Web-Mining Agents Probabilistic Reasoning over Sequential Structures

Prof. Dr. Ralf Möller Universität zu Lübeck Institut für Informationssysteme

Tanya Braun (Übungen)



IM FOCUS DAS LEBEN

Literature



Stuart Russell • Peter Norvig Prentice Hall Series in Artificial Intelligence

http://aima.cs.berkeley.edu

Structures (ordinal, interval, ratio scale)

- Temporal structures
- Sentence structures
- ...
- Genomic structures
- ...



Temporal Probabilistic Agent



Probabilistic Temporal Models

- Dynamic Bayesian Networks (DBNs)
- Hidden Markov Models (HMMs)
- Kalman Filters



Time and Uncertainty

- The world changes, we need to track and predict it
- Examples: diabetes management, traffic monitoring
- Basic idea: copy state and evidence variables for each time step
- **X**_t set of unobservable state variables at time t
 - e.g., BloodSugar_t, StomachContents_t
- **E**_t set of evidence variables at time t
 - e.g., MeasuredBloodSugar_t, PulseRate_t, FoodEaten_t
- Assumes discrete time steps



States and Observations

- Process of change viewed as series of snapshots, each describing the state of the world at a particular time
- Time slice involves a set of random variables indexed by t:
 - the set of unobservable state variables X_t
 - the set of observable evidence variable **E**_t
- The observation at time t is $\mathbf{E}_t = \mathbf{e}_t$ for some set of values \mathbf{e}_t
- The notation **X**_{a:b} denotes the set of variables from **X**_a to **X**_b



Dynamic Bayesian Networks

- How can we model dynamic situations with a Bayesian network?
- Example: Is it raining today?

$$X_t = \{R_t\}$$
$$E_t = \{U_t\}$$

 \Rightarrow next step: specify dependencies among the variables.

The term "dynamic" means we are modeling a dynamic system, not that the network structure changes over time.



DBN - Representation

 $P(U_t | Parent(U_t))$

- Problem:
 - 1. Necessity to specify an unbounded number of conditional probability tables, one for each variable in each slice,
 - 2. Each one might involve an unbounded number of parents.
- Solution:
 - 1. Assume that changes in the world state are caused by a stationary process (unmoving process over time).

is the same for all t



Generalization of DBNs

- Time is just one sequential structures
- Can generalize to any dynamic structure "expansion"
 - Sentences
 - Spatial structures
 - ...



DBN - Representation

- Solution cont.:
 - 2. Use **Markov assumption** The current state depends on only in a finite history of previous states.

Using the first-order Markov process:

$$P(X_t / X_{0:t-1}) = P(X_t / X_{t-1})$$
Transition
Model

In addition to restricting the parents of the state variable *Xt*, we must restrict the parents of the evidence variable *Et*

$$P(E_t / X_{0:t}, E_{0:t-1}) = P(E_t / X_t)$$
 Sensor
Model



Stationary Process/Markov Assumption

- Markov Assumption: X_t depends on some previous X_is
- First-order Markov process:

 $P(X_t | X_{0:t-1}) = P(X_t | X_{t-1})$

- kth order: depends on previous k time steps
- Sensor Markov assumption:

 $P(E_t|X_{0:t}, E_{0:t-1}) = P(E_t|X_t)$

- Assume stationary process: transition model:
 - $P(X_t|X_{t-1})$ and sensor model $P(E_t|X_t)$ are the same for all t
 - Changes in the world state governed by laws not changing over time



Dynamic Bayesian Networks

- There are two possible fixes if the approximation is too inaccurate:
 - Increasing the order of the Markov process model. For example, adding $Rain_{t-2}$ as a parent of $Rain_t$, which might give slightly more accurate predictions.
 - Increasing the set of state variables. For example, adding $Season_t$ to allow to incorporate historical records of rainy seasons, or adding $Temperature_t$, $Humidity_t$ and $Pressure_t$ to allow to use a physical model of rainy conditions.



Dynamic Bayesian Network



Bayesian network structure corresponding to a first-order of Markov process with state defined by the variables Xt.



A second order of Markov process



Complete Joint Distribution

- Given:
 - Transition model: $P(X_t|X_{t-1})$
 - Sensor model: $P(E_t|X_t)$
 - Prior probability: $P(X_0)$
- Then we can specify complete joint distribution:

$$P(X_0, X_1, ..., X_t, E_1, ..., E_t) = P(X_0) \prod_{i=1}^t P(X_i | X_{i-1}) P(E_i | X_i)$$



Example





Inference Tasks

- **Filtering:** What is the probability that it is raining today, given all the umbrella observations up through today?
- **Prediction:** What is the probability that it will rain the day after tomorrow, given all the umbrella observations up through today?
- **Smoothing:** What is the probability that it rained yesterday, given all the umbrella observations through today?
- Most likely explanation / most probable explanation: if the umbrella appeared the first three days but not on the fourth, what is the most likely weather sequence to produce these umbrella sightings?



• Filtering or Monitoring:

Compute the belief state - the posterior distribution over the *current* state, given all evidence to date.

 $P(X_t / e_{1:t})$

Filtering is what a rational agent needs to do in order to keep track of the current state so that the rational decisions can be made.



• Filtering cont.

Given the results of filtering up to time *t*, one can easily compute the result for t+1 from the new evidence e_{t+1}

$$P(X_{t+1} / e_{1:t+1}) = f(e_{t+1,} P(X_t / e_{1:t+1}))$$

$$= P(X_{t+1} / e_{1:t,} e_{t+1})$$

$$= \alpha P(e_{t+1} / X_{t+1,} e_{1:t}) P(X_{t+1} / e_{1:t})$$

$$= \alpha P(e_{t+1} / X_{t+1}) P(X_{t+1} / e_{1:t})$$

$$= \alpha P(e_{t+1} / X_{t+1}) P(X_{t+1} / e_{1:t})$$

$$(for some function f)$$

$$(dividing up the evidence)$$

$$(using Bayes' Theorem)$$

$$(by the Markov property$$

$$of evidence)$$

 α is a normalizing constant used to make probabilities sum up to 1.



• Filtering cont.

The second term $P(X_{t+1} / e_{1:t})$ represents a one-step prediction of the next step, and the first term $P(e_{t+1} / X_{t+1})$ updates this with the new evidence.

Now we obtain the one-step prediction for the next step by conditioning on the current state Xt:

$$P(X_{t+1} / e_{1:t+1}) = \alpha P(e_{t+1} / X_{t+1}) \sum_{X_t} P(X_{t+1} / x_t, e_{1:t}) P(x_t / e_{1:t})$$

$$= \alpha P(e_{t+1} / X_{t+1}) \sum_{X_t} P(X_{t+1} / X_t) P(x_t / e_{1:t})$$
(using the Markov property)



$f_{1:t+1} = FORWARD(f_{1:t}, e_{t+1}) \text{ where } f_{1:t} = P(\mathbf{X}_t | e_{1:t})$ Time and space **constant** (independent of t)



Example





Illustration for two steps in the Umbrella example:

• On day 1, the umbrella appears so U1=true. The prediction from t=0 to t=1 is

$$P(R_1) = \sum_{r_0} P(R_1 / r_0) P(r_0)$$

and updating it with the evidence for t=1 gives

$$P(R_1 / u_1) = \alpha P(u_1 / R_1) P(R_1)$$

• On day 2, the umbrella appears so U2=true. The prediction from t=1 to t=2 is

$$P(R_2 / u_1) = \sum_{r_1} P(R_2 / r_1) P(r_1 / u_1)$$

and updating it with the evidence for t=2 gives

$$P(R_2 / u_1, u_2) = \alpha P(u_2 / R_2) P(R_2 / u_1)$$

UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME

Example cntd.





• Prediction:

Compute the posterior distribution over the *future* state, given all evidence to date.

$$P(X_{t+k} / e_{1:t})$$
 for some

The task of prediction can be seen simply as filtering without the addition of new evidence.



k > 0

• Smoothing or hindsight:

Compute the posterior distribution over the *past* state, given all evidence up to the present.

$$P(X_k / e_{1:t}) \quad f^{c}$$

or some k such that $0 \le k < t$.

Hindsight provides a better estimate of the state than was available at the time, because it incorporates more evidence.



