Advanced Artificial Intelligence

Part II. Statistical NLP

Probabilistic Logic Learning

Wolfram Burgard, Luc De Raedt, Bernhard Nebel, Lars Schmidt-Thieme

Many slides taken from Kristian Kersting and for Logic From Peter Flach's Simply Logical

Overview

- Expressive power of PCFGs, HMMs, BNs still limited
 - First order logic is more expressive
- Why not combine logic with probabilities ?
 - Probabilistic logic learning
- Short recap of logic (programs)
- Stochastic logic programs
 - Extend PCFGs
- Bayesian logic programs
 - Extend Bayesian Nets
- Logical HMMs
 - Extend HMMs

Context

One of the key open questions of artificial intelligence concerns

"probabilistic logic learning",

i.e. the integration of probabilistic reasoning with first order logic representations and machine learning.



Sometimes called Statistical Relational Learning

So far

- We have largely been looking at probabilistic representations and ways of learning these from data
 - BNs, HMMs, PCFGs
- Now, we are going to look at their expressive power, and make traditional probabilistic representations more expressive using logic
 - Probabilistic First Order Logics
 - Lift BNs, HMMs, PCFGs to more expressive frameworks
 - Upgrade also the underlying algorithms

London Underground example



London Underground in Prolog (1)

connected(bond_street,oxford_circus,central). connected(oxford_circus,tottenham_court_road,central). connected(bond_street,green_park,jubilee). connected(green_park,charing_cross,jubilee). connected(green_park,piccadilly_circus,piccadilly). connected(piccadilly_circus,leicester_square,piccadilly). connected(green_park,oxford_circus,victoria). connected(oxford_circus,piccadilly_circus,bakerloo). connected(piccadilly_circus,charing_cross,bakerloo). connected(tottenham_court_road,leicester_square,northern). connected(leicester_square,charing_cross,northern).

Symmetric facts now shown !!!

London Underground in Prolog (2)

Two stations are nearby if they are on the same line with at most one other station in between (symmetric facts not showr

nearby(bond_street,oxford_circus).
nearby(oxford_circus,tottenham_court_road).
nearby(bond_street,tottenham_court_road).
nearby(bond_street,green_park).
nearby(green_park,charing_cross).
nearby(bond_street,charing_cross).
nearby(green_park,piccadilly_circus).

or better

nearby(X,Y):-connected(X,Y,L).
nearby(X,Y):-connected(X,Z,L),connected(Z,Y,L).

Facts: unconditional truths Rules/Clauses: conditional truths

Both definitions are equivalent.



:- denotes implication

Recursion (2)

A station is reachable from another if they are on the same line, or with one, two, ... changes:

or better

```
reachable(X,Y):-connected(X,Y,L).
reachable(X,Y):-connected(X,Z,L),reachable(Z,Y).
```

Substitutions

- A *substitution* maps variables to terms:
 - {S->maria}
- A substitution can be *applied* to a clause:
 - likes(peter,maria):-student_of(maria,peter).
- The resulting clause is said to be an *instance* of the original clause, and a *ground instance* if it does not contain variables.
- Each instance of a clause is among its logical consequences.



?-reachable(oxford_circus, charing_cross, R).

- R = route(tottenham_court_road,route(leicester_square,noroute));
- R = route(piccadilly_circus, noroute);
- R = route(picadilly circus,route(leicester square,noroute))



Lists (3)

?-reachable(oxford_circus, charing_cross, R).

- R = [tottenham_court_road,leicester_square];
- R = [piccadilly_circus];
- R = [picadilly_circus,leicester_square]



Answering queries (1)

- Query: which station is nearby Tottenham Court Road?
 - ?- nearby(tottenham_court_road, W).
- Prefix ?- means it's a query and not a fact.
- Answer to query is:
 - {W -> leicester_square}
 - a so-called substitution.
- When nearby defined by facts, substitution found by unification.



Recall from AI course

- Unification to unify two different terms
- Resolution inference rule
- Refutation proofs, which derive the empty clause
- SLD-tree, which summarizes all possible proofs (left to right) for a goal



SLD-tree: one path for each proof-tree

The least Herbrand model

Definition:

- The set of all ground facts that are logically entailed by the program
- All ground facts not in the LHM are false ...
- LHM be computed as follows:
 - M0 = {}; M1 = { true }; i:=0
 - while Mi =\= Mi+1 do
 - i := i +1;
 - Mi := { h θ | h:- b1, ..., bn is clause and there is a substitution θ such that all bi θ ∈ Mi-1 }
 - Mi contains all true facts, all others are false

Example LHM

KB: p(a,b). a(X,Y) := p(X,Y). p(b,c). a(X,Y) := p(X,Z), a(Z,Y). M0 = emtpy; M1 = { true } M2 = { true, p(a,b), p(b,c) } M3 = M2 U {a(a,b), a(b,c) } M4 = M3 U { a(a,c) } M5 = M4

. . .

