

Interactive HMM construction based on interesting sequences

Szymon Jaroszewicz

National Institute of Telecommunications
Warsaw, Poland

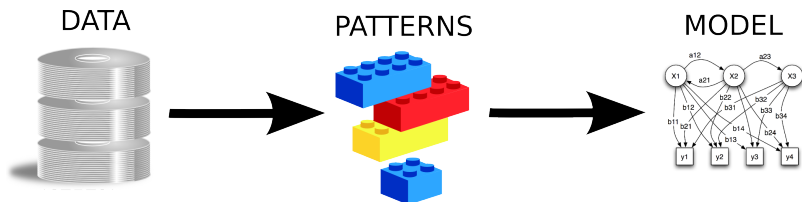
LeGo 2008

- Building models interactively based on interesting patterns
- Hidden Markov Models
- Interesting patterns w.r.t. Hidden Markov Models
- Experimental evaluation: web server log
- Conclusions and Future research

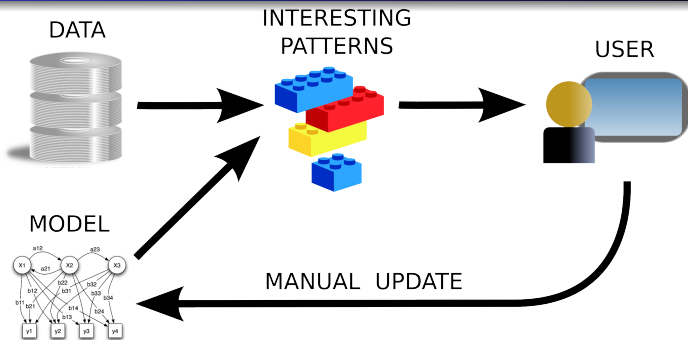
Typical approach: Automatic model construction



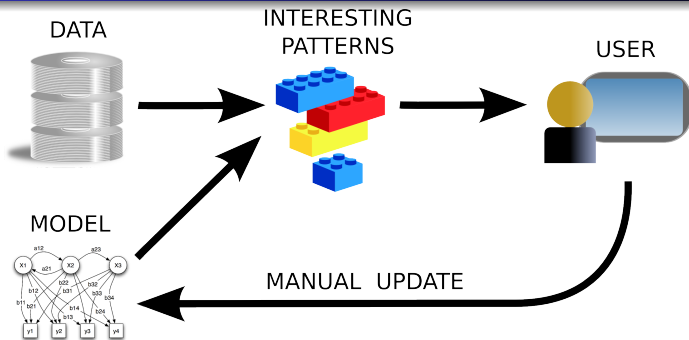
Or:



Here: Interactive model construction



Here: Interactive model construction



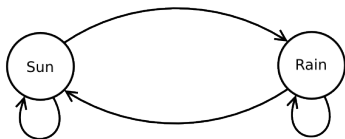
- + Understandable models
- + Learn while building models
- Have to do 'manual' work :(

Scalable pattern mining with Bayesian networks as background knowledge

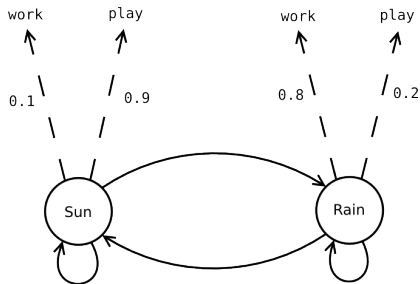
S. Jaroszewicz, T. Scheffer, D. Simovici
KDD'04, KDD'05, DMKD (to appear)

- Bayesian networks used as background model
- Exact and approximate algorithms given
- Models much closer to real relationships than automatically built models

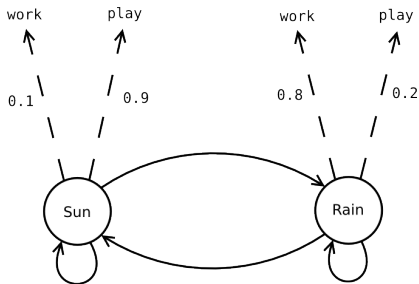
Hidden Markov Models (HMMs)



