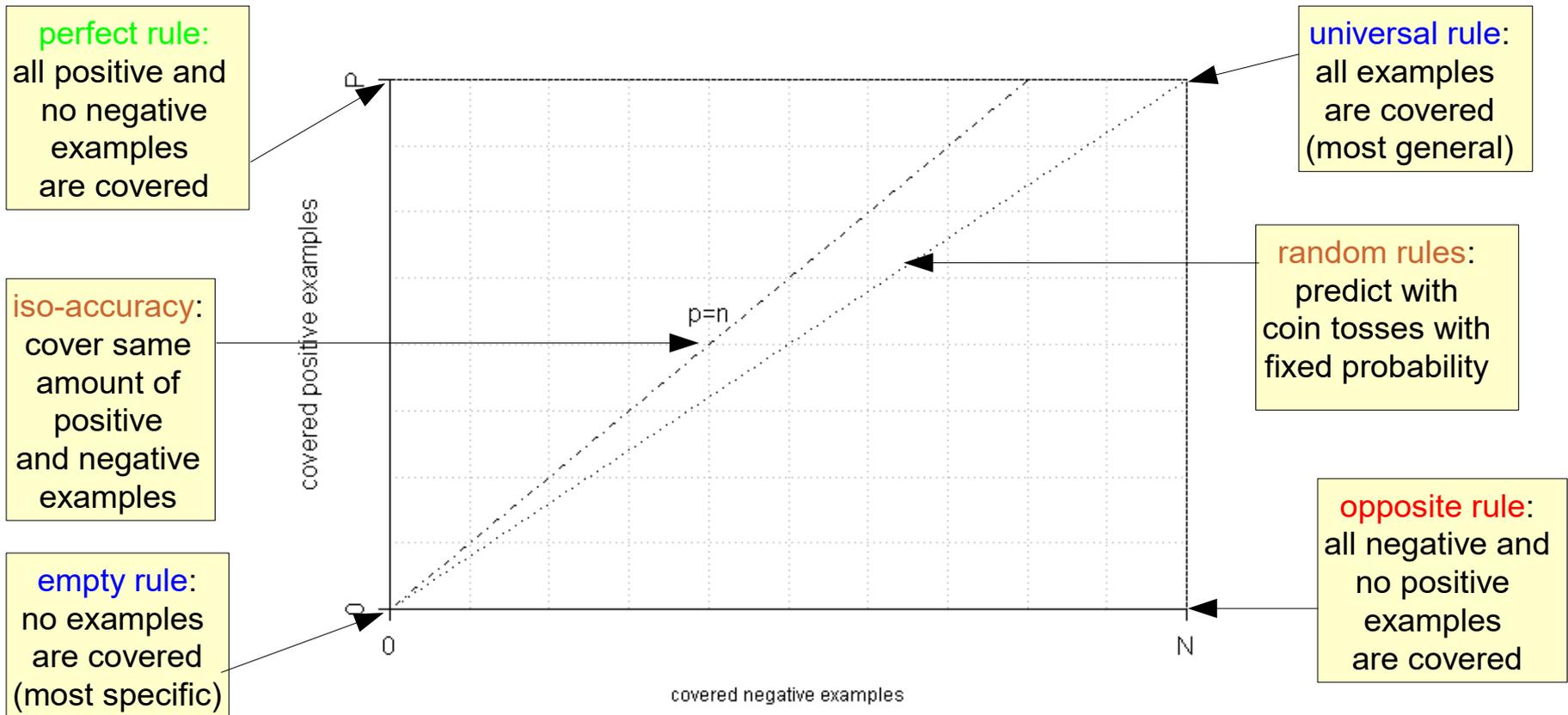
	predicted +	predicted -	
class +	p (true positives)	$P-p$ (false negatives)	P
class -	n (false positives)	$N-n$ (true negatives)	N
	$p + n$	$P+N - (p+n)$	$P+N$



Coverage Spaces

(Fürnkranz & Flach 2005)

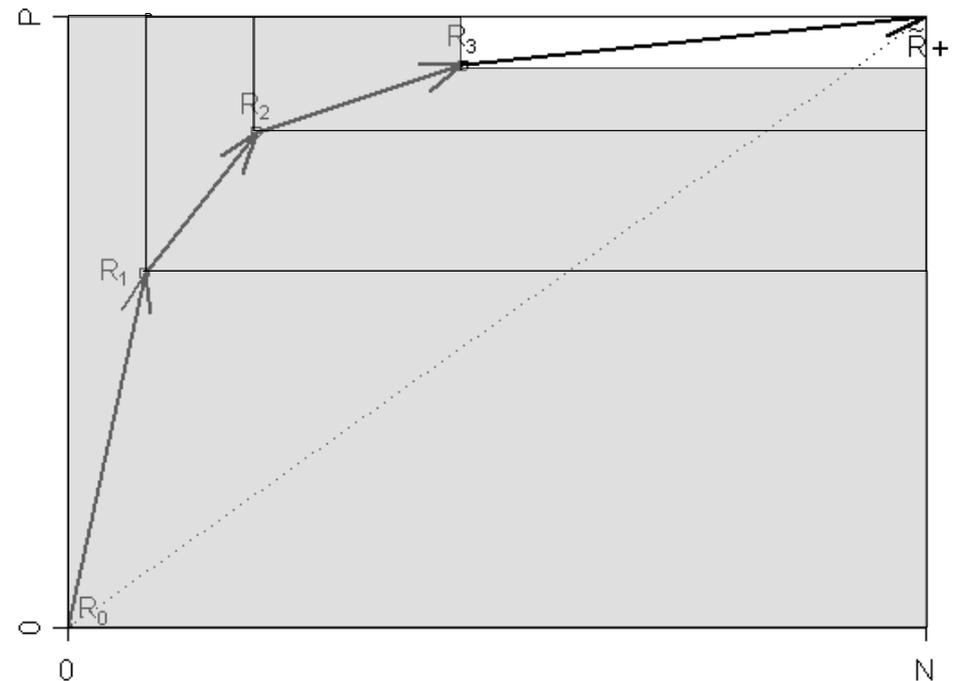
- good tool for visualizing properties of rule evaluation heuristics
 - each point is a rule covering p positive and n negative examples



Rule Selection: Covering Strategy

(survey → Fürnkranz 1999)

- **Covering** or **Separate-and-Conquer** rule learning algorithms learn one rule at a time
 - and then removes the examples covered by this rule
- This corresponds to a path in coverage space:
 - The **empty theory** R_0 (no rules) corresponds to $(0,0)$
 - Adding one rule **never decreases p or n** because adding a rule covers *more* examples (generalization)
 - The **universal theory** R_+ (all examples are positive) corresponds to (N,P)



Learning one Rule: Subgroup Discovery

■ Definition

“Given a population of individuals and a property of those individuals that we are interested in, **find population subgroups** that are statistically 'most interesting', e.g., are **as large as possible** and have the most **unusual distributional characteristics** with respect to the property of interest”

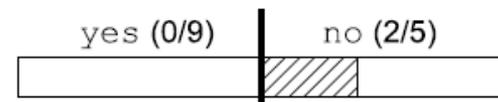
(Klösgen 1996; Wrobel 1997)

■ Examples

no :- MaritalStatus = single,
Income = Low.

yes :- MaritalStatus = married.

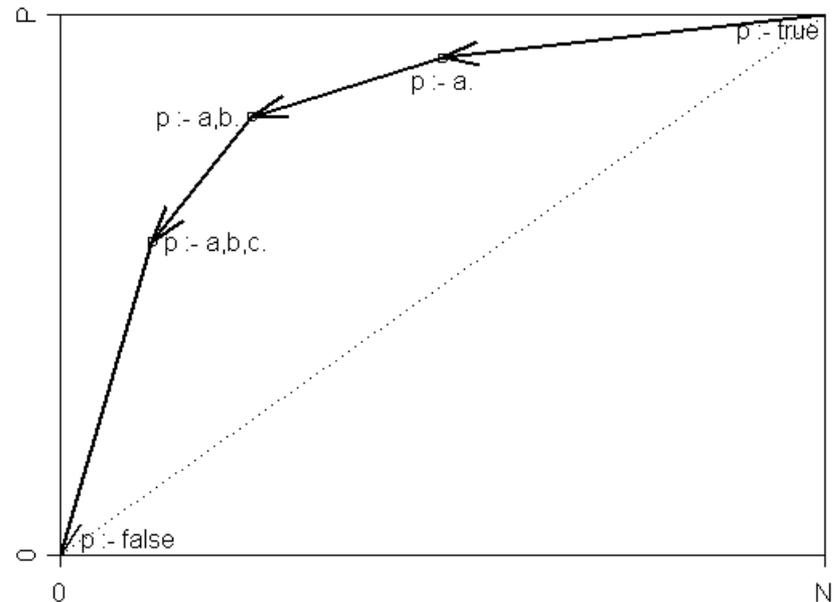
No :- MaritalStatus = divorced,
Children = yes.