Smoothing

Divide evidence $\mathbf{e}_{1:t}$ into $\mathbf{e}_{1:k}$, $\mathbf{e}_{k+1:t}$:

$$\mathbf{P}(\mathbf{X}_{k}|\mathbf{e}_{1:t}) = \mathbf{P}(\mathbf{X}_{k}|\mathbf{e}_{1:k}, \mathbf{e}_{k+1:t})$$

= $\alpha \mathbf{P}(\mathbf{X}_{k}|\mathbf{e}_{1:k})\mathbf{P}(\mathbf{e}_{k+1:t}|\mathbf{X}_{k}, \mathbf{e}_{1:k})$
= $\alpha \mathbf{P}(\mathbf{X}_{k}|\mathbf{e}_{1:k})\mathbf{P}(\mathbf{e}_{k+1:t}|\mathbf{X}_{k})$
= $\alpha \mathbf{f}_{1:k}\mathbf{b}_{k+1:t}$

Backward message computed by a backwards recursion:

$$\begin{aligned} \mathbf{P}(\mathbf{e}_{k+1:t}|\mathbf{X}_k) &= \sum_{\mathbf{x}_{k+1}} \mathbf{P}(\mathbf{e}_{k+1:t}|\mathbf{X}_k, \mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1}|\mathbf{X}_k) \\ &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1:t}|\mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1}|\mathbf{X}_k) \\ &= \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1}|\mathbf{x}_{k+1}) P(\mathbf{e}_{k+2:t}|\mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1}|\mathbf{X}_k) \end{aligned}$$

Forward-backward algorithm: cache forward messages along the way Time linear in t (polytree inference), space $O(t|\mathbf{f}|)$



Example contd.





• Most likely explanation:

Compute the sequence of states that is most likely to have generated a given sequence of observation.

$$\arg \max_{x_{1:t}} P(X_{1:t} | e_{1:t})$$

Algorithms for this task are useful in many applications, including, e.g., speech recognition.



Web-Mining Agents Probabilistic Reasoning over Time

Prof. Dr. Ralf Möller Universität zu Lübeck Institut für Informationssysteme

Tanya Braun (Übungen)



IM FOCUS DAS LEBEN

Most-likely explanation

Most likely sequence \neq sequence of most likely states!!!!

Most likely path to each \mathbf{x}_{t+1}

= most likely path to some \mathbf{x}_t plus one more step

 $\max_{\mathbf{x}_{1}...\mathbf{x}_{t}} \mathbf{P}(\mathbf{x}_{1},\ldots,\mathbf{x}_{t},\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1}) = \mathbf{P}(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_{t}} \left(\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_{t}) \max_{\mathbf{x}_{1}...\mathbf{x}_{t-1}} P(\mathbf{x}_{1},\ldots,\mathbf{x}_{t-1},\mathbf{x}_{t}|\mathbf{e}_{1:t}) \right)$

Identical to filtering, except $\mathbf{f}_{1:t}$ replaced by

$$\mathbf{m}_{1:t} = \max_{\mathbf{x}_1 \dots \mathbf{x}_{t-1}} \mathbf{P}(\mathbf{x}_1, \dots, \mathbf{x}_{t-1}, \mathbf{X}_t | \mathbf{e}_{1:t}),$$

I.e., $\mathbf{m}_{1:t}(i)$ gives the probability of the most likely path to state i. Update has sum replaced by max, giving the Viterbi algorithm:

 $\mathbf{m}_{1:t+1} = \mathbf{P}(\mathbf{e}_{t+1} | \mathbf{X}_{t+1}) \max_{\mathbf{X}_t} \left(\mathbf{P}(\mathbf{X}_{t+1} | \mathbf{x}_t) \mathbf{m}_{1:t} \right)$



Rain/Umbrella Example





The occasionally dishonest casino

- A casino uses a fair die most of the time, but occasionally switches to a loaded one
 - Fair die: Prob(1) = Prob(2) = ... = Prob(6) = 1/6
 - Loaded die: Prob(1) = Prob(2) = ... = Prob(5) = 1/10, $Prob(6) = \frac{1}{2}$
 - These are the *emission* probabilities
- Transition probabilities
 - Prob(Fair \rightarrow Loaded) = 0.01
 - Prob(Loaded \rightarrow Fair) = 0.2
 - Transitions between states modeled by a Markov process



Transition model for the casino





The occasionally dishonest casino

- Known:
 - The structure of the model
 - The transition probabilities
- Hidden: What the casino did
 - FFFFFLLLLLLFFFF...
- Observable: The series of die tosses
 - 3415256664666153...
- What we must infer:
 - When was a fair die used?
 - When was a loaded one used?
 - The answer is a sequence FFFFFFFLLLLLFFF...



Making the inference

- Model assigns a probability to each explanation of the observation: P(326|FFL) $= P(3|F) \cdot P(F \rightarrow F) \cdot P(2|F) \cdot P(F \rightarrow L) \cdot P(6|L)$
 - $= 1/6 \cdot 0.99 \cdot 1/6 \cdot 0.01 \cdot \frac{1}{2}$
- **Maximum Likelihood:** Determine which explanation is most likely
 - Find the path *most likely* to have produced the observed sequence
- **Total probability:** Determine probability that observed sequence was produced by the model
 - Consider *all* paths that could have produced the observed sequence



Notation

- x is the sequence of symbols emitted by model
 - x_i is the symbol emitted at time *i*
- A *path*, π , is a sequence of states
 - The *i*-th state in π is π_i
- *a_{kr}* is the probability of making a transition from state *k* to state *r*:

$$a_{kr} = \Pr(\pi_i = r \mid \pi_{i-1} = k)$$

e_k(b) is the probability that symbol *b* is emitted when in state *k*

$$e_k(b) = \Pr(x_i = b \mid \pi_i = k)$$


A "parse" of a sequence





The occasionally dishonest casino

$$x = \langle x_1, x_2, x_3 \rangle = \langle 6, 2, 6 \rangle$$

$$Pr(x, \pi^{(1)}) = a_{0F}e_F(6)a_{FF}e_F(2)a_{FF}e_F(6)$$

$$= 0.5 \times \frac{1}{6} \times 0.99 \times \frac{1}{6} \times 0.99 \times \frac{1}{6}$$

$$\approx 0.00227$$

$$Pr(x, \pi^{(2)}) = a_{0L}e_L(6)a_{LL}e_L(2)a_{LL}e_L(6)$$

= 0.5 × 0.5 × 0.8 × 0.1 × 0.8 × 0.5
= 0.008

$$\pi^{(3)} = LFL$$

 $\pi^{(2)} = LLL$

$$Pr(x, \pi^{(3)}) = a_{0L}e_L(6)a_{LF}e_F(2)a_{FL}e_L(6)a_{L0}$$
$$= 0.5 \times 0.5 \times 0.2 \times \frac{1}{6} \times 0.01 \times 0.5$$
$$\approx 0.0000417$$



The most probable path

The most likely path
$$\pi^*$$
 satisfies
 $\pi^* = \arg \max_{\pi} \Pr(x, \pi)$
To find π^* , consider all possible ways the last symbol
of x could have been emitted
Let
 $v_k(i) = \operatorname{Prob. of path} \langle \pi_1, \dots, \pi_i \rangle \text{ most likely}$
to emit $\langle x_1, \dots, x_i \rangle$ such that $\pi_i = k$
Then
 $v_k(i) = e_k(x_i) \max_r (v_r(i-1)a_{rk})$



The Viterbi Algorithm

• Initialization (i = 0)

$$v_0(0) = 1$$
, $v_k(0) = 0$ for $k > 0$

• Recursion (i = 1, ..., L): For each state k

$$v_k(i) = e_k(x_i) \max_r \left(v_r(i-1)a_{rk} \right)$$

• Termination:

$$\Pr(x,\pi^*) = \max_k \left(v_k(Length)a_{k0} \right)$$

To find π^* , use trace-back, as in dynamic programming



Viterbi: Example





Viterbi gets it right more often than not

Rolls	315116246446644245321131631164152133625144543631656626566666
Die	FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi	FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Rolls	6511664531326512456366646316366631623264552352666666625151631
Die	LLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi	LLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Rolls	222555441666566563564324364131513465146353411126414626253356
Die	FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi	FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Rolls	366163666466232534413661661163252562462255265252266435353336
Die	LLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi	LLLLLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Rolls	23312162536441443233516324363366556246666626326666612355245242
Die	FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi	FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF



Dynamic Bayesian Networks

- *Learning* requires the full smoothing inference, rather than filtering, because it provides better estimates of the state of the process.
- Learning the parameters of a BN is done using Expectation – Maximization (EM) Algorithms. Iterative optimization method to estimate some unknowns parameters.