Hidden Markov Models (HMMs)



Hidden Markov Models (HMMs)



User gives the structure of the HMM:

- internal states
- which transitions are possible (not probabilities)
- which emission symbols are possible for each state (not probabilities)

Interestingness of sequences w.r.t. an HMM

$$\text{Inter}(\text{seq}) = \left| \text{Prob}^{\text{HMM}}\{\text{seq}\} - \text{Prob}^{\text{Data}}\{\text{seq}\} \right|$$

Algorithm for finding all ε -interesting sequences

- 1 Train *HMM* parameters based on *Data* (Baum-Welch)
- 2 Find all *seq* such that $\text{Prob}^{Data}\{seq\} > \varepsilon$
- 3 Find all *seq* such that $\text{Prob}^{HMM}\{seq\} > \varepsilon$
- 4 Compute Prob^{Data} for *seq* frequent in *HMM* but not in *Data*
- 5 Compute Prob^{HMM} for *seq* frequent in *Data* but not in *HMM*
- 6 Compute $Inter(seq)$ for all sequences
- 7 Output ε -interesting sequences

Inference in Hidden Markov Models

- Probability that sequence seq (starting at $t = 0$) is emitted and HMM ends in state s_i

$$\alpha(seq, s_i)$$

- Efficient recursive updating:

$$\alpha(seq + o^{n+1}, s_i) = \sum_j \alpha(seq, s_j) \mathbf{P}_{ji} \mathbf{E}_{io^{n+1}}$$

- $\text{Prob}^{HMM}\{seq\} = \sum_i \alpha(seq, s_i)$

Finding frequent sequences in Hidden Markov Models

- Monotonicity property holds

$$\text{Prob}^{HMM}\{seq + o\} \leq \text{Prob}^{HMM}\{seq\}$$

- Standard depth-first frequent pattern mining works
alpha probabilities used instead of support counting
- Very efficient: probability updating is fast

Web log format:

```
195.205.118.10 [01/Jan/2007:00:04:33 +0100] "GET  
/journal/paper_1.pdf" 200 8833 "http://www.google.pl/"
```

```
65.55.208.68 [01/Jan/2007:00:04:45] "GET /robots.txt" 200  
51 "-" "msnbot/1.0"
```

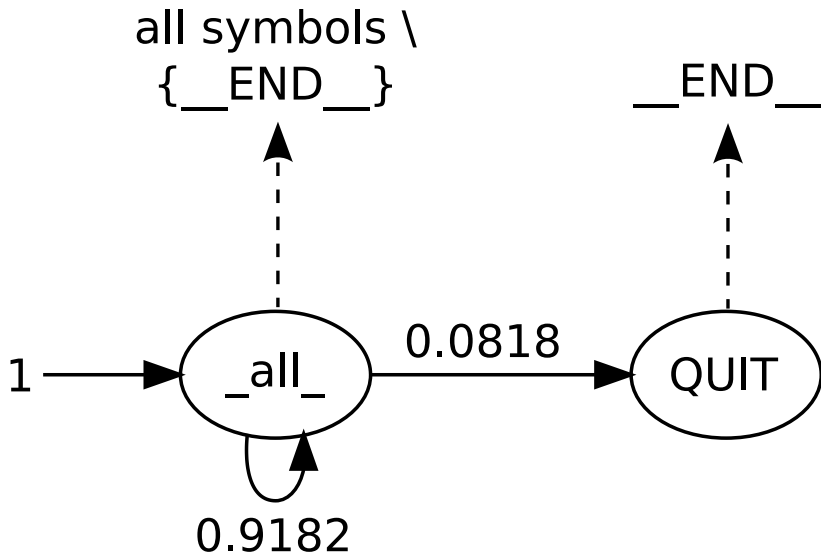
Preprocessing:

- keep only top level directory
- sessionizing

Result: sessions:

```
journal/, journal/, __END__  
robots.txt, index.html, journal/, ..., __END__  
exchweb/, exchange/, exchange/, ..., __END__  
...
```