Rule Refinement: Top-Down Hill-Climbing

- successively extends a rule by adding conditions

- This corresponds to a path in coverage space:
 - The rule $p :- \text{true}$ covers all examples (universal theory)
 - Adding a condition never increases p or n (specialization)
 - The rule $p :- \text{false}$ covers no examples (empty theory)



- which conditions are selected depends on a *heuristic function* that estimates the quality of the rule



Rule Learning Heuristics

- How can we measure the quality of a rule?
 - should cover as few negative examples as possible (*consistency*)
 - should cover as many positive examples as possible (*completeness*)
- An evaluation heuristic should therefore trade off these two properties
 - Example: **Laplace heuristic** $h_{Lap} = \frac{p+1}{p+n+2}$
 - grows with $p \rightarrow \infty$
 - grows with $n \rightarrow 0$
 - Example: **Precision** $h_{Prec} = \frac{p}{p+n}$
 - is not a good heuristic. Why?



Example

Condition		p	n	Precision	Laplace	p-n
Education =	Primary	2	2	0.5000	0.5000	0
	University	3	1	0.7500	0.6667	2
	Secondary	4	2	0.6667	0.6250	2
Marital =	Single	2	3	0.4000	0.4286	-1
	Married	4	0	1.0000	0.8333	4
	Divorced	3	2	0.6000	0.5714	1
Income =	Low	3	4	0.4286	0.4444	-1
	High	6	1	0.8571	0.7778	5
Children =	yes	3	3	0.5000	0.5000	0
	no	6	2	0.7500	0.7000	4

Heuristics Precision and Laplace

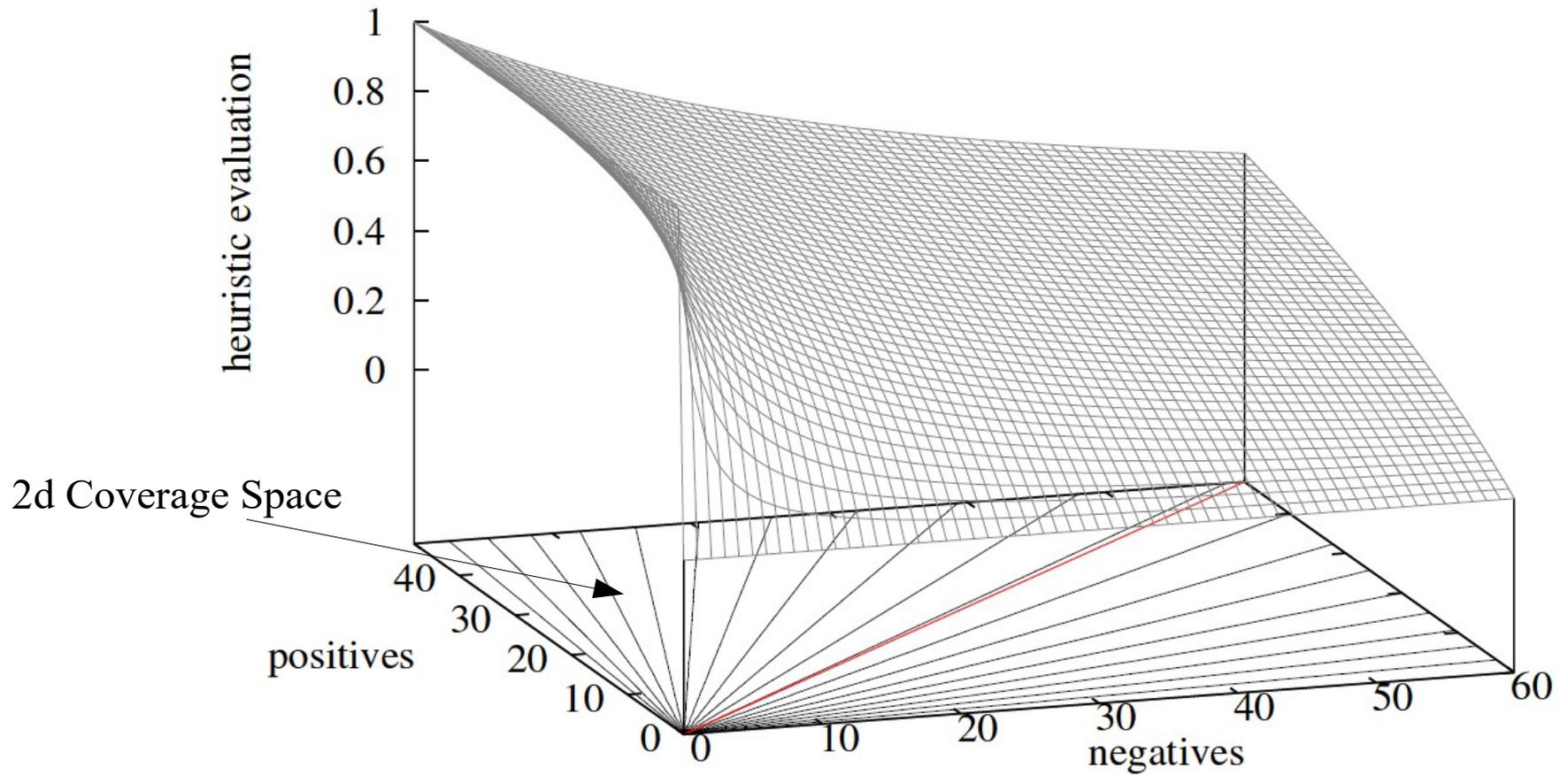
- add the condition **Marital = Married** to the (empty) rule
- stop and try to learn the next rule

Heuristic Accuracy / $p - n$

- adds Humidity = Normal
- continue to refine the rule (until no covered negative)

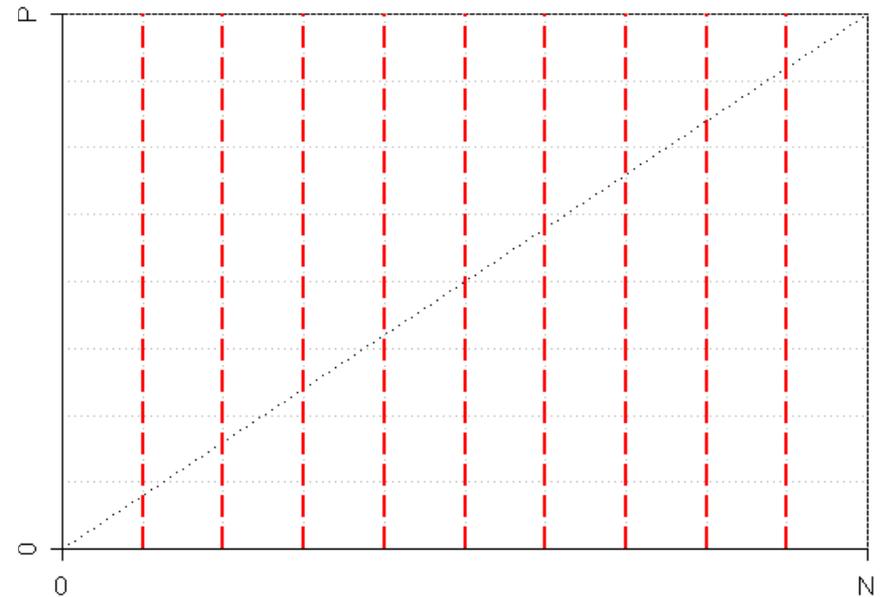
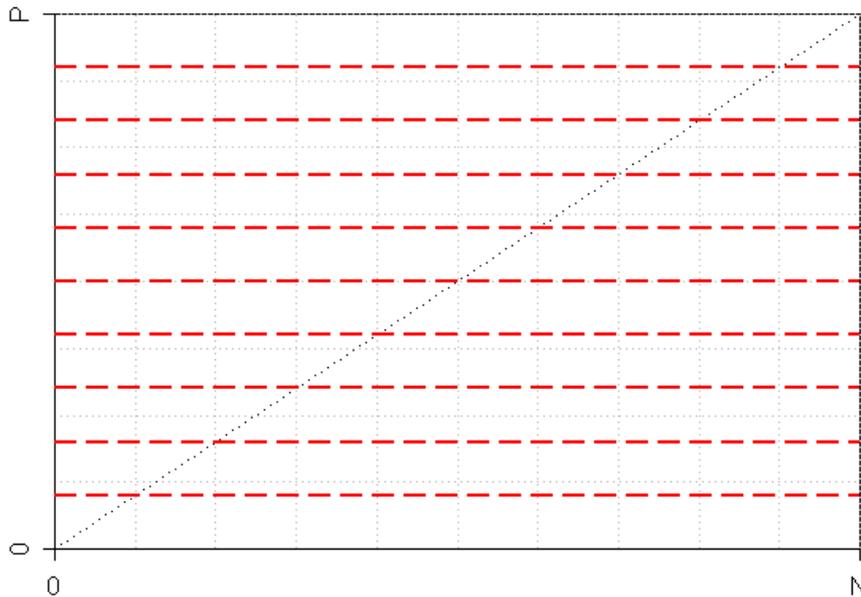