Acknowledgements:

The following slides are taken from CS 388 Natural Language Processing: Part-Of-Speech Tagging, Sequence Labeling, and Hidden Markov Models (HMMs)

> **Raymond J. Mooney** University of Texas at Austin



Application 1: Part Of Speech Tagging

- Annotate each word in a sentence with a part-of-speech marker.
- Lowest level of syntactic analysis.

John saw the saw and decided to take it to the table. NNP VBD DT NN CC VBD TO VB PRP IN DT NN

• Useful for subsequent syntactic parsing and word sense disambiguation.



English Parts of Speech

- Noun (person, place or thing)
 - Singular (NN): dog, fork
 - Plural (NNS): dogs, forks
 - Proper (NNP, NNPS): John, Springfields
 - Personal pronoun (PRP): I, you, he, she, it
 - Wh-pronoun (WP): who, what
- Verb (actions and processes)
 - Base, infinitive (VB): eat
 - Past tense (VBD): ate
 - Gerund (VBG): eating
 - Past participle (VBN): eaten
 - Non 3rd person singular present tense (VBP): eat
 - 3rd person singular present tense: (VBZ): eats
 - Modal (MD): should, can
 - To (TO): to (to eat)



English Parts of Speech (cont.)

- Adjective (modify nouns)
 - Basic (JJ): red, tall
 - Comparative (JJR): redder, taller
 - Superlative (JJS): reddest, tallest
- Adverb (modify verbs)
 - Basic (RB): quickly
 - Comparative (RBR): quicker
 - Superlative (RBS): quickest
- Preposition (IN): on, in, by, to, with
- Determiner:
 - Basic (DT) a, an, the
 - WH-determiner (WDT): which, that
- Coordinating Conjunction (CC): and, but, or,
- Particle (RP): off (took off), up (put up)



In General: Sequence Labeling Problem

- Many NLP problems can viewed as sequence labeling.
- Each token in a sequence is assigned a label.
- Labels of tokens are dependent on the labels of other tokens in the sequence, particularly their neighbors (not i.i.d).





Application 2: Information Extraction

- Identify phrases in language that refer to specific types of entities and relations in text.
- Named entity recognition is the task of identifying names of people, places, organizations, etc. in text.
 people organizations places
 - Michael Dell is the CEO of Dell Computer Corporation and lives in Austin Texas.
- Extract pieces of information relevant to a specific application, e.g. used car ads:

make model year mileage price

For sale, 2002 Toyota Prius, 20,000 mi, \$15K or best offer.
 Available starting July 30, 2006.



Semantic Role Labeling

- For each clause, determine the semantic role played by each noun phrase that is an argument to the verb.
 agent patient source destination instrument
 - John drove Mary from Austin to Dallas in his Toyota Prius.
 - The hammer broke the window.
- Also referred to a "case role analysis," "thematic analysis," and "shallow semantic parsing"



Application 3: Bioinformatics

• Sequence labeling also valuable in labeling genetic sequences in genome analysis.

extron intron

– AGCTAACGTTCGATACGGATTACAGCCT



Back to App1: Sequence Labeling as Classification

























 Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to the table. classifier TO















 Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it





 Classify each token independently but use as input features, information about the surrounding tokens (sliding window).

John saw the saw and decided to take it to





Using Outputs as Inputs

- Better input features are usually the categories of the surrounding tokens, but these are not available yet.
- Can use category of either the preceding or succeeding tokens by going forward or back and using previous output.


















































Disambiguating "to" in this case would be even easier backward.





NN

Disambiguating "to" in this case would be even easier backward.

John saw the saw and decided to take it to the













































DBN – Special Cases

• Hidden Markov Model (HMMs):

Temporal probabilistic model in which the state of the process is described by a single discrete random variable. (The simplest kind of DBN)

• Kalman Filter Models (KFMs):

Estimate the state of a physical system from noisy observations over time. Also known as linear dynamical systems (LDSs).



DBN – Basic Inference

• Filtering

$$P(X_{t+1} / e_{1:t+1}) = \alpha P(e_{t+1} / X_{t+1}) \sum_{X_t} P(X_{t+1} / x_t) P(x_t / e_{1:t})$$

Smoothing

 $\mathbf{P}(\mathbf{X}_k|\mathbf{e}_{1:t}) = \alpha \mathbf{P}(\mathbf{X}_k|\mathbf{e}_{1:k})\mathbf{P}(\mathbf{e}_{k+1:t}|\mathbf{X}_k)$

 $\mathbf{P}(\mathbf{e}_{k+1:t}|\mathbf{X}_k) = \sum_{\mathbf{x}_{k+1}} P(\mathbf{e}_{k+1}|\mathbf{x}_{k+1}) P(\mathbf{e}_{k+2:t}|\mathbf{x}_{k+1}) \mathbf{P}(\mathbf{x}_{k+1}|\mathbf{X}_k)$

Most likely sequence

 $\max_{\mathbf{x}_{1}...\mathbf{x}_{t}} \mathbf{P}(\mathbf{x}_{1},\ldots,\mathbf{x}_{t},\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1})$ = $\mathbf{P}(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) \max_{\mathbf{x}_{t}} \left(\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_{t}) \max_{\mathbf{x}_{1}...\mathbf{x}_{t-1}} P(\mathbf{x}_{1},\ldots,\mathbf{x}_{t-1},\mathbf{x}_{t}|\mathbf{e}_{1:t}) \right)$





Forward and backward messages as column vectors:

 $\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^{\top} \mathbf{f}_{1:t}$ $\mathbf{b}_{k+1:t} = \mathbf{T} \mathbf{O}_{k+1} \mathbf{b}_{k+2:t}$

Forward-backward algorithm needs time $O(S^2t)$ and space O(St)

UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEMI

Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \alpha \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\alpha'(\mathbf{T}^{\mathsf{T}})^{-1} \mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \mathbf{f}_{1:t}$$





Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \alpha \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\alpha'(\mathbf{T}^{\mathsf{T}})^{-1} \mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \mathbf{f}_{1:t}$$

Algorithm: forward pass computes \mathbf{f}_t , backward pass does \mathbf{f}_i , \mathbf{b}_i





Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \alpha \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\alpha'(\mathbf{T}^{\mathsf{T}})^{-1} \mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \mathbf{f}_{1:t}$$





Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \alpha \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\alpha'(\mathbf{T}^{\mathsf{T}})^{-1} \mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \mathbf{f}_{1:t}$$





Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \alpha \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\alpha'(\mathbf{T}^{\mathsf{T}})^{-1} \mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \mathbf{f}_{1:t}$$





Can avoid storing all forward messages in smoothing by running forward algorithm backwards:

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \alpha \mathbf{T}^{\mathsf{T}} \mathbf{f}_{1:t}$$
$$\alpha'(\mathbf{T}^{\mathsf{T}})^{-1} \mathbf{O}_{t+1}^{-1} \mathbf{f}_{1:t+1} = \mathbf{f}_{1:t}$$

Algorithm: forward pass computes \mathbf{f}_t , backward pass does \mathbf{f}_i , \mathbf{b}_i





Example



$$A^{-1} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{\det(A)} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \qquad U_2 = true \quad O_2 = \begin{pmatrix} 0.9 & 0 \\ 0 & 0.2 \end{pmatrix}$$
$$O_2^{-1} = 5.5 \quad \begin{pmatrix} 0.2 & 0 \\ 0 & 0.9 \end{pmatrix} = \begin{pmatrix} 1.1 & 0 \\ 0 & 4.95 \end{pmatrix} \qquad T^{\mathsf{T}} = \begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix}$$
$$(\mathsf{T}^{\mathsf{T}})^{-1} = 2.5 \quad \begin{pmatrix} 0.7 & -0.3 \\ -0.3 & 0.7 \end{pmatrix} = \begin{pmatrix} 1.75 & -0.75 \\ -0.75 & 1.75 \end{pmatrix}$$
$$\begin{pmatrix} 1.925 & -3.7125 \\ -0.825 & 8.6625 \end{pmatrix} \begin{pmatrix} 0.883 \\ 0.117 \end{pmatrix} = \begin{pmatrix} 1.265 \\ 0.285 \end{pmatrix} \qquad \alpha = 0.64497 \rightarrow \begin{pmatrix} 0.817 \\ 0.183 \end{pmatrix}$$







Set of states: $\{s_1, s_2, ..., s_N\}$

Process moves from one state to another generating a sequence of states : $s_{i1}, s_{i2}, \dots, s_{ik}, \dots$ Markov chain property: probability of each subsequent state depends only on what was the previous state:

$$P(s_{ik} | s_{i1}, s_{i2}, \dots, s_{ik-1}) = P(s_{ik} | s_{ik-1})$$

States are not visible, but each state randomly generates one of M observations (or visible states)

$$\{v_1, v_2, ..., v_M\}$$



State View: Hidden Markov models

To define a hidden Markov model, the following probabilities have to be specified:

- Matrix of transition probabilities $A=(a_{ij}), a_{ij}=P(s_j | s_i)$,
- Matrix of observation probabilities $B=(b_i(v_m))$, $b_i(v_m) = P(v_m \mid s_i)$ and a
- Vector of initial probabilities $\pi = (\pi_i), \pi_i = P(s_i)$.