Stochastic Logic Programs

- Recall :
 - Prob. Regular Grammars
 - Prob. Context-Free Grammars
- What about Prob. Turing Machines ? Or Prob. Grammars ?
 - Stochastic logic programs combine probabilistic reasoning in the style of PCFGs with the expressive power of a programming language.

Recall PCFGs



We defined

- $P(tree|G) = \prod_i p_i^{c_i}$ where *i* ranges over all rules *i* used to derive tree, and c_i is the number of times they were applied
- $P(w_{1m}|G) = \sum_{j} P(tree_j|G)$ where j ranges over all possible parse trees for w_{1m} .
- Key concept : probabilities over derivations !!!

Stochastic Logic Programs

- Correspondence between CFG SLP
 - Symbols Predicates
 - Rules Clauses
 - Derivations SLD-derivations/Proofs
- **So**,
 - a stochastic logic program is an annotated logic program.
 - Each clause has an associated probability label. The sum of the probability labels for clauses defining a particular predicate is equal to 1.

An Example

```
1: \mathsf{card}(\mathsf{X},\mathsf{Y}) \gets \mathsf{rank}(\mathsf{X}), \mathsf{suit}(\mathsf{Y})
                                                      0.25 : suit(c) \leftarrow
0.25 : suit(d) \leftarrow
                                                      0.25 : suit(s) \leftarrow
0.25 : suit(h) \leftarrow
                                                      0.125 : rank(10) \leftarrow
0.125 : rank(a) \leftarrow
                                                      0.125 : rank(k) \leftarrow
0.125 : rank(7) \leftarrow
                                                      0.125 : rank(q) \leftarrow
0.125 : rank(8) \leftarrow
                                                      0.125 : rank(f) \leftarrow
0.125 : rank(9) \leftarrow
                                 :-card(a,s)
                                                                  Prob derivation
                                   :-rank(a), suit(s)
                                                                = 1.0.125.0.25
                                       :-suit(s)
```

Example

 $\begin{array}{l} 1: \mathsf{sentence}(\mathsf{A},\mathsf{B}) \leftarrow \mathsf{noun} - \mathsf{phrase}(\mathsf{C},\mathsf{A},\mathsf{D}), \mathsf{verb} - \mathsf{phrase}(\mathsf{C},\mathsf{D},\mathsf{B}).\\ 1: \mathsf{noun} - \mathsf{phrase}(\mathsf{A},\mathsf{B},\mathsf{C}) \leftarrow \mathsf{article}(\mathsf{A},\mathsf{B},\mathsf{D}), \mathsf{noun}(\mathsf{A},\mathsf{D},\mathsf{C}).\\ 1: \mathsf{verb} - \mathsf{phrase}(\mathsf{A},\mathsf{B},\mathsf{C}) \leftarrow \mathsf{intransitive} - \mathsf{verb}(\mathsf{A},\mathsf{B},\mathsf{C}).\\ 1/3: \mathsf{article}(\mathsf{singular},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{a},\mathsf{B}).\\ 1/3: \mathsf{article}(\mathsf{singular},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{the},\mathsf{B}).\\ 1/3: \mathsf{article}(\mathsf{plural},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{the},\mathsf{B}).\\ 1/2: \mathsf{noun}(\mathsf{singular},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{turtle},\mathsf{B}).\\ 1/2: \mathsf{noun}(\mathsf{plural},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{turtle},\mathsf{B}).\\ 1/2: \mathsf{intransitive} - \mathsf{verb}(\mathsf{singular},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{sleeps},\mathsf{B}).\\ 1/2: \mathsf{intransitive} - \mathsf{verb}(\mathsf{plural},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{sleep},\mathsf{B}).\\ 1/2: \mathsf{intransitive} - \mathsf{verb}(\mathsf{plural},\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{sleep},\mathsf{B}).\\ 1: \mathsf{terminal}([\mathsf{A}|\mathsf{B}],\mathsf{A},\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{sleep},\mathsf{B}).\\ 1: \mathsf{terminal}(\mathsf{A},\mathsf{A}|\mathsf{B}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{A}|\mathsf{A},\mathsf{A}) \leftarrow \mathsf{terminal}(\mathsf{A},\mathsf{A},\mathsf{A})).\\ 1: \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A})).\\ 1: \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{A})).\\ 1: \mathsf{A}(\mathsf{A},\mathsf{A}) \leftarrow \mathsf{A}(\mathsf{A},\mathsf{$

s([the,turtle,sleeps],[])?