3d-Visualization of Precision



Isometrics in Coverage Space

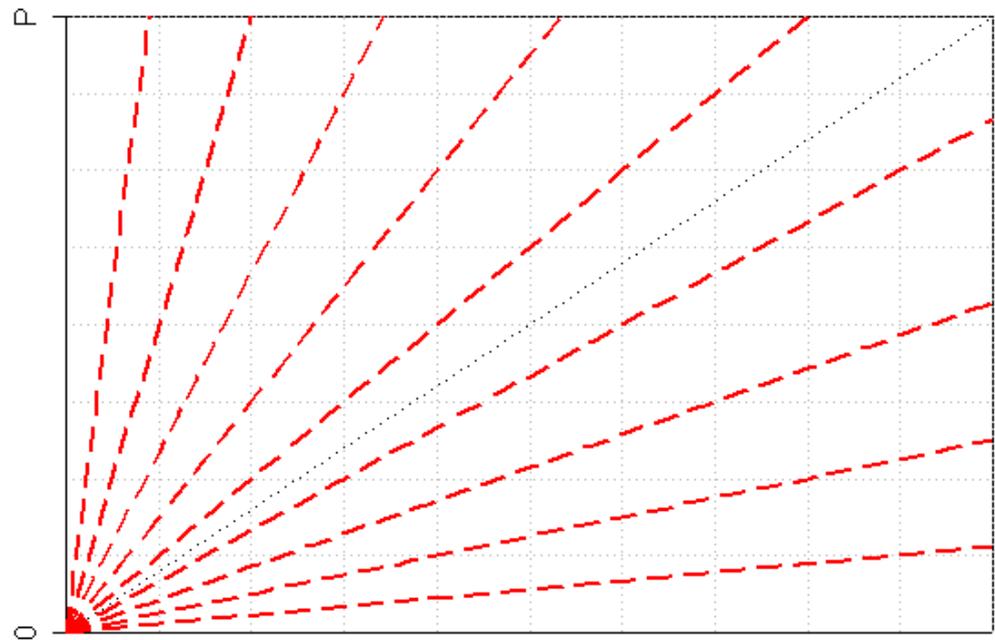
- Isometrics are lines that connect points for which a function in p and n has equal values
 - *Examples:*
Isometrics for heuristics $h_p = p$ and $h_n = -n$



Precision

- *basic idea:*
 - percentage of positive examples among covered examples
- *effects:*
 - rotation around origin (0,0)
 - all rules with same angle equivalent
 - in particular, all rules on P/N axes are equivalent
- typically **overfits**

$$h_{Prec} = \frac{p}{p+n}$$

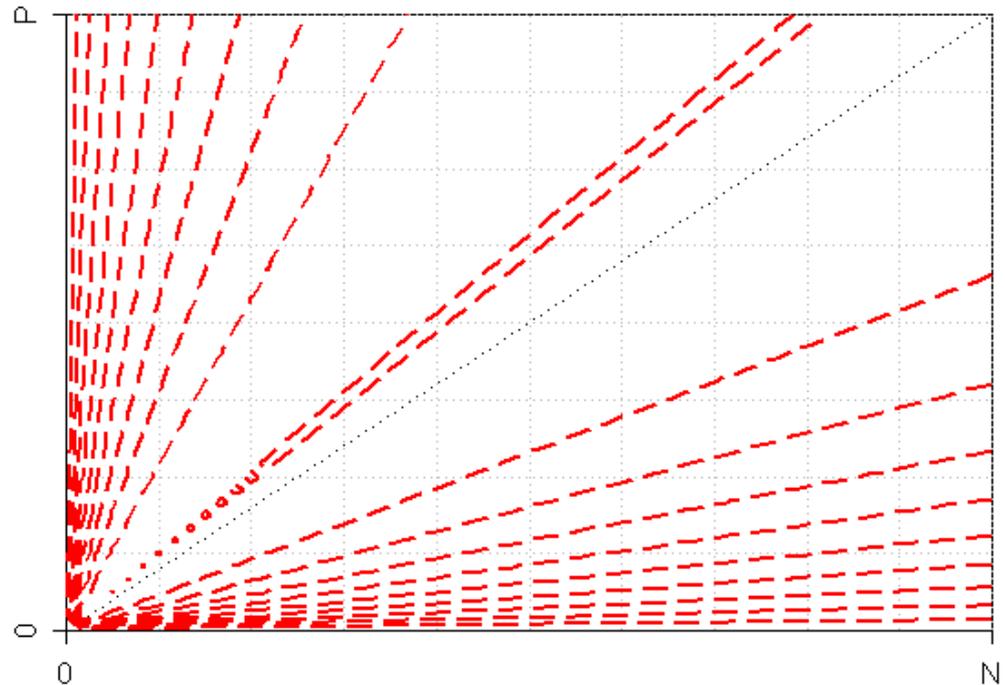


Entropy and Gini Index

$$h_{Ent} = -\left(\frac{p}{p+n} \log_2 \frac{p}{p+n} + \frac{n}{p+n} \log_2 \frac{n}{p+n}\right)$$

$$h_{Gini} = 1 - \left(\frac{p}{p+n}\right)^2 - \left(\frac{n}{p+n}\right)^2 \simeq \frac{pn}{(p+n)^2}$$

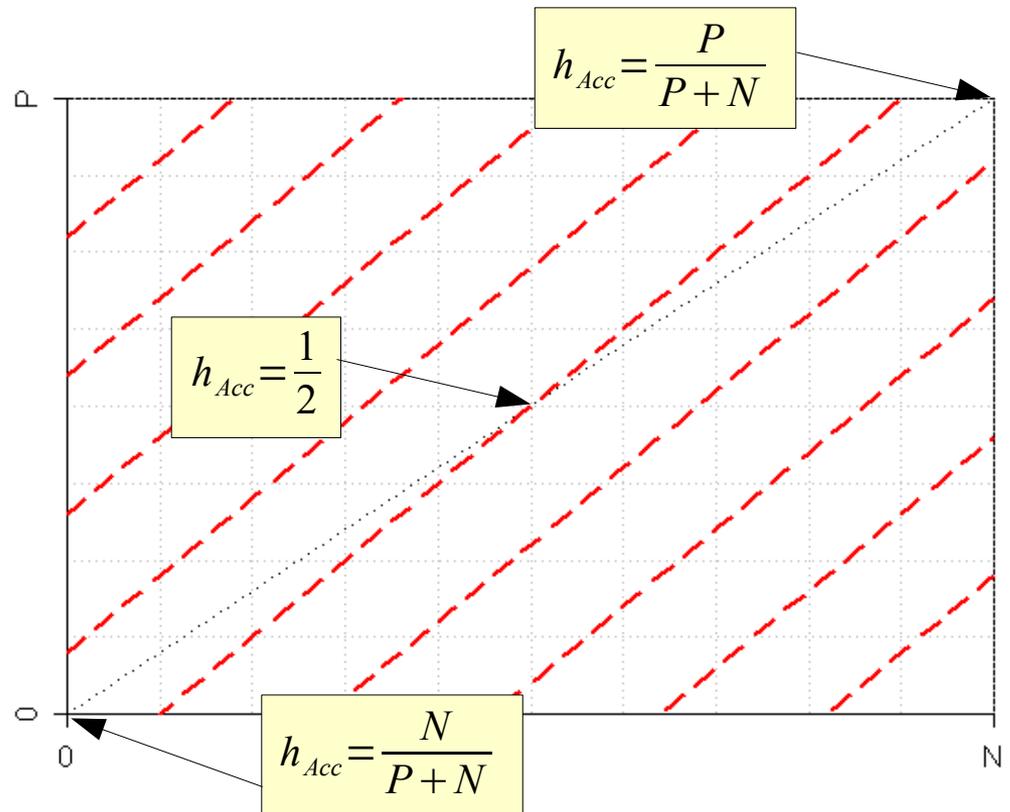
- *effects:*
 - entropy and Gini index are equivalent
 - like precision, isometrics rotate around (0,0)
 - isometrics are symmetric around 45° line
 - a rule that only covers negative examples is as good as a rule that only covers positives



