

Probabilistic Graphical Models

Course 2 : Markov Logic Networks

Hoel Le Capitaine

Academic year 2017-2018

Introduction

Background

Markov Logic

Inference

Learning

Introduction

Information and Knowledge – 1988

Databases

SQL

Datalog

Free text

Information retrieval

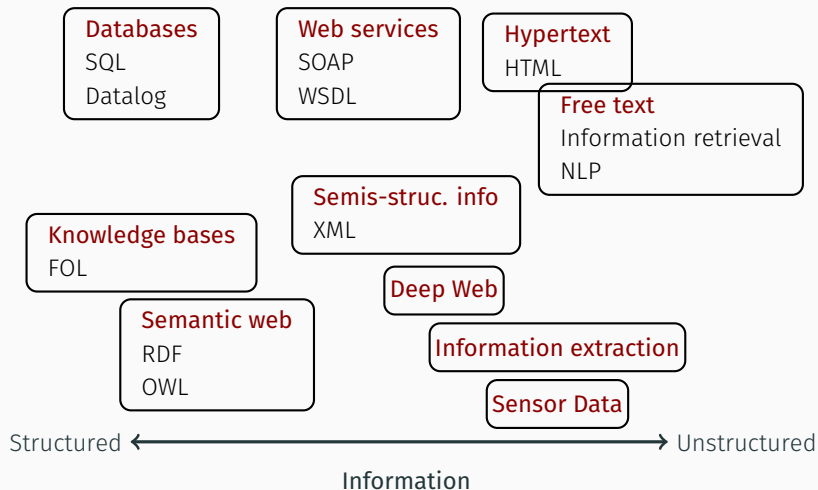
NLP

Knowledge bases

FOL

Structured ←————→ Unstructured
Information

Information and Knowledge – today



We need languages that can handle

- Structured information
- Unstructured information
- Any variation or combination of them

We need efficient algorithms for them

- Inference
- Machine learning

Background

Unifies first-order logic and probabilistic graphical models

- First-order logic handles structured information

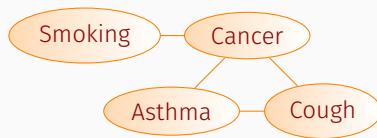
Unifies first-order logic and probabilistic graphical models

- First-order logic handles structured information
- Probability handles unstructured information

Unifies first-order logic and probabilistic graphical models

- First-order logic handles structured information
- Probability handles unstructured information
- No separation between the two
- Builds on previous work (BN, PRM, ...)

Undirected graphical models



Potential functions

- potential functions defined on **cliques**
- joint distribution of the MN :

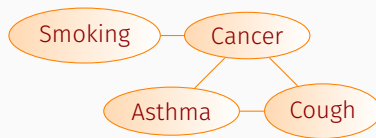
$$P(X = x) = \frac{1}{Z} \prod_k \phi_k(x_{\{k\}})$$

- $x_{\{k\}}$ is the state of the k th clique
- Z is the partition function

$$Z = \sum_x \prod_k \phi_k(x_{\{k\}})$$

Smoking	Cancer	$\phi(S, C)$
false	false	4.5
false	true	4.5
true	false	2.7
true	true	4.5

Undirected graphical models

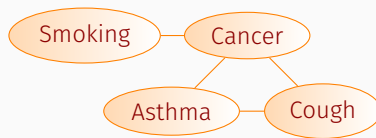


Log-linear models

- each clique replaced by an exponentiated weighted sum of (binary) features of the state

$$P(X = x) = \frac{1}{Z} \exp \left(\sum_j w_j \text{ (} f_j(x) \text{)} \right)$$

Undirected graphical models



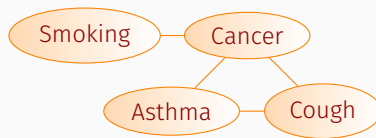
Log-linear models

- each clique replaced by an exponentiated weighted sum of (binary) features of the state

$$P(X = x) = \frac{1}{Z} \exp \left(\sum_j w_j \text{ } \textcircled{f_j(x)} \right)$$