SLPs : Key Ideas

- let us consider goals g (atoms, queries) of the form $p(X_1, ..., X_n)$ with the X_i different variables; corresponds to a non-terminal in a PCFG
- $P_D(der \ g|SLP) = \prod_i p_i^{c_i}$ where *i* ranges over all clauses *i* used in the derivation der *g* for goal *g* and c_i is the number of times *i* was applied; so far this is similar as for PCFGs
- key difference with PCFGs:
 - derivations in PCFGs always succeed,
 - proofs in SLPs can fail;
 - *refutations* are successful proofs
 - we are interested in the conditional probability of a derivation given that we know it is succesful / a refutation
- Therefore, define $P_R(ref \ g|SLP) = \frac{P_D(ref \ g)}{\sum_j P_D(ref_j \ g|SLP)}$, the probability P_R of refutation ref of g; where j ranges over all refutations ref_j of the goal g
- For ground atoms $g\theta$ (where θ is the substitution grounding g), define : $P_A(g\theta|SLP) = \sum_j P_R(ref_j \ g\theta|SLP)$ where ref_j ranges over all possible refutations for $g\theta$.

Example

• Cards :

- card(R,S) no proof with R in {a,7,8,9...} and S in { d,h,s,c} fails
- For each card, there is a unique refutation
- So,

 $P_A(card(r,s)) = P_D(derivation \ card(r,s))$ $= P_R(refutation \ card(r,s)) = 1/32$

Consider

same_suit(S,S) :suit(S), suit(S).

 In total 16 possible derivations, only 4 will succeed, so

 $P_D(derivation \ samesuit(S1, S2)) = 1/16$ $P_R(refutation \ samesuit(s, s)) = 1/4$ $P_A(samesuit(s, s)) = 1/4$

Another example (due to Cussens)

0.4:s(X) := p(X), p(X).0.3:p(a).0.2:q(a).0.6:s(X) := q(X).0.7:p(b).0.8:q(b).



Questions we can ask (and answer) about SLPs

- Compute the probability $P_A(a)$ of a ground atom a
- Find the most likely refutation r of a goal g, i.e. $argmax_{ref}P_R(ref g)$
- For a given SLP and set of atoms Atfor a goal g, compute the maximum likelihood parameters λ , i.e. $argmax_{\lambda}(\prod_{a \in At} P_A(a|SLP(\lambda)))$

Answers

- The algorithmic answers to these questions, again extend those of PCFGs and HMMs, in particular,
 - Tabling is used (to record probabilities of partial proofs and intermediate atoms)
 - Failure Adjusted EM (FAM) is used to solve parameter re-estimation problem
 - Additional hidden variables range over
 - Possible refutations and derivations for observed atoms
 - Topic of recent research
 - Freiburg : learning from refutations (instead of atoms), combined with structure learning

Sampling

- PRGs, PCFGs, and SLPs can also be used for sampling sentences, ground atoms that follow from the program
- Rather straightforward. Consider SLPs:
 - Probabilistically explore SLD-tree
 - At each step, select possible resolvents using the probability labels attached to clauses
 - If derivation succeeds, return corresponding (ground) atom
 - If derivation fails, then restart.

Bayesian Networks [Pearl 91]

Compact representation of joint probability distributions P(E,B,A,M,J)



Together:

Define a unique distribution in a compact, factored form

Quantitative part: Set of conditional probability distributions

P(E,B,A,M,J)=P(E) * P(B) * P(A|E,B) * P(M|A) * P(J|A)

Bayesian Networks [Pearl 91]



P(j) = P(j|a) * P(m|a) * P(a|e,b) * P(e) * P(b)+ P(j|a) * P(m|a) * P(a|e,b) * P(e) * P(b)...+ P(j|a) * P(m|a) * P(a|e,b) * P(e) * P(b)

Expressiveness Bayesian Nets

- A Bayesian net defines a probability distribution over a propositional logic
- Essentially, the possible states (worlds) are propositional interpretations
- But propositional logic is severely limited in expressive power, therefore consider combining BNs with logic programs
 - Bayesian logic programs
 - Actually, a BLP + some background knowledge generates a BN
 - So, BLP is a kind of BN template !!!