Model is represented by $M=(A, B, \pi)$.



Example of Hidden Markov Model





IM FOCUS DAS LEBEN

Given some training observation sequences $O=O_1O_2...O_k$ and general structure of HMM (numbers of hidden and visible states), determine HMM parameters $M=(A, B, \pi)$ that best fit training data, that is maximizes $P(O \mid M)$.



State View: Learning problem (2)

If training data has information about sequence of hidden states, then use maximum likelihood estimation of parameters:

 $a_{ij} = P(s_j | s_i) = \frac{\text{Number of transitions from state } S_i \text{ to state } S_j}{\text{Number of transitions out of state } S_i}$ $b_i(v_m) = P(v_m | s_i) = \frac{\text{Number of transitions out of state } S_i}{\text{Number of times observation } V_m \text{ occurs in state } S_i}$

Otherwise: Use iterative expectation-maximization algorithm to find local maximum of $P(O \mid M)$: **Baum-Welch Algorithm.**



IM FOCUS DAS LEBEN

Baum-Welch algorithm

General idea:

$$a_{ij} = P(s_j | s_i) =$$

Expected number of transitions from state \mathbf{S}_i to state \mathbf{S}_j

Expected number of transitions out of state \mathbf{S}_i

 $b_{\scriptscriptstyle i}(v_{\scriptscriptstyle m}) = P(v_{\scriptscriptstyle m} | s_{\scriptscriptstyle i}) =$

Expected number of times observation V_m occurs in state S_i

Expected number of times in state S_i

 $\pi_i = P(s_i) = Expected frequency in state S_i at time k=1.$



IM FOCUS DAS LEBEN

Baum-Welch algorithm: Expectation step(1)

Define variable $\xi_{k}(i,j)$ as the probability of being in state S_{i} at time k and in state S at time k+1, given the observation sequence $O_1 O_2 \dots O_T$ with k < T $\xi_{k}(i,j) = P(q_{k} = S_{i}, q_{k+1} = S_{i} | O_{1}O_{2}...O_{T})$ $\xi_{k}(i,j) = \frac{P(q_{k} = s_{i}, q_{k+1} = s_{j}, o_{1} o_{2} \dots o_{T})}{P(o_{1} o_{2} \dots o_{T})}$ $P(q_k = s_i, o_1 o_2 \dots o_k) a_{ij} b_j(o_{k+1}) P(o_{k+2} \dots o_T | q_{k+1} = s_i)$ $P(O_1 O_2 ... O_T)$ α forward_k(i) a_{ii} b_i(o_{k+1}) backward_{k+1}(j)



Baum-Welch algorithm: Expectation step(2)

Define variable $\gamma_k(i)$ as the probability of being in state S_i at time k, given the observation sequence $O_1 O_2 \dots O_T$.

$$\gamma_{k}(i) = P(\mathbf{q}_{k} = \mathbf{S}_{i} | \mathbf{O}_{1} \mathbf{O}_{2} \dots \mathbf{O}_{T})$$

$$\gamma_{k}(i) = \frac{P(q_{k} = s_{i}, o_{1} o_{2} ... o_{T})}{P(o_{1} o_{2} ... o_{T})} =$$

 α forward_k(i) backward_k(i)



IM FOCUS DAS LEBEN

Baum-Welch algorithm: Expectation step(3)

We calculated
$$\xi_{k}(i,j) = P(\mathbf{q}_{k} = \mathbf{S}_{i}, \mathbf{q}_{k+1} = \mathbf{S}_{j} | \mathbf{O}_{1}\mathbf{O}_{2}...\mathbf{O}_{T})$$

and $\gamma_{k}(i) = P(\mathbf{q}_{k} = \mathbf{S}_{i} | \mathbf{O}_{1}\mathbf{O}_{2}...\mathbf{O}_{T})$

Expected number of transitions from state S_i to state $S_j =$

$$= \, \sum_k \, \xi_k({\rm i},{\rm j}) \,$$

Expected number of transitions out of state $S_i = \sum_k \gamma_k(i)$

Expected number of times observation V_m occurs in state $S_i = \sum_k \gamma_k(i)$, k is such that $O_k = V_m$ Expected frequency in state S_i at time k=1 : $\gamma_1(i)$.


Baum-Welch algorithm: Maximization step

$$a_{ij} = \frac{\text{Expected number of transitions from state } S_j \text{ to state } S_i}{\text{Expected number of transitions out of state } S_j} = \frac{\sum_k \xi_k(i,j)}{\sum_k \gamma_k(i)}$$
$$b_i(v_m) = \frac{\text{Expected number of times observation } v_m \text{ occurs in state } S_i}{\text{Expected number of times of times in state } S_i} = \frac{\sum_k \xi_k(i,j)}{\sum_{k,o_k = v_m} \gamma_k(i)}$$

$$\pi_i$$
 = (Expected frequency in state S_i at time $k=1$) = $\gamma_1(i)$.



DBN – Special Cases

• Hidden Markov Model (HMMs):

Temporal probabilistic model in which the state of the process is described by a single discrete random variable. (The simplest kind of DBN)

• Kalman Filter Models (KFMs):

Estimate the state of a physical system from noisy observations over time. Also known as linear dynamical systems (LDSs).



Kalman Filters

Modelling systems described by a set of continuous variables,

e.g., tracking a bird flying— $\mathbf{X}_t = X, Y, Z, \dot{X}, \dot{Y}, \dot{Z}$.

Airplanes, robots, ecosystems, economies, chemical plants, planets,



Gaussian prior, linear Gaussian transition model and sensor model



Updating Gaussian Distributions

Prediction step: if $P(\mathbf{X}_t | \mathbf{e}_{1:t})$ is Gaussian, then prediction

 $\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t}) = \int_{\mathbf{X}_t} \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) P(\mathbf{x}_t|\mathbf{e}_{1:t}) \, d\mathbf{x}_t$

is Gaussian. If $\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t})$ is Gaussian, then the updated distribution

 $\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1}) = \alpha \mathbf{P}(\mathbf{e}_{t+1}|\mathbf{X}_{t+1})\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t})$

is Gaussian

Hence $\mathbf{P}(\mathbf{X}_t | \mathbf{e}_{1:t})$ is multivariate Gaussian $N(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$ for all t

General (nonlinear, non-Gaussian) process: description of posterior grows unboundedly as $t \to \infty$



Simple 1-D Example



General Kalman Update

Transition and sensor models:

$$P(\mathbf{x}_{t+1}|\mathbf{x}_t) = N(\mathbf{F}\mathbf{x}_t, \mathbf{\Sigma}_x)(\mathbf{x}_{t+1})$$

$$P(\mathbf{z}_t|\mathbf{x}_t) = N(\mathbf{H}\mathbf{x}_t, \mathbf{\Sigma}_z)(\mathbf{z}_t)$$

Left for your studies

F is the matrix for the transition; Σ_x the transition noise covariance H is the matrix for the sensors; Σ_z the sensor noise covariance