Accuracy

$$h_{Acc} = \frac{p + (N - n)}{P + N} \simeq p - n$$

- *basic idea:*
percentage of correct classifications
(covered positives plus uncovered negatives)
- *effects:*
 - isometrics are parallel to 45° line
 - covering one positive example is as good as not covering one negative example



Weighted Relative Accuracy

- Two Basic ideas:
 - Precision Gain:** compare precision to precision of a rule that classifies all examples as positive

$$\frac{p}{p+n} - \frac{P}{P+N}$$

- Coverage:** Multiply with the percentage of covered examples

$$\frac{p+n}{P+N}$$

- Resulting formula:

$$h_{WRA} = \frac{p+n}{P+N} \cdot \left(\frac{p}{p+n} - \frac{P}{P+N} \right)$$

- one can show that this sorts rules in exactly the same way as

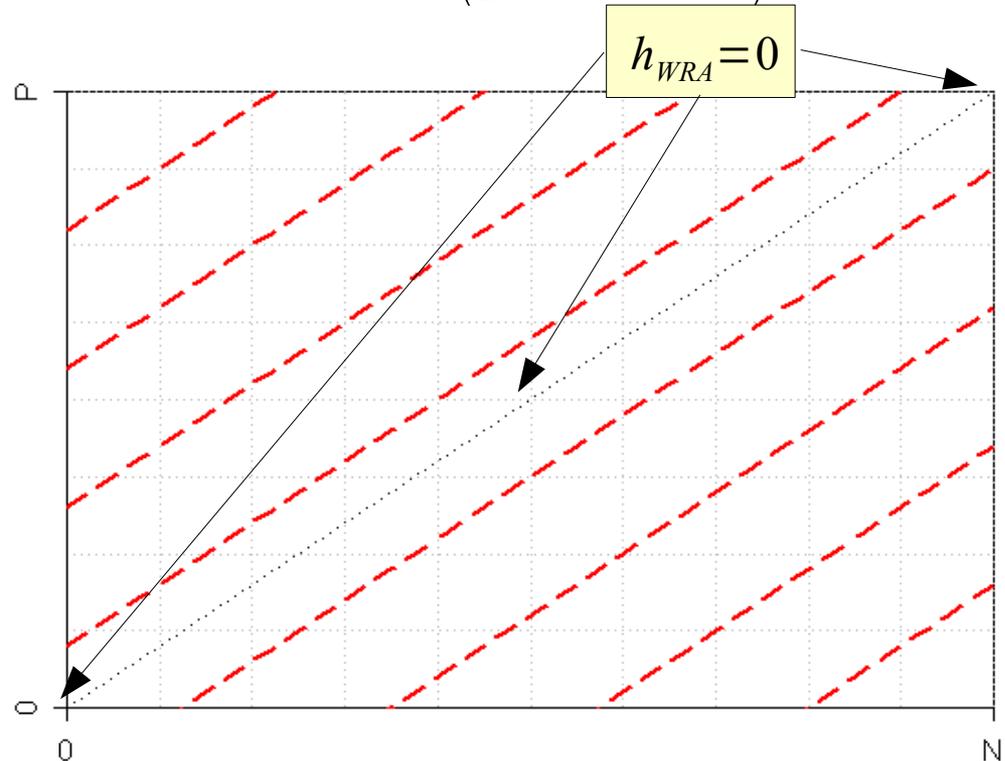
$$h_{WRA}' = \frac{p}{P} - \frac{n}{N}$$



Weighted relative accuracy

- *basic idea:*
compute the distance from the diagonal (i.e., from random rules)
- *effects:*
 - isometrics are parallel to diagonal
 - covering $x\%$ of the positive examples is considered to be as good as not covering $x\%$ of the negative examples
 - typically **over-generalizes**

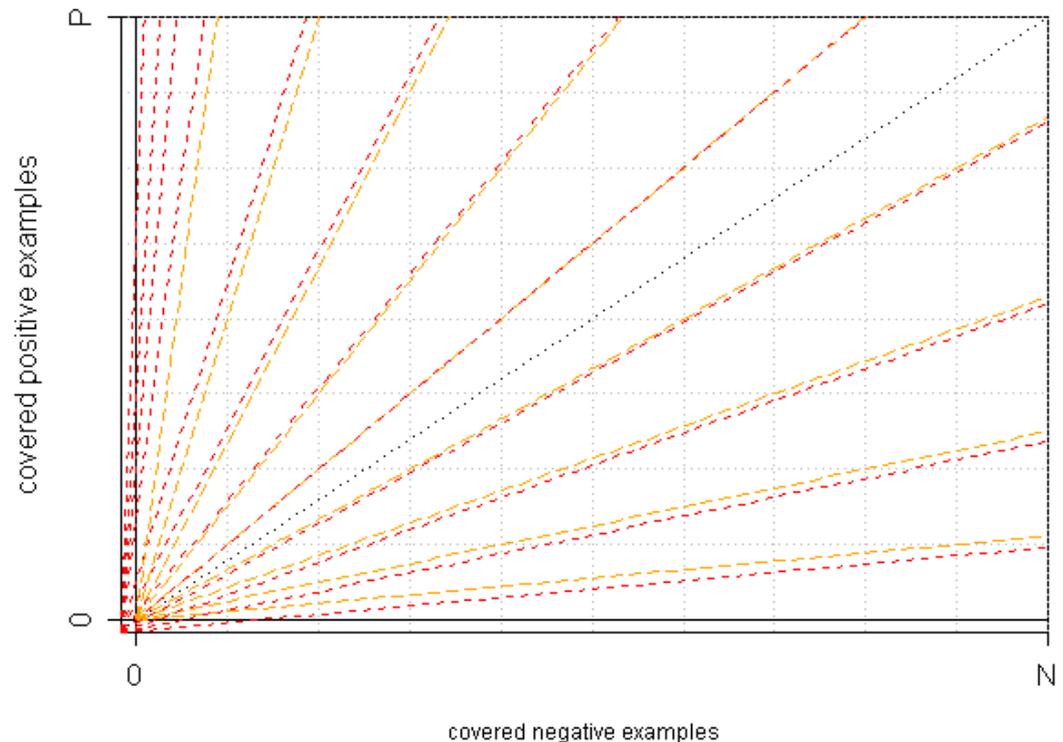
$$h_{WRA} = \frac{p+n}{P+N} \left(\frac{p}{p+n} - \frac{P}{P+N} \right) \approx \frac{p}{P} - \frac{n}{N}$$



Laplace-Estimate

- *basic idea:*
precision, but count coverage for positive and negative examples starting with 1 instead of 0
- *effects:*
 - origin at (-1,-1)
 - different values on $p=0$ or $n=0$ axes
 - not equivalent to precision