Bayesian Logic Programs (BLPs)



[Kersting, De Raedt] Bayesian Logic Programs (BLPs)



Bayesian Logic Programs (BLPs)



mc(Person) | mother(Mother, Person), pc(Mother), mc(Mother).
pc(Person) | father(Father, Person), pc(Father), mc(Father).
bt(Person) | pc(Person), mc(Person).





Bayesian logic programs

- Computing the ground BN (the BN that defines the semantics)
 - Compute the least Herbrand Model of the BLP
 - For each clause H | B1, ... BN with CPD
 - if there is a substitution θ such that {H θ, B1 θ, ..., BN θ} subset LHM, then H θ's parents include B1 θ, ..., BN θ, and with CPD specified by the clause
 - Delete logical atoms from BN (as their truth-value is known) - e.g. mother, father in the example
 - Possibly apply aggregation and combining rules
- For specific queries, only part of the resulting BN is necessary, the support net, cf. Next slides

Procedural Semantics

P(bt(ann))?



Procedural Semantics

Bayes' rule

P(bt(ann) | bt(fred)) =

P(bt(ann), bt(fred))

P(bt(ann), bt(fred))?

P(bt(fred))



Combining Rules



- Any algorithm which
 - has an empty output if and only if the input is empty
 - combines a set of CPDs into a single (combined) CPD
- E.g. noisy-or

Noisy-or: $P(A = true | V_1, ..., V_n) = 1 - \prod_{V_i = true} P(A = false | V_i = true)$ all causes can be independently inhibited

Combining Partial Knowledge



- variable # of parents for prepared/2 due to read/2
 - whether a student prepared a topic depends on the books she read
- CPD only for one book-topic pair

Summary BLPs



Joint probability distribution over the least
 Herbrand interpretation

Bayesian Logic Programs - Examples



Bayesian Logic Programs (BLPs) Unique probability distribution over Herbrand

- Unique probability distribution over Herbrand interpretations
 - Finite branching factor, finite proofs, no selfdependency
- Highlight
 - Separation of qualitative and quantitative parts
 - Functors
- Graphical Representation
- Discrete and continuous RV
- BNs, DBNs, HMMs, SCFGs, Prolog ...
- Turing-complete programming language
- Learning

Learning BLPs from Interpretations



Learning BLPs from Interpretations

Data case:

• Random Variable + States = (partial) Herbrand interpretation

Model(1) pc(brian)=b, bt(ann)=a, bt(brian)=?, bt(dorothy)=a

Bloodtype example



Model(2) bt(cecily)=ab, pc(henry)=a,

mc(fred)=?, bt(kim)=a,

pc(bob)=b

Parameter Estimation - BLPs



Parameter Estimation – BLPs

- Estimate the CPD θ entries that best fit the data
- "Best fit": ML parameters θ^*

 $θ^* = \operatorname{argmax}_{θ} P(\text{ data } | \text{ logic program, } θ)$ = argmax_θ log P(data | logic program, θ)

 Reduces to problem to estimate parameters of a Bayesian networks:

given structure,

partially observed random variables

Parameter Estimation – BLPs

		Background					
Model(1)		m(ann.dorothy).					
pc(brian)=b,		f(brian dorothy)					
bt(ann)=a.		ilv fred)					
bt(Model(2)				(med),		
bt	ot(cecily)=ab,			y	,trea),		
	bt(henry)=	а,		0	ob),		
	bt(fred)=?	,	<u>Mo</u>	o	lel(3)		
	bt(kim)=a,		pc(rex)=b,				
	bt(bob)=b		bt(d	d	oro)=a,		
			bt(ł	DI	rian)=?		
					/		

+





Parameter Estimation – BLPs

_		B	Background			
Model(1)		m(ann,dorothy),				
pc(brian)=b,		f(brian,dorothy),				
bt(bt(ann)=a.		ily,fred),			
bt(Model(2)		v.fred).			
bt(bt(cecily)=	=ab	o, pob)			
	bt(henry)=	a,				
	bt(fred)=?,		Model(3)			
	bt(kim)=a,		pc(rex)=b,			
	bt(bob)=b		bt(doro)=a,			
			bt(brian)=?			
			- + -			
			_			