- weight of feature j

Undirected graphical models



Log-linear models

- each clique replaced by an exponentiated weighted sum of (binary) features of the state

$$P(X = x) = \frac{1}{Z} \exp \left(\sum_j w_j \text{ } \text{ } f_j(x) \right)$$

- weight of feature j
- feature j

Some recalls

- constants, variables, functions, predicates
Sarah, x, MotherOf(x), Friends(x,y)
- grounding : replace all variables by constants
Friends(Sarah, Bernard)
- **World** (model, interpretation) : assignment of truth valuser to all ground predicates

Markov Logic

From logic to probability

- a logical knowledge base is a set of **hard constraints** on the set of possible worlds
- transform them into **soft constraints** : if a world violates a formula, it becomes less probable, not impossible
- each formula is associated to a **weight**, corresponding to the strength of the constraint

$$P(world) \approx \exp \left(\sum \text{weights of satisfied formulas} \right)$$

Definition

A **Markov Logic Network (MLN)** is a set of pairs (F, w) where F is a 1st-order logic formula, and w is a real value.

It defines a Markov network with

- one node for each grounding of each predicate
- one feature for each grounding of each formula F , with corresponding weight w

A first KB

- smoking causes cancer

A first KB

- smoking causes cancer $\forall x \text{ Smokes}(x) \rightarrow \text{Cancer}(x)$

Example : friends and smokers

A first KB

- smoking causes cancer $\forall x \text{ Smokes}(x) \rightarrow \text{Cancer}(x)$
- if 2 people are friends, either both smokes or neither does

Example : friends and smokers

A first KB

- smoking causes cancer $\forall x \text{ Smokes}(x) \rightarrow \text{Cancer}(x)$
- if 2 people are friends, either both smokes or neither does
 $\forall x \forall y \text{ Friends}(x, y) \rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y))$

Example : friends and smokers

A first KB

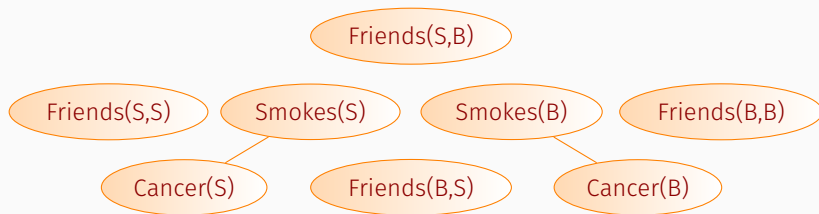
- **smoking causes cancer**

$$\forall x \text{ Smokes}(x) \rightarrow \text{Cancer}(x) \quad w = 1.5$$

- if 2 people are friends, either both smokes or neither does

$$\forall x \forall y \text{ Friends}(x, y) \rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y)) \quad w = 1.1$$

Two constants : Sarah (S) and Bernard (B)



Example : friends and smokers

A first KB

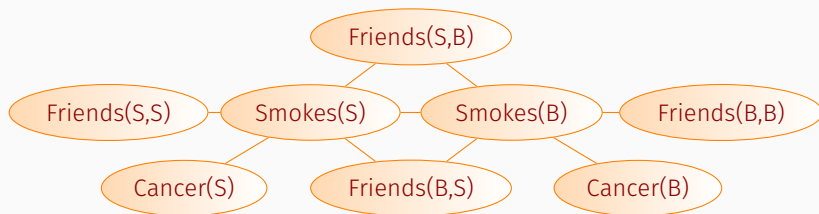
- smoking causes cancer

$$\forall x \text{ Smokes}(x) \rightarrow \text{Cancer}(x) \quad w = 1.5$$

- **if 2 people are friends, either both smokes or neither does**

$$\forall x \forall y \text{ Friends}(x, y) \rightarrow (\text{Smokes}(x) \Leftrightarrow \text{Smokes}(y)) \quad w = 1.1$$

Two constants : Sarah (S) and Bernard (B)