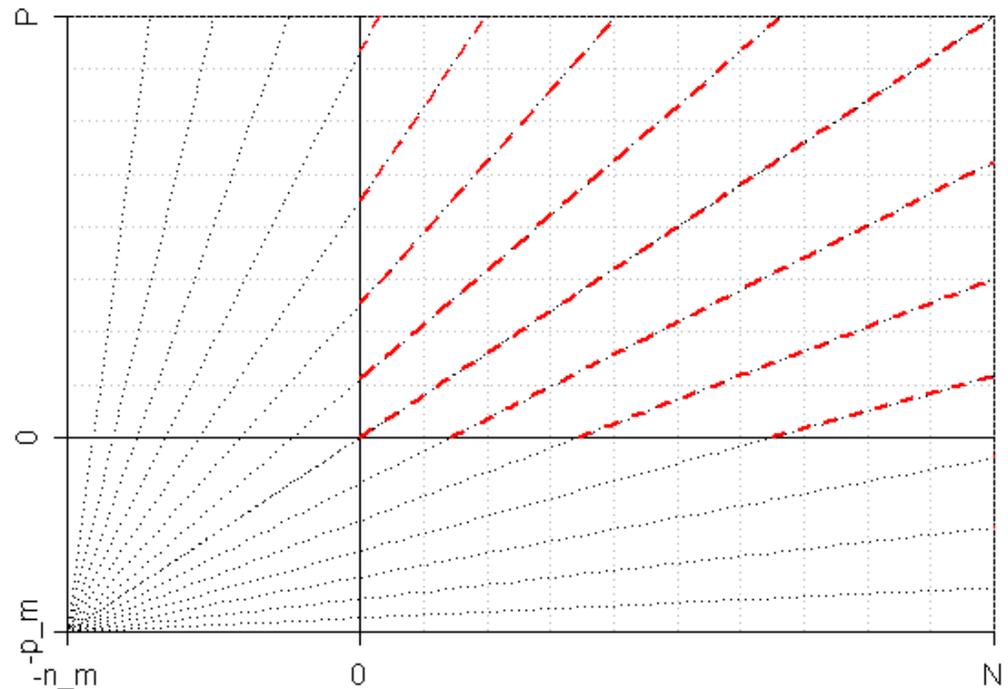
$$h_{Lap} = \frac{p+1}{(p+1)+(n+1)} = \frac{p+1}{p+n+2}$$



m-estimate

- *basic idea:*
initialize the counts with m examples in total,
distributed according to the
prior distribution $P/(P+N)$ of
 p and n .
- *effects:*
 - origin shifts to
 $(-mP/(P+N), -mN/(P+N))$
 - with increasing m , the
lines become more and
more parallel
- can be re-interpreted as a
**trade-off between WRA
and precision/confidence**

$$h_m = \frac{p + m \frac{P}{P+N}}{\left(p + m \frac{P}{P+N}\right) + \left(n + m \frac{N}{P+N}\right)} = \frac{p + m \frac{P}{P+N}}{p + n + m}$$

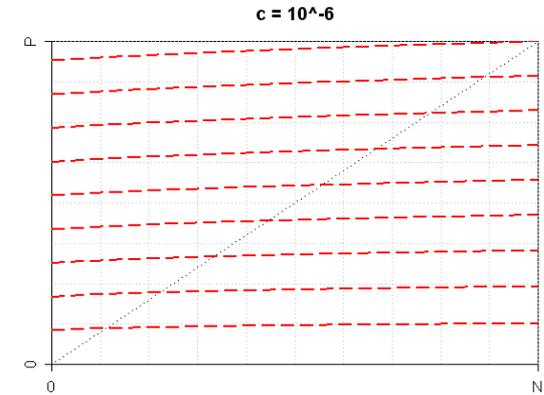
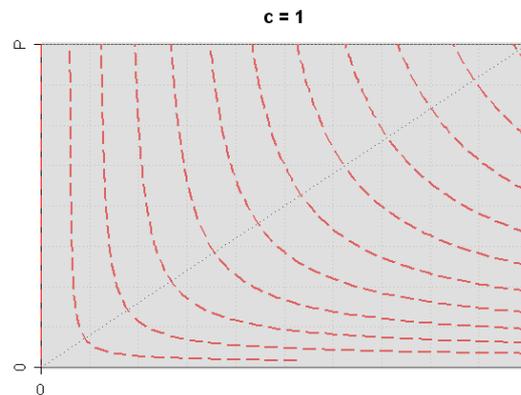
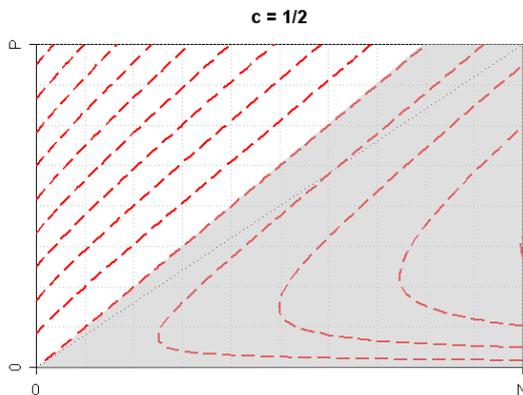
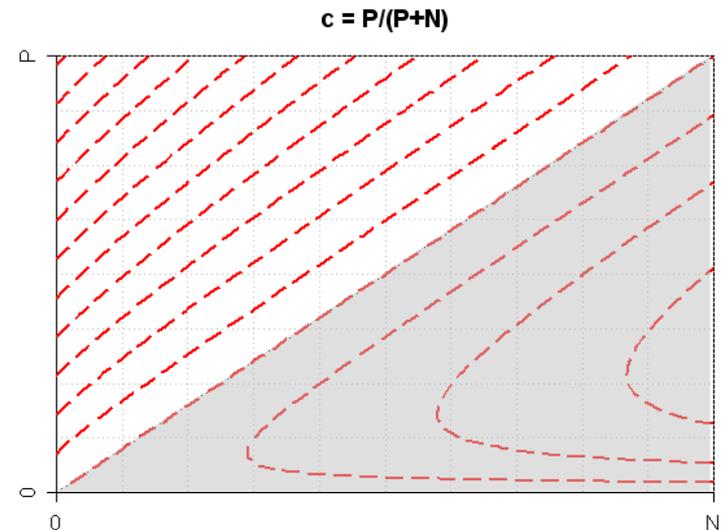


Non-Linear Isometrics – Foil's Information Gain

(Quinlan 1991)

$$h_{foil} = -p \left(\log_2 c - \log_2 \frac{p}{p+n} \right)$$

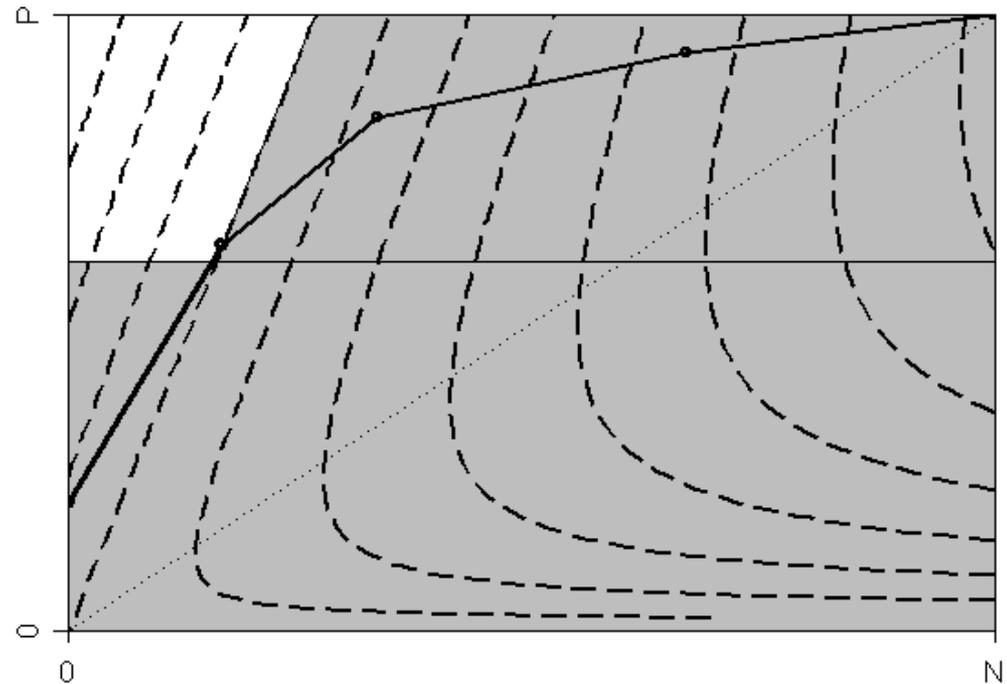
(c is the precision of the parent rule)



How Foil Works

→ Foil (almost) implements Support/Confidence Filtering

- filtering of rules with no information gain
 - after each refinement step, the region of acceptable rules is adjusted as in precision/confidence filtering
- filtering of rules that exceed rule length
 - after each refinement step, the region of acceptable rules adjusted as in support filtering



Overview

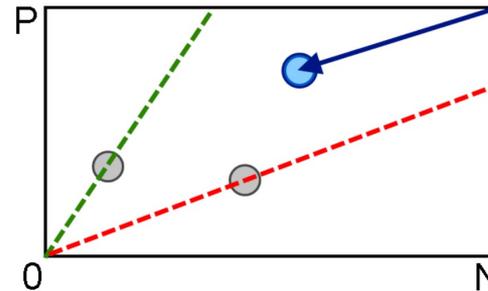
- Motivation
 - Understandability has not received much attention
- Understandability
 - Conjunctive Fallacy
 - Gambler's Fallacy
 - Representativeness Heuristic
- Different Types of Rules
 - Discriminative vs. Characteristic Rules
 - Formal Concepts
 - Closed Itemsets
- Heuristic Rule Learning
 - Concept Learning
 - Coverage Spaces
 - Rule Learning Heuristics
- Inverted Heuristics
- Explain-A-LOD
 - Semantic Coherence
 - Representation Heuristics
- Algorithmic Enhancements
 - Structured theories
 - More complex problems
- Conclusions