Initial HMM

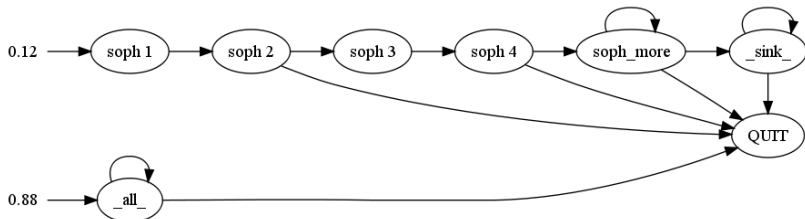


Top sequences:

- sophos/,sophos/
 $\text{Prob}^{HMM} = 1.17\%$
 $\text{Prob}^{Data} = 11.48\%$
- sophos/,sophos/,sophos/,sophos/
 $\text{Prob}^{HMM} = 0.013\%$
 $\text{Prob}^{Data} = 9.29\%$
- Update of the Sophos antivirus
- **Always** accessed 2, 4 or more times

The Sophos antivirus: update to the model

- The new model is:



- Each `soph` state only emits the `sophos/` symbol
- `sophos/` symbol removed from `_all_` state

- Sequence: journals/, journals/, `favicon.ico`

$$\text{Prob}^{HMM} \approx 0$$

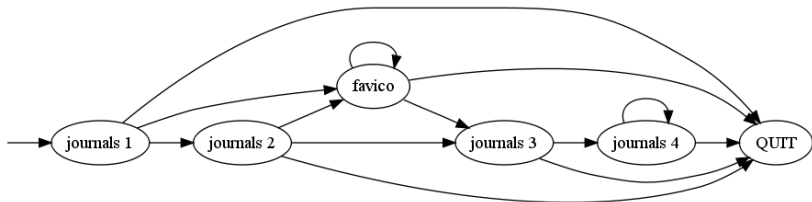
$$\text{Prob}^{Data} \approx 2\%$$

- `favicon.ico` small icon next to web address



- Default location: main directory
At the Institute: `img/` directory
- HTML header contains the other location; PDF can't
- Browser tries the default location and fails
- Fixed: icon appears now

- Added the following segment to the model:



- The same PDF file often accessed twice; unable to explain:
 - accelerators?
 - browser errors?
 - server errors?

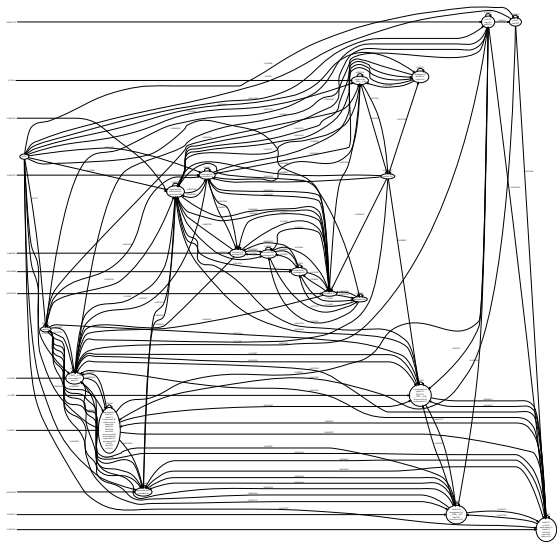
- Exchange mail web reader
- robots: Google / MSN / Yahoo
- RSS readers
- ...

- Quickly built a model of high level user behavior
- **Accuracy:** probability of all sequences modeled with error < 0.01
Every sequence is either:
 - uninteresting (modeled well)
 - infrequent
- **Understandability:** the model is easily understandable
- **Learnt** a lot about the data while modeling

Comparison with automatically learned models

- 20 hidden states + Baum Welch algorithm
- only transitions with prob. > 0.01
- all transitions with prob. > 0.001

All transitions with prob. > 0.001



Conclusions and Future work

Conclusions:

- Interactive model construction based on interesting patterns = **Understandability** + **Accuracy** + **Learning** about the data

Future work:

- Patterns starting at arbitrary time
- More general models: Dynamic Bayesian Networks, models of biological systems
- Automatic model updating (?)