Filter computes the following update:

$$\boldsymbol{\mu}_{t+1} = \mathbf{F}\boldsymbol{\mu}_t + \mathbf{K}_{t+1}(\mathbf{z}_{t+1} - \mathbf{H}\mathbf{F}\boldsymbol{\mu}_t)$$

$$\boldsymbol{\Sigma}_{t+1} = (\mathbf{I} - \mathbf{K}_{t+1})(\mathbf{F}\boldsymbol{\Sigma}_t\mathbf{F}^\top + \boldsymbol{\Sigma}_x)$$

where $\mathbf{K}_{t+1} = (\mathbf{F} \boldsymbol{\Sigma}_t \mathbf{F}^\top + \boldsymbol{\Sigma}_x) \mathbf{H}^\top (\mathbf{H} (\mathbf{F} \boldsymbol{\Sigma}_t \mathbf{F}^\top + \boldsymbol{\Sigma}_x) \mathbf{H}^\top + \boldsymbol{\Sigma}_z)^{-1}$ is the Kalman gain matrix

 Σ_t and \mathbf{K}_t are independent of observation sequence, so compute offline



2-D Tracking: Filtering





2-D Tracking: Smoothing





Where it breaks

Cannot be applied if the transition model is nonlinear

Extended Kalman Filter models transition as locally linear around $\mathbf{x}_t = \boldsymbol{\mu}_t$ Fails if systems is locally unsmooth

Standard solution: switching Kalman filter



Keeping track of many objects: Identity uncertainty



Web-Mining Agents Probabilistic Reasoning over Sequential Structures

Prof. Dr. Ralf Möller Universität zu Lübeck Institut für Informationssysteme

Tanya Braun (Übungen)



Every HMM is a single-variable DBN; every discrete DBN is an HMM





Consider the transition model

Sparse dependencies \Rightarrow exponentially fewer parameters;

e.g., 20 state variables, three parents each

DBN has $20 \times 2^3 = 160$ parameters, HMM has $2^{20} \times 2^{20} \approx 10^{12}$



Learning (1)

- The techniques for learning DBN are mostly straightforward extensions of the techniques for learning BNs
- Parameter learning
 - The transition model $P(X_t | X_{t-1})$ / The observation model $P(Y_t | X_t)$
 - Offline learning
 - Parameters must be tied across time-slices
 - The initial state of the dynamic system can be learned independently of the transition matrix
 - Online learning
 - Add the parameters to the state space and then do online inference (filtering)
 - The usual criterion is maximum-likelihood(ML)
- The goal of parameter learning is to compute
 - $\theta^{*}_{ML} = argmax_{\theta}P(Y|\theta) = argmax_{\theta}log P(Y|\theta)$
 - $\theta^*_{MAP} = \operatorname{argmax}_{\theta} \log P(Y | \theta) + \log P(\theta)$
 - Two standard approaches: gradient ascent and EM(Expectation Maximization)



Learning (2)

- Structure learning
 - Intra-slice connectivity: Structural EM
 - Inter-slice connectivity:
 For each node in slice t, we must choose its parents from slice t-1
 - Given structure is unrolled to a certain extent,
 the inter-slice connectivity is identical for all pairs of slices:
 - Constraints on Structural EM

121/29



Constructing Dynamic Bayesian Networks

 X_t , E_t contain arbitrarily many variables in a replicated Bayes net





DBNs transient failure

For simplicity we assume that BMeter_t and Battery_t are taken from 0..5.



Generic gaussian error model produces "overreaction"

Explicit transient failure model required:

$$P(BMeter_t = 0 | Battery_t = 5) = 0.03$$



DBNs persistent failure



Additional variable required: BMBroken_t

Upper curve: transient failure with different observation sequences

Lower curve: persistent failure with two different observation sequences



Learning DBN pattern structures?

- Difficult
- Need "deep" domain knowledge



Recap: Exact Inference in DBNs

Naive method: unroll the network and run any exact algorithm for the whole BN



Problem: inference cost for each update grows with t

Rollup filtering: add slice t + 1, "sum out" slice t using variable elimination

Largest factor is $O(d^{n+1})$, update cost $O(d^{n+2})$ (cf. HMM update cost $O(d^{2n})$)

d = possible values for variables n = number of states Example: 20 state variables with 4 values each means 4²⁰⁺¹=4.398.046.511.104

Employ forward chaining (constant per time update but exponential in the number of variables per state)



Inference: Algorithms

- Exact Inference algorithms
 - Forwards-backwards smoothing algorithm (on any discrete-state DBN)
 - Kalman filtering and smoothing (for continuous variables)
 - The Frontier Algorithm (sweep a Markov blanket, the frontier set F, across the DBN, first forwards and then backwards)
 - The Interface Algorithm (use only the set of nodes with outgoing arcs to the next time slice to d-separate the past from the future)
- Approximate algorithms:
 - The Boyen-Koller (BK) algorithm (approximate the joint distribution over the interface as a product of marginals)
 - Factored Frontier (FF) Algorithm / Loopy propagation algorithm (LBP)
 - Stochastic sampling algorithm:
 - Importance sampling or MCMC (offline inference)
 - Particle filtering (PF) (online)



Approximate inference in DBNs

Central idea:

- Create N initial-state examples (from prior dist $P(\mathbf{X}_0)$)
- Based on the transition model, each sample is propagated forward by sampling the next state value x_{t+1} given the current value x_t
- Each sample is weighted by the likelihood it assigns to the new evidence $P(\mathbf{e}_{t+1} \mid \mathbf{x}_{t+1})$
- Resample to generate new population of N samples
- Select new sample based on its weight

Samples are called particles (\rightarrow Particle Filtering)



Particle Filtering

Basic idea: ensure that the population of samples ("particles") tracks the high-likelihood regions of the state-space

Replicate particles proportional to likelihood for \mathbf{e}_t



Widely used for tracking nonlinear systems, esp. in vision

Also used for simultaneous localization and mapping in mobile robots 10^5 -dimensional state space



Example



 $N(r_{t+1}|e) = \sum_{x_t} P(x_{t+1}|x_t) N(x_t|e)$ For rain = 0.7*8+0.3*2= 6.2 => 6 For not rain = 0.3 *8 + 0.7*2= 3.8 => 4

Suppose no umbrella for t+1 total weight(rain particles) = 0.1 * 6= 0.6 total weight(not rain) = 0.8 * 4= 3.2 Normalized =<0.17, 0.83>



function PARTICLE-FILTERING(e, N, dbn) **returns** a set of samples for the next time step **inputs**: e, the new incoming evidence

N, the number of samples to be maintained

dbn, a DBN with prior $\mathbf{P}(\mathbf{X}_0)$, transition model $\mathbf{P}(\mathbf{X}_1|\mathbf{X}_0)$, sensor model $\mathbf{P}(\mathbf{E}_1|\mathbf{X}_1)$ persistent: S, a vector of samples of size N, initially generated from $\mathbf{P}(\mathbf{X}_0)$ local variables: W, a vector of weights of size N

for i = 1 to N do $S[i] \leftarrow \text{sample from } \mathbf{P}(\mathbf{X}_1 \mid \mathbf{X}_0 = S[i])$ /* step 1 */ $W[i] \leftarrow \mathbf{P}(\mathbf{e} \mid \mathbf{X}_1 = S[i])$ /* step 2 */ $S \leftarrow \text{WEIGHTED-SAMPLE-WITH-REPLACEMENT}(N, S, W)$ /* step 3 */ return S



Particle Filtering (cntd.)

Assume consistent at time t: $N(\mathbf{x}_t | \mathbf{e}_{1:t}) / N = P(\mathbf{x}_t | \mathbf{e}_{1:t})$

Propagate forward: populations of \mathbf{x}_{t+1} are

 $N(\mathbf{x}_{t+1}|\mathbf{e}_{1:t}) = \sum_{\mathbf{x}_t} P(\mathbf{x}_{t+1}|\mathbf{x}_t) N(\mathbf{x}_t|\mathbf{e}_{1:t})$

Weight samples by their likelihood for e_{t+1} :

 $W(\mathbf{x}_{t+1}|\mathbf{e}_{1:t+1}) = P(\mathbf{e}_{t+1}|\mathbf{x}_{t+1})N(\mathbf{x}_{t+1}|\mathbf{e}_{1:t})$

Resample to obtain populations proportional to W:

 $N(\mathbf{x}_{t+1}|\mathbf{e}_{1:t+1})/N = \alpha W(\mathbf{x}_{t+1}|\mathbf{e}_{1:t+1}) = \alpha P(\mathbf{e}_{t+1}|\mathbf{x}_{t+1})N(\mathbf{x}_{t+1}|\mathbf{e}_{1:t})$ $= \alpha P(\mathbf{e}_{t+1}|\mathbf{x}_{t+1})\Sigma_{\mathbf{x}_t}P(\mathbf{x}_{t+1}|\mathbf{x}_t)N(\mathbf{x}_t|\mathbf{e}_{1:t})$ $= \alpha' P(\mathbf{e}_{t+1}|\mathbf{x}_{t+1})\Sigma_{\mathbf{x}_t}P(\mathbf{x}_{t+1}|\mathbf{x}_t)P(\mathbf{x}_t|\mathbf{e}_{1:t})$ $= P(\mathbf{x}_{t+1}|\mathbf{e}_{1:t+1})$



Summary

Temporal models use state and sensor variables replicated over time

Markov assumptions and stationarity assumption, so we need

- transition model $P(\mathbf{X}_t | \mathbf{X}_{t-1})$
- sensor model $\mathbf{P}(\mathbf{E}_t | \mathbf{X}_t)$

Tasks are filtering, prediction, smoothing, most likely sequence; all done recursively with constant cost per time step

Hidden Markov models have a single discrete state variable; used for speech recognition

Kalman filters allow n state variables, linear Gaussian, $O(n^3)$ update

Dynamic Bayes nets subsume HMMs, Kalman filters; exact update intractable

Particle filtering is a good approximate filtering algorithm for DBNs



HMM-LDA

- In traditional topic modeling, such as LDA, we remove most syntactic words (e.g., stopwords) since we are only interested in *meaning*.
- In doing so, we discard much of the structure, and all of the order the original author intended.
- In topic modeling, we are concerned long-range topic dependencies rather document structure.



Thomas L. Griffiths, Mark Steyvers, David M. Blei, and Joshua B. Tenenbaum. Integrating topics and syntax. In *Proc. of* NIPS'04, pp. 537-544. **2004**.

Introduction

- HMMs are useful for segmenting documents into different types of words, regardless of meaning.
- For example, all nouns will be grouped together because they play the same role in different passages/documents.
- Syntactic dependencies last at most for a sentence.
- The standardized nature of grammar means that it stays fairly constant across different contexts.



Combining syntax and semantics 1

- All words (both syntactic and semantic) exhibit short range dependencies.
- Only content words exhibit long range semantic dependencies.
- This leads to the HMM-LDA.
- HMM-LDA is a composite model, in which an HMM decides the parts of speech, and a topic model (LDA) extracts topics only those words which are deemed semantic.



Generative Process 1

Definitions

<u>Words</u> $\mathbf{w} = \{w_1, \dots, w_n\}$ form document d where each word w_i is one of \mathbf{W} words <u>Topic assignments</u> $\mathbf{z} = \{z_1, \dots, z_n\}$ for each word, where each z_i taking one of \mathbf{T} topics

 $\frac{1}{2}$ $\frac{1}$

<u>Class assignments</u> $\mathbf{c} = \{c_1, \ldots, c_n\}$ for each word, where each c_i taking one of **C** word classes

 $\theta^{(d)}$ Multinomial distribution over topics for document d

 $\phi^{(z)}$ Multinomial distribution over **semantic** words for topic indicated by z.

 $\phi^{(c)}$ Multinomial distribution over **non-semantic** words for class indicated by *class c*. $\pi^{(c_{i-1})}$ Transition probability from c_{i-1} to c_i



Generative Process 2

 $\begin{array}{l} \theta^{(d)} \sim \operatorname{Dirichlet}(\alpha) \\ \phi^{(z)} \sim \operatorname{Dirichlet}(\beta) \\ \pi^{(c)} \sim \operatorname{Dirichlet}(\gamma) \quad \text{Where} \pi^{(c)} \text{ is the row of the transition matrix indicated by c.} \\ \phi^{(c)} \sim \operatorname{Dirichlet}(\delta) \end{array}$



Graphical Model 1





Simple Example 1



Preposition class

Verb Class

 The HMM allocates words which vary across context to the semantic class, since grammar is fairly standardized but content is not.



Model Inference 1

MCMC inference

Topic indicators

$$P(z_i | \mathbf{z}_{-i}, \mathbf{c}, \mathbf{w}) \propto P(z_i | \mathbf{z}_{-i}) \qquad P(w_i | \mathbf{z}, \mathbf{c}, \mathbf{w}_{-i})$$

$$\propto \begin{cases} n_{z_i}^{(d_i)} + \alpha & c_i \neq 1 \\ (n_{z_i}^{(d_i)} + \alpha) & \frac{n_{w_i}^{(z_i)} + \beta}{n^{(z_i)} + W\beta} & c_i = 1 \end{cases}$$

 $n_{z_i}^{(d_i)}$ is the number of words in document d_i signed to topic z_i

 $n_{w_i}^{(z_i)}$ is the number of words in topic z_i that are the same as w_i

All counts include only words for which $c_i = 1$ and exclude word i_i



Model Inference 2

Class indicators

/ERSITÄT ZU LÜBECK

$$P(c_{i}|\mathbf{c}_{-i}, \mathbf{z}, \mathbf{w}) \propto P(w_{i}|\mathbf{c}, \mathbf{z}, \mathbf{w}_{-i}) \qquad P(c_{i}|\mathbf{c}_{-i}) \\ \propto \begin{cases} \frac{n_{w_{i}}^{(c_{i})} + \delta}{n_{\cdot}^{(c_{i})} + W\delta} & \frac{(n_{c_{i}}^{(c_{i-1})} + \gamma)(n_{c_{i+1}}^{(c_{i})} + I(c_{i-1} = c_{i}) \cdot I(c_{i} = c_{i+1}) + \gamma)}{n_{\cdot}^{(c_{i})} + I(c_{i-1} = c_{i}) + C\gamma} & c_{i} \neq 1 \\ \frac{n_{w_{i}}^{(z_{i})} + \beta}{n_{\cdot}^{(z_{i})} + W\beta} & \frac{(n_{c_{i}}^{(c_{i-1})} + \gamma)(n_{c_{i+1}}^{(c_{i})} + I(c_{i-1} = c_{i}) \cdot I(c_{i} = c_{i+1}) + \gamma)}{n_{\cdot}^{(c_{i})} + I(c_{i-1} = c_{i}) + C\gamma} & c_{i} = 1 \end{cases}$$

 $n_{z_i}^{(d_i)}$ is the number of words in document d_i signed to topic z_i $n_{w_i}^{(z_i)}$ is the number of words in topic z_i at are the same as w_i $n_{w_i}^{(c_i)}$ is the number of words in class c_i at are the same as w_i $n_{c_i}^{(c_{i-1})}$ is the number of transitions from class c_{i-1} : lass c_i $I(\cdot)$ is an indicator variable which equals 1 if argument is true

All counts exclude transitions to and from *c*_i

- If we set the number of semantic topics to T = 1, then the model reduces to an HMM parts of speech tagger.
- If we set the number of HMM classes to C = 2, where one state is for punctuation, the the model reduces to LDA.



Results 1

Brown + TASA corpus: 38,151 documents; Vocab Size = 37,202; number of word tokens = 13,328,397 words

the	the	the	the	the	а	the	the	the	-
blood	,	,	of	a	the	,	2	,	
,	and	and	,	of	of	of	a	a	
of	of	of	to	2	2	a	of	in	
body	а	in	in	in	in	and	and	game	LDA only
heart	in	land	and	to	water	in	drink	ball	·
and	trees	to	classes	picture	is	story	alcohol	and	
in	tree	farmers	government	film	and	is	to	team	
to	with	for	a	image	matter	to	bottle	to	
is	on	farm	state	lens	are	as	in	play	
									-
blood	forest	farmers	government	lıght	water	story	drugs	ball	-
heart	trees	land	state	eye	matter	stories	drug	game	
pressure	forests	crops	federal	lens	molecules	poem	alcohol	team	
body	land	farm	public	image	liquid	characters	people	*	
lungs	soil	food	local	mirror	particles	poetry	drinking	baseball	HMM-I DA Semantic
oxygen	areas	people	act	eyes	gas	character	person	players	
vessels	park	farming	states	glass	solid	author	effects	football	Topics
arteries	wildlife	wheat	national	object	substance	poems	marijuana	player	•
*	area	farms	laws	objects	temperature	life	body	field	
breathing	rain	com	department	lenses	changes	poet	use	basketball	_
the	in	he	*	be	said	can	time	,	-
а	for	it	new	have	made	would	way	;	
his	to	you	other	see	used	will	years	(
this	on	they	first	make	came	could	day	:	HMM-I DA Syntactic
their	with	i	same	do	went	may	part)	
these	at	she	great	know	found	had	number		Classes
your	by	we	good	get	called	nust	kind		
her	from	there	small	go		do	place		
my	as	this	little	take		have			1 FOCUS DAS LEBEN
some	into	who	old	find		did			
Results 2