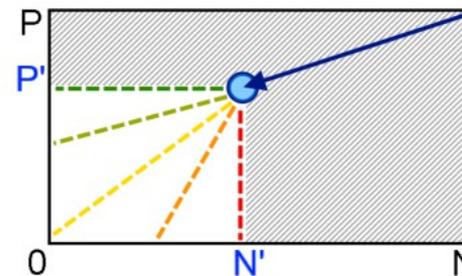
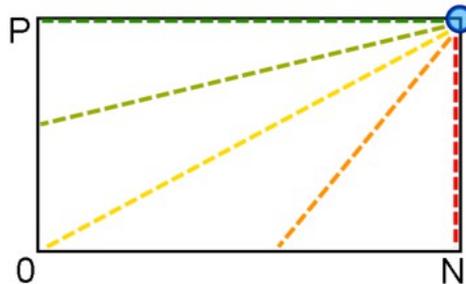
Inverted Heuristics – Motivation

(Stecher, Janssen, Fürnkranz 2014)

- While the search proceeds top-down
- the evaluation of refinements happens from the point of view of the origin (bottom-up)

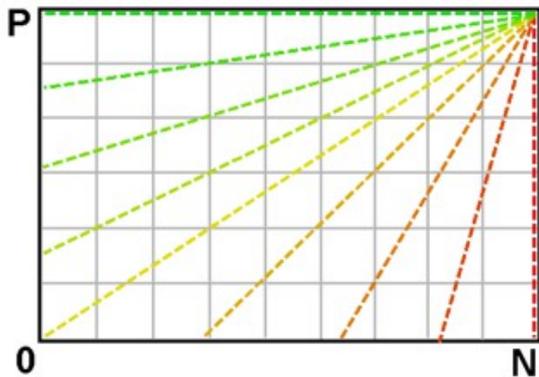


- Instead, we want to evaluate the refinement from the point of view of the predecessor

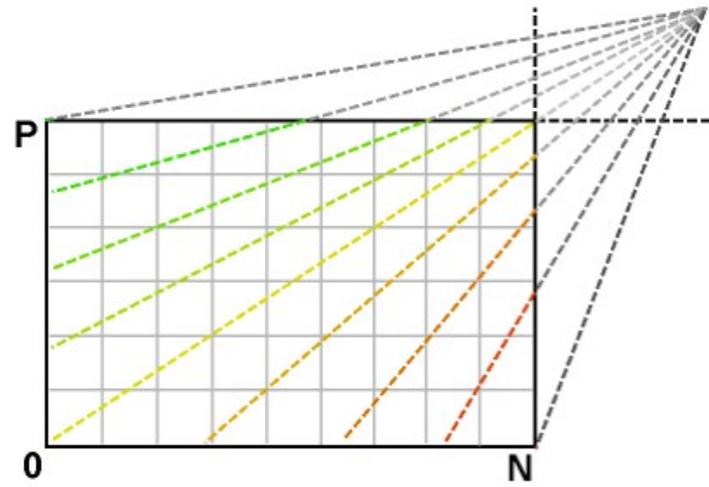


Inverted Heuristics

- Many heuristics can be “inverted” by replacing changing their angle point from the origin to the current rule



$$h'_{precision}(p, n, P, N) = \frac{N - n}{(P + N) - (p + n)}$$



$$h'_{m-Estimate}(p, n, P, N) = \frac{N - n + m \cdot \frac{P}{P + N}}{(P + N) - (p + n - m)}$$

- Note:** not all heuristics can be inverted
 - e.g. WRA is invariant w.r.t. inversion (because of symmetry)

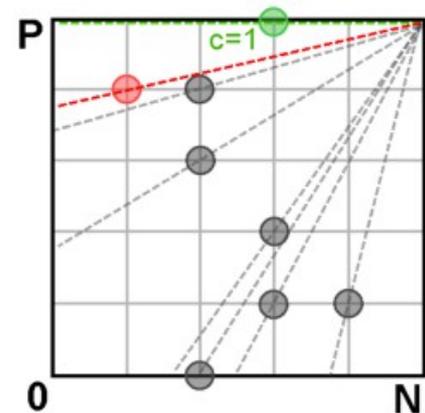
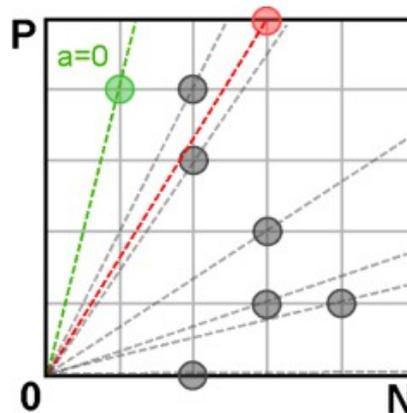


Inverted Heuristics – Example

First refinement step in small example dataset

- 4 Attributes, 10 data points, binary-class

a	b	c	d	C
0	1	1	1	+
0	1	1	1	+
0	0	1	0	-
1	1	1	0	-
1	0	0	1	-
0	1	1	0	+
0	0	1	1	+
1	1	1	0	-
1	0	1	1	+
1	0	0	1	-

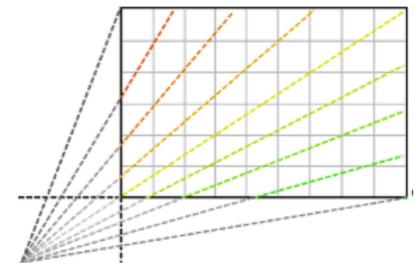
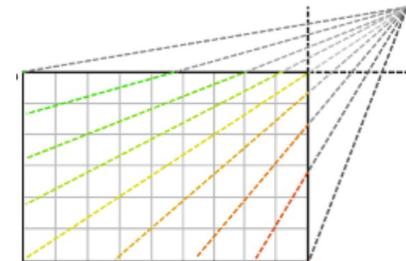


Inverted heuristic function (right image) selects preferable refinement condition $c=1$ with coverage of $(p, n)=(5,3)$



Implementation

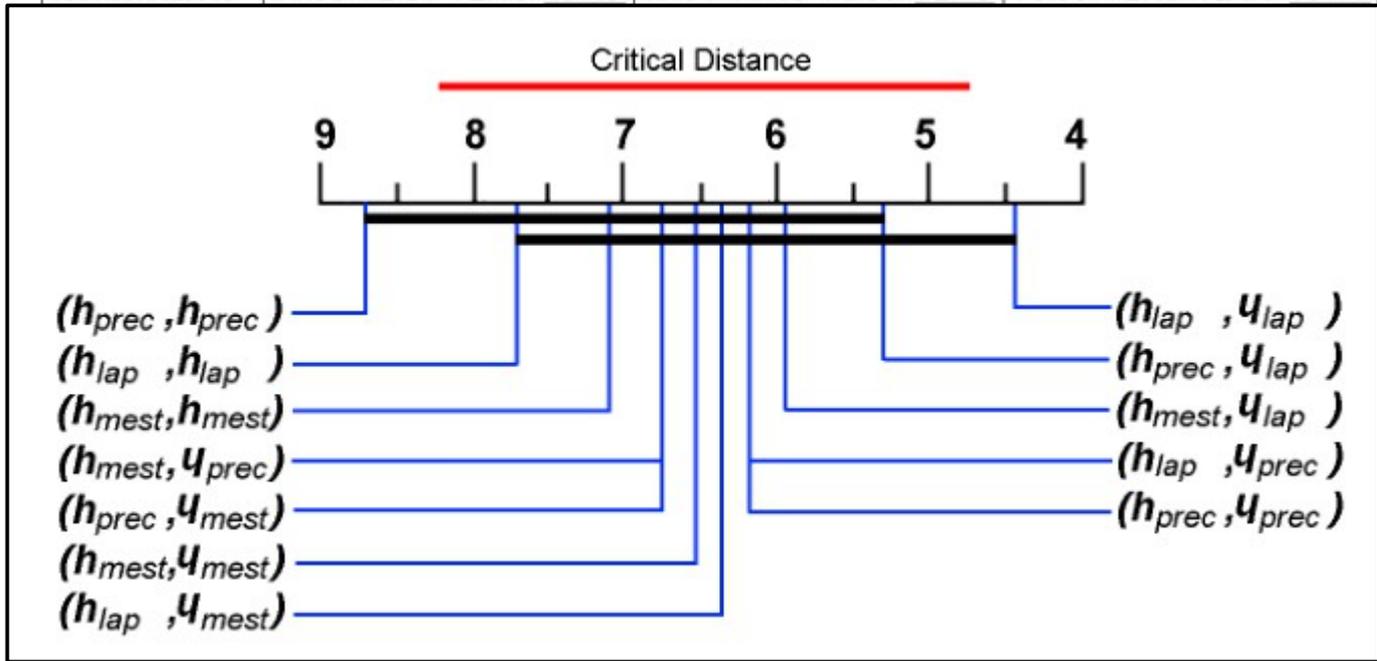
- Modification of a conventional covering algorithm
 - CN2-like
 - No pruning, no significance test
- **Rule refinement** proceeds with **inverted heuristics**
 - In each iteration, the best condition is added to the rule until the rule covers no more examples
- **Rule selection** proceeds with **regular heuristics**
 - Among all refinements on the path, the best rule is selected using a regular heuristic