- MLN is a template for ground Markov networks
- probability of a world X

$$P(X) = \frac{1}{Z} \exp \left(\sum_i w_i n_i(X) \right)$$

$n_i(X)$: number of true groundings of formula i in X

- **typed** variables and constants reduce size of the net
- functions, existential quantifiers
- both infinite and continuous domains

- infinite weights exactly correspond to first-order logic
- satisfiable KB, positive weights : Satisfying assignments = Modes of distribution
- Markov logic allows contradictions between formulas

Inference

- we want to find the most likely state of world, given our observation
 $\text{argmax}_y P(y|x)$

- we want to find the most likely state of world, given our observation

$$\text{argmax}_y P(y|x)$$

- y : query
- x : evidence

- we want to find the most likely state of world, given our observation
 $\operatorname{argmax}_y P(y|x)$
 - y : query
 - x : evidence
- using definition, **$\operatorname{argmax}_y \frac{1}{Z_x} \exp(\sum_i w_i n_i(x, y))$**

- we want to find the most likely state of world, given our observation
 $\operatorname{argmax}_y P(y|x)$
 - y : query
 - x : evidence
- using definition, **$\operatorname{argmax}_y \frac{1}{Z_x} \exp(\sum_i w_i n_i(x, y))$**
- due to indep. and monotonicity, **$\operatorname{argmax}_y \sum_i w_i n_i(x, y)$**
- reduced to a weighted maxSAT problem (eg walkSAT)
- may be faster than logical inference !

```
for i  $\leftarrow$  1 to max-tries do
  solution = random truth assignment
  for j  $\leftarrow$  1 to max-flips do
    if all clauses satisfied then
      return solution
    c  $\leftarrow$  random unsatisfied clause
    with probability p
      flip a random variable in c
    else
      flip variable in c that maximizes number of satisfied clauses
  return failure
```

```
for i ← 1 to max-tries do
  solution = random truth assignment
  for j ← 1 to max-flips do
    if  $\sum \text{weights of satisfied clauses} > t$  then
      return solution
    c ← random unsatisfied clause
    with probability p
      flip a random variable in c
    else
      flip variable in c that maximizes  $\sum \text{weights of satisfied clauses}$ 
  return failure, best solution found
```


- memory explosion ...
- if there are n constants, and highest clause arity is c , ground network requires $O(n^c)$ memory
- we need to exploit sparseness : LazySAT (2006)

Learning

- Data is a relational database
- Closed world assumption (if not: EM)
- Learning parameters (weights)
 - Generatively
 - Discriminatively
- Learning structure (formulas)

- based on maximum likelihood
- uses gradient ascent (or more complicated ones)
- no local maxima $\Delta_{w_i} \log P_w(x) = n_i(x) - E_w[n_i(x)]$
 - $E_w[n_i(x)]$ expected number of true groundings according to the model
- but requires inference at each step of the optimization (slow)

- Likelihood of each variable given its neighbors in the data

$$PL(x) = \prod_i P(x_i | N(x_i))$$

- Does not require inference at each step
- Consistent estimator
- Widely used in vision, spatial statistics, etc.
- But PL parameters may not work well for long inference chains

- maximize conditional likelihood of query y given observation x

$$\Delta_{w_i} \log P_w(y|x) = n_i(x, y) - E_w[n_i(x, y)]$$

- approximation of expectation counts by counts in MAP state of y given x (voted perceptron)

- Generalizes feature induction in Markov nets
- Any inductive logic programming approach can be used, but . . .
- Goal is to induce any clauses, not just Horn
- Evaluation function should be likelihood
- Requires learning weights for each candidate
- Turns out not to be bottleneck
- Bottleneck is counting clause groundings
- Solution: Subsampling

- **Initial state** : Unit clauses or hand-coded KB
- **Operators** : Add/remove literal, flip sign
- **Evaluation function** : Pseudo-likelihood + Structure prior
- **Search** : Beam, Shortest-first, Bottom-up (2005–2007)