NIPS Papers 1713 documents; Vocabulary Size: 17268; Number of word tokens = 4,321,614

image	data	state	membrane	chip	experts	kernel	network	-
images	gaussian	poncy	synaptic	anatog	expert	support	neurai	
object	mixture	varue	cell	neuron	gating	vector	networks	
objects	likelihood	function	*	dıgıtal	hme	svm	output	
feature	posterior	action	current	synapse	architecture	kernels	input	Suptactic Morde
recognition	prior	reinforcement	dendritic	neural	mixture	#	training	Syntactic words
views	distribution	learning	potential	hardware	learning	space	inputs	
#	em	classes	neuron	weight	mixtures	function	weights	
pixel	bayesian	optimal	conductance	#	function	machines	Ŧ	
visual	parameters	8	channels	vlsi	gate	set	outputs	_
m	15	see	used	model	networks	however	#	-
with	was	show	trained	algorithm	values	also	*	
for	has	note	obtained	system	results	then	i	
on	becomes	consider	described	case	models	thus	x	
from	denotes	assume	given	problem	parameters	therefore	t	
at	being	present	found	network	units	first	n	Semantic Words
using	remains	need	presented	method	data	here	-	
into	represents	propose	defined	approach	functions	now	с	
over	exists	describe	generated	paper	problems	hence	r	
within	seems	suggest	shown	process	algorithms	finally	р	-



Results 2 (cont'd)

NIPS Papers 1713 documents; Vocabulary Size: 17268; Number of word tokens = 4,321,614

In contrast to this approach, we study here how the overall network activity can control single cell parameters such as input resistance, as well as time and space constants, parameters that are crucial for excitability and spariotemporal (sic) integration.

1.

2

The integrated architecture in this paper combines feed forward control and error feedback adaptive control using neural networks.

In other words, for our proof of convergence, we require the softassign algorithm to return a doubly stochastic matrix as *sinkhorn theorem guarantees that it will instead of a matrix which is merely close to being doubly stochastic based on some reasonable metric.

The aim is to construct a portfolio with a maximal expected return for a given risk level and time horizon while simultaneously obeying *institutional or *legally required constraints.

The left graph is the standard experiment the right from a training with # samples.
The graph G is called the *guest graph, and H is called the host graph.

Black words are semantic, Graylevel words are syntactic. Boxed words are semantic on one passsage and syntactic in another. Asterisked words have low frequency and not considered.

Results 3

Log Marginal probabilities of the data



Figure 5: Log marginal probabilities of each corpus under different models. Labels on horizontal axis indicate the order of the HMM.



Results 4

Parts of speech tagging

Black bars indicate performance on a fine tagset (297 word types), white bars indicate performance on coarse tagset (10 word types).





Conclusion

- HMM-LDA is a composite topic model which considers both long range semantic dependencies and short range syntactic dependencies.
- The model is quite competitive with a traditional HMM parts of speech tagger, and outperforms LDA when stopwords and punctuation are not removed.



- In LDA the order of documents does not matter
- Not appropriate for sequential corpora (e.g., that span hundreds of years)
- Further, we may want to track how language changes over time
- Let the topics *drift* in a sequence.



My fellow citizens: I stand here today humbled by the task before us, grateful for the trust you have bestowed, mindful of the sacrifices borne by our ancestors...





AMONG the vicissitudes incident to life no event could have filled me with greater anxieties than that of which the notification was transmitted by your order...



David M. Blei and John D. Lafferty. Dynamic topic models. In Proc. ICML '06. pp. 113-120. **2006**.

IM FOCUS DAS LEBEN 150

2009

Recap: Smoothed LDA Model



- Give a different word distribution to each topic
 - β is K×V matrix (V vocabulary size), each row denotes word distribution of a topic
- For each document d
 - Choose $\theta_d \sim \text{Dirichlet}(\cdot \mid \alpha)$
 - Choose $\beta_k \sim \text{Dirichlet}(\eta \cdot |)$
 - For each position $i = 1, ..., N_d$
 - Generate a topic $z_i \sim Mult(\cdot \mid \theta_d)$
 - Generate a word $w_i \sim Mult(\cdot | z_{i'}\beta_k)$





Topics drift through time

2



- Use a logistic normal distribution to model topics evolving over time
- Embed it in a state-space model on the log of the topic distribution

$$\beta_{t,k} | \beta_{t-1,k} \sim \mathcal{N}(\beta_{t-1,k}, I\sigma^2)$$

 $p(w | \beta_{t,k}) \propto \exp{\{\beta_{t,k}\}}$

• Lets us make inferences about sequences of documents



Logit Normal Distribution

The probability density function (PDF) of a logit-normal distribution, for $0 \le x \le 1$, is:

$$f_X(x;\mu,\sigma) = rac{1}{\sigma\sqrt{2\pi}} \, rac{1}{x(1-x)} \, e^{-rac{(\log ext{i}(x)-\mu)^2}{2\sigma^2}}$$

where μ and σ are the mean and standard deviation of the variable's logit (by definition, the variable's logit is normally distributed).



[Wikipedia]

UNIVERSITÄT ZU LÜBECK INSTITUT FÜR INFORMATIONSSYSTEME



Original article

Topic proportions

TECHVIEW: DNA SEQUENCING

Sequencing the Genome, Fast James C. Mullikin and Amanda A. McMurray

100 -. ...

...

0 100 200

Genome sequencing projects reveal the genetic makeup of an organism provide the genetic makeup of the sequence of the DNA bases, which encodes all of the infor-are aiming to sequence I Ob of human semation necessary for the life of the organism. The base sequence contains four nu-cleotides-adenine, thymidine, guanosine, and cytosine-which are linked together into long double-helical chains. Over the last two decades, automated DNA se-quencers have made the process of obtaining the base-by-base sequence of DNA easier. By application of an electric field across a gel matrix, these sequencers sepa-rate fluorescently labeled DNA molecules that differ in size by one base. As the molecules move past a given point in the gel, laser excitation of a fluorescent dye specific to the base at the end of the molecule yields a base-specific signal that can be automatically recorded.

The latest sequencer to be launched is Perkin-Elmer's much-anticipated ABI Prism 3700 DNA Analyzer which, like the Molecular Dynamics MegaBACE 1000 launched last year, incorporates a capillary tube to hold the sequence gel rather than a traditional slab-shaped gel apparatus. Extra interest in the ABI 3700 has been generated because Craig Venter of Celera Ge-nomics Corporation anticipates that ~230 of these machines (1) will enable the company to produce raw sequence for the enin 3 years. The specifications of the ABI 3700 machine say that, with less than 1 hour of human labor per day, it can se-

ples from the plates into wells that open in-to the capillaries. This and the rest of the sequencing operation is fully automatic. The machine can currently process four 96-well plates of DNA samples unattended, taking approximately 16 hours before oper-ator intervention is required. This rate falls

the use of a sheath flow fluorescence detec tion system (4). Detection of the DNA frag-ments occurs 300 µm past the end of the cap-illary within a fused silica cuvette. A laminar fluid flows over the ends of the capillaries drawing the DNA fragments as they emergi from the capillaries through a fixed lase beam that simultaneously intersects with all gel sequencers from Perkin-Elmer plus 6 Molecular Dynamics MegaBACE 1000 capillary sequencers, allowing a maximum of the samples. The emitted fluorescence is detected with a spectral CCD (charge-cou-pled device) detector. This arrangement throughput of 32,000 samples per day. Two ABI 3700 capillary sequencers-delivered means that there are no moving parts in the detection system, other than a shutter in front



Fig. 1. Comparison of read-length histograms for h

quence 768 samples per day. Assuming to the Sanger Centre in December 1998— are desirable. In fact, a system that could that each sample prives an average of 400 are in our Research and Development de read twice as many bases but at half the base pairs (bp) of usable sequence data (its partment for evaluation. Thus, the ABI speed of another system is preferable, if read length) and any section from the en 3700 will ultimately be added to our pre- both systems cost the same. This is beread length) and any section from the error the human genome is covered by an aver-age of 10 overlapping independent reads (2), the 75 millions samples that Celerar into a floor-standing cabinet, which con-bing many short ones. So, read length is