Results:

Inverted heuristics tend to work better

Dataset	(h_{prec}, \cdot)				(h_{lap}, \cdot)				(h_{mest}, \cdot)			
	h_{prec}	u_{prec}	u_{lap}	u_{mest}	h_{lap}	u_{prec}	u_{lap}	u_{mest}	h_{mest}	u_{prec}	u_{lap}	u_{mest}
breast-cancer	68.53	72.38	72.03	73.43	69.58	70.63	71.33	72.73	71.33	72.03	72.38	73.78



soybean	90.84	91.84	92.24	91.88	90.84	91.88	92.27	90.88	91.84	90.92	90.18	91.88
tic-tac-toe	97.39	98.02	97.60	97.81	97.60	98.02	97.60	97.91	98.12	98.02	97.60	97.81
vote	94.94	93.56	94.25	94.48	95.40	94.25	94.25	94.94	93.33	93.56	94.71	96.09
zoo	84.16	88.12	92.08	90.01	86.14	88.12	92.08	90.10	89.11	88.12	92.08	90.10
average rank	3.075	2.400	1.975	2.550	3.000	2.500	1.975	2.525	2.700	2.625	2.225	2.450



Inverted Heuristics – Rule Length

- Inverted Heuristics tend to learn longer rules
 - If there are conditions that can be added without decreasing coverage on the positive examples, inverted heuristics will add them first (before adding discriminative conditions)

Dataset	(h_{lap}, h_{lap})		(h_{lap}, h'_{lap})		Dataset	(h_{lap}, h_{lap})		(h_{lap}, h'_{lap})	
	R	L	R	L		R	L	R	L
breast-cancer	25	67	38	173	ionosphere	17	25	8	42
car	107	495	107	506	labor	5	7	3	12
contact-lenses	5	14	5	15	lymphography	18	42	11	47
futebol	4	7	2	5	monk3	13	38	11	32
glass	50	103	14	83	mushroom	11	13	7	35
hepatitis	13	26	7	46	primary-tumor	80	319	72	518
horse-colic	44	114	19	111	soybean	62	134	45	195
hypothyroid	27	65	9	69	tic-tac-toe	22	84	16	69
iris	7	15	5	17	vote	13	48	12	58
idh	4	5	2	5	zoo	19	19	6	14
averages						28.2	85.6	20.6	106.2



Example: Mushroom dataset

- The best three rules learned with conventional heuristics

```
poisonous :- odor = foul. (2160,0)
poisonous :- gill-color = buff. (1152,0)
poisonous :- odor = pungent. (256,0)
```



- The best three rules learned with inverted heuristics

```
poisonous :- veil-color = white, gill-spacing = close,
no bruises, ring-number = one,
stalk-surface-above-ring = silky. (2192,0)
poisonous :- veil-color = white, gill-spacing = close,
gill-size = narrow, population = several,
stalk-shape = tapering. (864,0)
poisonous :- stalk-color-below-ring = white,
ring-type = pendant, ring-number = one,
stalk-color-above-ring = white,
cap-surface = smooth, stalk-root = bulbuous,
gill-spacing = close. (336,0)
```



Overview

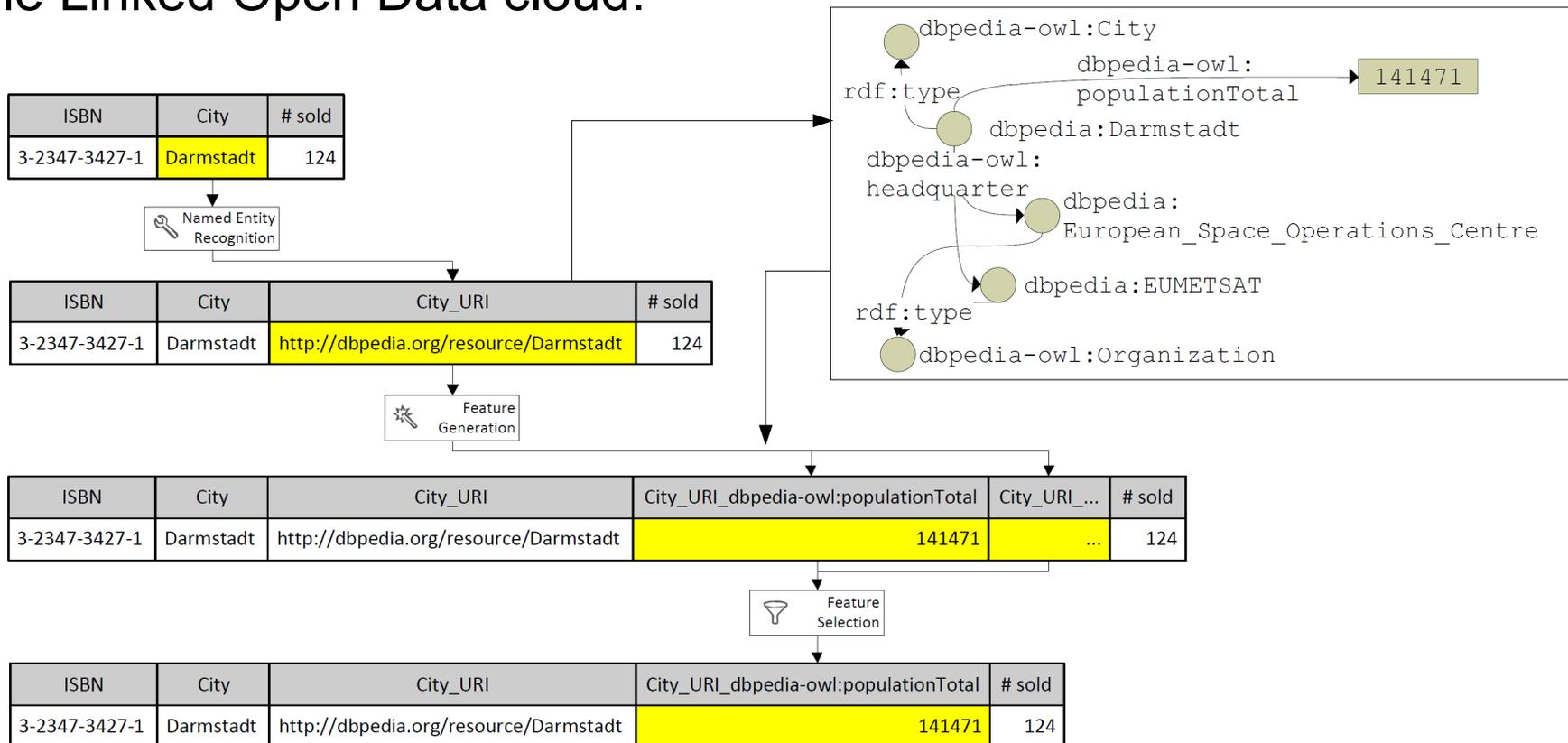
- Motivation
 - Understandability has not received much attention
- Understandability
 - Conjunctive Fallacy
 - Gambler's Fallacy
 - Representativeness Heuristic
- Different Types of Rules
 - Discriminative vs. Characteristic Rules
 - Formal Concepts
 - Closed Itemsets
- Heuristic Rule Learning
 - Concept Learning
 - Coverage Spaces
 - Rule Learning Heuristics
- Inverted Heuristics
- Explain-A-LOD
 - Semantic Coherence
 - Representation Heuristics
- Algorithmic Enhancements
 - Structured theories
 - More complex problems
- Conclusions



Explain-A-LOD

(Paulheim & Fürnkranz 2012)

- Generates features for data mining using features derived from the Linked Open Data cloud.

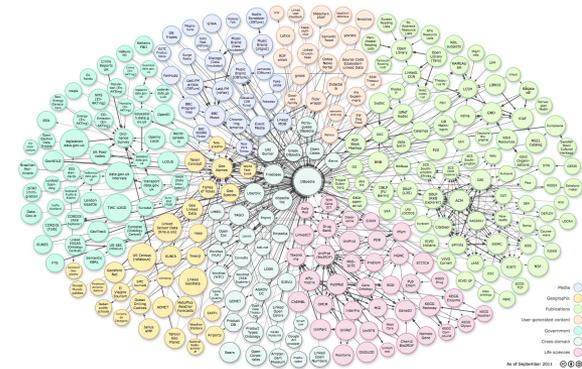


Zero-Knowledge Data Mining

(Paulheim 2012)

Mine a database without explicit background knowledge

City	Country	Index 2010
Vienna	 Austria	108.6
Zürich	 Switzerland	108.0
Auckland	 New Zealand	107.4
Munich	 Germany	107.0
Vancouver	 Canada	107.4
Düsseldorf	 Germany	107.2
Frankfurt	 Germany	107.0
Geneva	 Switzerland	107.9
Copenhagen	 Denmark	106.2
Sydney	 Australia	106.3



Quality-of-living
Index

LOD

QOL = High :-
European capital of culture



Some more rules

(Paulheim 2012)

Good discriminative rules, highly rated by users:

- `QOL = High :- Many events take place.`
- `QOL = High :- Host City of Olympic Summer Games.`
- `QOL = Low :- African Capital.`

Good discriminative rules, but lowly rated by users:

- `QOL = High :- # Records Made >= 1,
Companies/Organisations >= 22.`
- `QOL = High :- # Bands >= 18,
Airlines founded in 2000 > 1.`
- `QOL = Low :- # Records Made = 0,
Average January Temp <= 16.`



Semantic Coherence

Rule discovery algorithms only check the discriminative power of a condition to be added

- First world / Third world would be a plausible distinction
- A distinction based on latitude is less plausible
- A distinction based on number of records made even less plausible

→ conditions that may **cover the same examples** may have a **different “degree of understandability”**.

Similarly, combinations of conditions that are semantically far, do not appear to be plausible.

- Number of records made and number of companies are **coherent**
- Number of companies and average temperature are **not coherent**



Recognition Heuristic

(Gigerenzer & Todd 1999)

Which of the two cities is larger?

Chongqing



Chengdu



Recognition Heuristic

(Gigerenzer & Todd 1999)

Which of the two cities is larger?

Hongkong



Chengdu



Recognition Heuristic

(Gigerenzer & Todd 1999)

“if one of two objects is recognized and the other is not, then infer that the **recognized object has the higher value** with respect to the criterion”

Hongkong



ca. 7,000,000

Chengdu



ca. 15,000,000

Chongqing



ca. 30,000,000



Overview

- Motivation
 - Understandability has not received much attention
- Understandability
 - Conjunctive Fallacy
 - Gambler's Fallacy
 - Representativeness Heuristic
- Different Types of Rules
 - Discriminative vs. Characteristic Rules
 - Formal Concepts
 - Closed Itemsets
- Heuristic Rule Learning
 - Concept Learning
 - Coverage Spaces
 - Rule Learning Heuristics
- Inverted Heuristics
- Explain-A-LOD
 - Semantic Coherence
 - Representation Heuristics
- Algorithmic Enhancements
 - Structured theories
 - More complex problems
- Conclusions



Structured Concepts

Most rule learning algorithms learn flat theories

- e.g., n-bit parity needs 2^n flat rules

```
+ :- x1, x2, x3, x4.  
+ :- x1, x2, not x3, not x4.  
+ :- x1, not x2, x3, not x4.  
+ :- x1, not x2, not x3, x4.  
+ :- not x1, x2, not x3, x4.  
+ :- not x1, x2, x3, not x4.  
+ :- not x1, not x2, x3, x4.  
+ :- not x1, not x2, not x2, not x4.
```

But structured concepts are often more interpretable

- e.g. only $O(n)$ rules with intermediate concepts

```
+           :- x1, not parity234.  
+           :- not x1, parity234.  
  
parity234  :- x2, not parity34.  
parity234  :- not x2, parity34.  
  
parity34   :- x3, x4.  
parity34   :- not x3, not x4.
```

Previous work in the 90s in ILP and restructuring knowledge bases was not successful

- new approaches could borrow ideas from Deep Learning



Understandability in more complex problem settings

(Loza & Janssen 2016)

Multilabel Rule Learning

- The key challenge in multi-label classification is to model the dependencies between the labels
 - much of current research in this area is devoted to this topic
- Rules can make these dependencies explicit and exploit them in the learning phase
 - regular rule: **university, female** → **quality, fashion**
 - label dependency: **fashion** ≠ **sports**
 - mixed rule: **university, tabloid** → **quality**



Steps Towards More Understandable Rule Learning Algorithms

Understandability of the learned rules should be explicitly considered in rule learning algorithms

1. Understand Understandability
 - Take a closer look at results from cognitive science
2. Develop heuristics that include understandability
 - Of course, discriminative power should not be ignored
3. Integrate them in Rule Learning Algorithms
 - Possibly also as a post-processor (“rule beautification”)
4. Develop better algorithms
 - E.g., for learning structured and multi-label concepts
5. Evaluate in user studies
 - Automatic evaluation would not be convincing



Conclusions

- Understandability is currently mostly defined via rule length
 - Occam's Razor: Shorter rules are better
 - On the other hand, longer rules are often more convincing
 - Characteristic rules, closed itemsets, formal concepts, rules learned with inverted heuristics, ...
 - Understandability is more than short rules, e.g.
 - **Representativeness**: a rule that is more typical to what we expect is more convincing
 - **Semantic coherence**: rules that have semantically similar conditions are better
 - **Recognition**: rules with well-recognized conditions are better
 - **Structure**: flat rules are not very natural
- these should be considered when evaluating understandability!



References

- Fürnkranz J.: Separate-and-Conquer Rule Learning. *Artificial Intelligence Review* 13(1): 3-54 (1999)
- Fürnkranz J., Flach, P.: ROC 'n' Rule Learning - Towards a Better Understanding of Covering Algorithms. *Machine Learning* 58(1): 39-77 (2005)
- Fürnkranz J., Gamberger D., Lavrac N.: Foundations of Rule Learning. Springer (2012)
- Fürnkranz J., Kliegr T.: A Brief Overview of Rule Learning. *Proceedings RuleML 2015*: 54-69 (2015)
- Gigerenzer G., Todd. P.M.: Simple Heuristics that make us smart. Oxford University Press (1999)
- Kahneman, D. and Tversky, A. (1972). "Subjective probability: A judgment of representativeness", *Cognitive Psychology* 3 (3): 430–454
- Loza Mencía, E. and Janssen, F. Learning rules for multi-label classification: a stacking and a separate-and-conquer approach, Machine Learning, 2016
- Michalski, R.S.: A Theory and Methodology of Inductive Learning. *Artificial Intelligence* 20(2): 111-161 (1983)
- Paulheim, H.: Generating Possible Explanations for Statistics from Linked Open Data. In: 9th Extended Semantic Web Conference (ESWC-12), (2012)
- Paulheim, H., Fürnkranz, J.: Unsupervised Feature Generation from Linked Open Data. In: International Conference on Web Intelligence, Mining, and Semantics (WIMS'12). (2012)
- Stecher J., Janssen F., Fürnkranz J.: Separating Rule Refinement and Rule Selection Heuristics in Inductive Rule Learning. *Proceedings ECML/PKDD (3) 2014*: 114-129 (2014)
- Tversky A. and Kahneman, D.: "Extensional versus Intuitive Reasoning: The Conjunction Fallacy in Probability Judgement" in *Psychological Review*, 91, pp. 293-315, (1984)
- Tversky, A. and Kahneman, D. "Extension versus intuitive reasoning: The conjunction fallacy in probability judgment". *Psychological Review* 90 (4): 293–315 (1983)
- Wille R.: Restructuring lattice theory: an approach based on hierarchies of concepts. In: I. Rival (Ed.): Ordered Sets, 445–470, Reidel, Dordrecht-Boston (1982)