(2), the 75 million samples that Celera into a floor-standing cabinet, which com-must process will require = 10000 ABI 3700 machine days. With ~230 machines, 434 days, which affords some margin of er-for fur uncepted developments. At the Sanger Centre, we have finished a data first a sumple are located for the second structure of the At the Sanger Centre, we have finished a data first a sumple are located and structure of the second structure of the sec 140 Mb of genomic sequence from a vari- The authors are at The Sanger Centre, Wellcome That Centre, Grande, Carlos, Carl

www.sciencemag.org SCIENCE VOL 283 19 MARCH 1999



short of the design specification of four 96-well plates in 12 hours. The main innovation of the ABI 3700 is quence in rough-draft form by 2001, with a finished version by 2003. Our sequenc-ing equipment includes 44 ABI 373XL, 61 ABI 377XL, and 31 ABI 377XL-96 slab

of the CCD detector. We have evaluated these machines for their performance, op-eration, ease of use, and reliabili-ty in comparison to the more commonly used slab gel se-

quencing machines. In automat-ed sequencers, there are two methods for containing the gel matrix. One is to polymerize a gel matrix between two finely separated glass plates (0.4 mm or less)—the slab gel method. The other is to inject a polymer ma-

DNA-that is, long read lengths

1867



Original article

SCIENCE S COMPASS + TECH.SIGHT

Most likely words from top topics

TECHVIEW: DNA SEQUENCING

Sequencing the Genome, Fast James C. Mullikin and Amanda A. McMurray

enome sequencing projects reveal the genetic makeup of an organism by reading off the sequence of the DNA bases, which encodes all of the infor-mation necessary for the life of the organism. The base sequence contains four nucleotides-adenine thymidine guanosine and cytosine-which are linked together into long double-helical chains. Over the last two decades, automated DNA sequencers have made the process of obtain-ing the base-by-base sequence of DNA easier. By application of an electric field across a gel matrix, these sequencers sepa-rate fluorescently labeled DNA molecules that differ in size by one base. As the molecules move past a given point in the gel, laser excitation of a fluorescent dye specific to the base at the end of the molecule yields a base-specific signal that can be automatically recorded. The latest sequencer to be launched is Perkin-Elmer's much-anticipated ABI

Prism 3700 DNA Analyzer which, like the Molecular Dynamics MegaBACE 1000 launched last year, incorporates a capillary tube to hold the sequence gel rather than a traditional slab-shaped gel apparatus. Extra interest in the ABI 3700 has been generated because Craig Venter of Celera Genomics Corporation anticipates that ~230 of these machines (1) will enable the company to produce raw sequence for the enin 3 years. The specifications of the ABI 3700 machine say that, with less than 1 hour of human labor per day, it can sequence 768 samples per day, it can ac-that each sample gives an average of 400 hase pairs (bp) of usable sequence data (its read length) and any section from the en-tire human genome is covered by an aver-age of 10 overlapping independent reads (2), the 75 million samples that Celera must process will require ~100,000 ABI 3700 machine days. With ~230 machines, that works out to less than 2 years or about 434 days, which affords some margin of er-ror for unexpected developments. At the Sanger Centre, we have finished

146 Mb of genomic sequence from a vari-

The authors are at The Sanger Centre, Wellcome Trust Genome Campus, Hinxton, Cambs, CB10 15A, Trust Genome Campus, Hinxt UK. E-mail: jcm@sanger.ac.uk

puter. A robotic arm transfers DNA samwww.sciencemag.org SCIENCE VOL 283 19 MARCH 1999

ples from the plates into wells that open in-to the capillaries. This and the rest of the equencing operation is fully automatic The machine can currently process four 96-well plates of DNA samples unattended, taking approximately 16 hours before operator intervention is required. This rate falls short of the design specification of four 96-well plates in 12 hours. The main innovation of the ABI 3700 is ety of genomes, including 81 Mb of se-quence from the human genome, the the use of a sheath flow fluorescence detection system (4). Detection of the DNA frag largest amount of any center so far (3). We ments occurs 300 µm past the end of the cap are aiming to sequence 1 Gb of human se-quence in rough-draft form by 2001, with a finished version by 2003. Our sequencillary within a fused silica cuvette. A lamina fluid flows over the ends of the capillaries, drawing the DNA fragments as they emerge from the capillaries through a fixed laser

ing equipment includes 44 ABI 373XL, 61 ABI 377XL, and 31 ABI 377XL-96 slab beam that simultaneously intersects with all of the samples. The emitted fluorescence is gel sequencers from Perkin-Elmer plus 6 Molecular Dynamics MegaBACE 1000 detected with a spectral CCD (charge-cou capillary sequencers, allowing a maximum throughput of 32,000 samples per day. Two ABI 3700 capillary sequencers—delivered pled device) detector. This arrangement means that there are no moving parts in the detection system, other than a shutter in front of the CCD detector



trix into a capillary (internal di-ameter <0.2 mm). Most sequenc-Fig. 1. Comparison of read-length histograms for se-quences collected with the ABI 3700 capillary machine and the ABI 377XL-96 slab gel machine. The capillary machine ing facilities use the slab gel method, because multicapillary under-performs the slab gel machine by about 200 bases. sequencers have only recently Both sets of reads are from runs with ABI Big Dye Terminaboth sets of reads are right furth ABI big type remma-tor chemistrics. Read length is computed as the number of the commercially available. bases per read where the predicted error rate is less than or equal to 1.0% (Q ≥ 20). The "phred" Q value was recali bated for each type of read. With either type of system the aim is to read as many bases as possible for a given sample of DNA-that is, long read lengths

to the Sanger Centre in December 1998-are in our Research and Development de-partment for evaluation. Thus, the ABI are desirable. In fact, a system that could read twice as many bases but at half the speed of another system is preferable, if 3700 will ultimately be added to our presboth systems cost the same. This is beent capacity to reach our goal. The ABI 3700 DNA sequencer is built cause assembling relatively fewer long-se quenced fragments is easier than assem into a floor-standing cabinet, which conbling many short ones. So, read length is tains in its base all the reagents required for its operation. The reagent containers are an important parameter when evaluating new sequencing technologies.

We have directly compared the ABI 3700 sequencer to the ABI 377XL slab gel sequencer by evaluating the sequence data obtained from both machines with human readily accessible for replenishment, which is required every day under high-through-put operation. At bench height within the cabinet is a four-position bed, on which microtiter plates of DNA samples are located. The operator places the prepared plates in-to position, closes the front of the machine DNA samples. These samples were subcloned into plasmid or m13 phage and pre-pared and sequenced with our standard protocols for Perkin-Elmer Big Dye Terand programs it by using a personal com-

minator chemistry

1867

sequence genome genes sequences human gene dna sequencing chromosome regions analysis data genomic number

devices device materials current high gate light silicon material technology electrical fiber power based

data information network web computer language networks time software system words algorithm number internet









"Theoretical Physics"

"Neuroscience"



Wang, Chong; Blei, David; Heckerman, David. "Continuous Time Dynamic Topic Models". *Proceedings of ICML*'08, **2008**.

- **Time-corrected similarity** shows a new way of using the posterior.
- Consider the expected Hellinger distance between the topic proportions of two documents,

$$d_{ij} = \mathrm{E}\left[\sum_{k=1}^{K} (\sqrt{\theta_{i,k}} - \sqrt{\theta_{j,k}})^2 | \mathbf{w}_i, \mathbf{w}_j \right]$$

- Uses the latent structure to define similarity
- Time has been factored out because the topics associated to the components are different from year to year.
- Similarity based only on topic proportions

For two discrete probability distributions $P=(p_1,\ldots,p_k)$ and $Q=(q_1,\ldots,q_k)$, their Hellinger distance is defined as



$$H(P,Q) = rac{1}{\sqrt{2}} \; \sqrt{\sum_{i=1}^k (\sqrt{p_i} - \sqrt{q_i})^2},$$
 [Wikipedia]

The Brain of the Orang (1880)





Representation of the Visual Field on the Medial Wall of Occipital-Parietal Cortex in the Owl Monkey (1976)





Dynamic Topic Models: Summary

- The Dirichlet assumption on topics and topic proportions makes strong conditional independence assumptions about the data.
- The **dynamic topic model** uses a logistic normal in a linear dynamic model to capture how topics change over time.













Source: Gartner (July 2016)





