Intelligent Agents Probabilistic Models for Sequential Structures

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



IM FOCUS DAS LEBEN

Motivation: Part Of Speech Tagging

- Annotate each word in a sentence with a part-ofspeech (POS) tags.
- 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
- Topic modeling as discussed before could be extended to better consider POS tags



Abbreviations: https://sites.google.com/site/partofspeechhelp/home

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"



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.



Forward Classification





IM FOCUS DAS LEBEN 6

Forward Classification





IM FOCUS DAS LEBEN 7





IM FOCUS DAS LEBEN 8











Time and Uncertainty

- The world changes, we need to track and predict it
- Examples: diabetes management, traffic monitoring
- Uncertainty is everywhere
- Need temporal probabilistic graphical models
- 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 $\mathbf{X}_{a:b}$ denotes the set of variables from \mathbf{X}_{a} to \mathbf{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.



Example





DBN - Representation

- Problem: all previous random variables could have an influence on those of the current timestamp
 - 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).

 $P(U_t / Parent(U_t))$ is the same for all t



Stationary Process/Markov Assumption

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

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

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

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

Sensor Model

- 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



Example





Complete Joint Distribution: Markov-1

- 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:

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



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?



DBN – Basic Inference

• Filtering or Monitoring:

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

 $P(X_t / e_{1 \cdot t})$

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.



DBN – Basic Inference

• 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}

$$\begin{split} P(X_{t+1} / e_{1:t+1}) &= f(e_{t+1,} P(X_t / e_{1:t+1})) & \text{(for some function f)} \\ &= P(X_{t+1} / e_{1:t,} e_{t+1}) & \text{(dividing up the evidence)} \\ &= \alpha P(e_{t+1} / X_{t+1,} e_{1:t}) P(X_{t+1} / e_{1:t}) & \text{(using Bayes' Theorem)} \\ &= \alpha P(e_{t+1} / X_{t+1}) P(X_{t+1} / e_{1:t}) & \text{(by the Markov property} \\ &= \alpha P(e_{t+1} / X_{t+1}) P(X_{t+1} / e_{1:t}) & \text{of evidence)} \end{split}$$

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



Bayes Rule

$P(A \mid B) = P(A, B) / P(B)$

P(A,B) = P(A | B) P(B) = P(B | A) P(A) = P(B, A)



IM FOCUS DAS LEBEN

Application of Bayes Rule

 $P(A \mid B, C) = P(A, B, C) / P(B, C)$ = P(C, A, B) / P(B, C) = P(C | A, B) P(A, B) / P(B, C) = P(C | A, B) P(A | B) P(B) / (P(C | B) P(B)) = α P(C | A, B) P(A | B)

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



DBN – Basic Inference

• Filtering cont.

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

$$\begin{split} P(X_{t+1} / e_{1:t+1}) &= f(e_{t+1,} P(X_t / e_{1:t+1})) & \text{(for some function f)} \\ &= P(X_{t+1} / e_{1:t,} e_{t+1}) & \text{(dividing up the evidence)} \\ &= \alpha P(e_{t+1} / X_{t+1,} e_{1:t}) P(X_{t+1} / e_{1:t}) & \text{(using Bayes' Theorem)} \\ &= \alpha P(e_{t+1} / X_{t+1}) P(X_{t+1} / e_{1:t}) & \text{(by the Markov property} \\ &= \alpha P(e_{t+1} / X_{t+1}) P(X_{t+1} / e_{1:t}) & \text{of evidence)} \end{split}$$

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



Application of Bayes Rule

 $P(A | B) = \Sigma_{c} P(A, c | B)$ $= \Sigma_{c} P(A, c, B) / P(B)$

 $= \Sigma_{c} P(A | c, B) P(c, B) / P(B)$

= $\Sigma_{c} P(A \mid c, B) P(c \mid B) P(B) / P(B)$

 $= \Sigma_{c} P(A \mid c, B) P(c \mid B)$

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



DBN – Basic Inference

• 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)



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



Example $P(Rain_0) = (0.5 \ 0.5)^T$





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.





DBN – Basic Inference

• Prediction:

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

$$P(X_{t+k} / e_{1:t}) \qquad \text{for some } k > 0$$

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



DBN – Basic Inference

• Smoothing or hindsight:

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

$$P(X_k / e_{1:t})$$
 for 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 $O(t|\mathbf{f}|)$


Application of Bayes Rule

 $P(A \mid B, C) = P(A, B, C) / P(B, C)$ = P(C, A, B) / P(B, C) = P(C \mid A, B) P(A, B) / P(B, C) = P(C \mid A, B) P(A \mid B) P(B) / (P(C \mid B) P(B)) = α P(C \mid A, B) P(A \mid B)

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



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 $O(t|\mathbf{f}|)$



Application of Bayes Rule

 $P(A | B) = \Sigma_{c} P(A, c | B)$

- $= \Sigma_{c} P(A, c, B) / P(B)$
- $= \Sigma_{c} P(A | c, B) P(c, B) / P(B)$
- = $\Sigma_{c} P(A \mid c, B) P(c \mid B) P(B) / P(B)$
- $= \Sigma_{c} P(A | c, B) P(c | B)$

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



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.





DBN – Basic Inference

• 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)



DBN – Basic Inference

• 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.



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

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

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





Consider special case of a dynamic Bayesian Network:

- Use vector of independent state variables **X**_t
- Use vector of independent evidence variables E_t
- This was already used in the rain-umbrella example
- For high-dimensional vectors the transition and sensor models become quite complex: O(d²) space

NB:

- In a general dynamic Bayesian network, state variables are not necessarily independent
- Even evidence variable might be dependent on one another (naïve Bayes does not work)



How to Incorporate Context into LDA?





- 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.

1789



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... Inaugural addresses



2009

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**.



Topics drift through time

Intelligent Agents Probabilistic Models for Sequential Structures

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



IM FOCUS DAS LEBEN



- LDA is a simple topic model.
- It can be used to find topics that describe a corpus.
- Each document exhibits multiple topics.
- How can we build on this simple model of text?



Using and Embedding LDA

- LDA model used to infer posterior distribution
 P(Z | wd)
- Based on Z one can find and rank related documents
 - Infer $P(Relevant_d | Z_q)$ for d being the documents in a repository and q being the query document
 - Previously introduced models for information retrieval can be extended with topic information
 - Works for books, articles, images, videos, and other media
- LDA can be embedded in more complicated models
 - Model further intuitions about the structure of texts
 - Links, citations ("relational" topic models), ...



- Traditional topic modeling (e.g., LDA):
 - Interested in *meaning*
 - Remove most syntactic words (e.g., stopwords)
 - Discard much of the structure, and all order information that the original author intended
 - Concerned about long-range topic dependencies rather document structure
- Not always easy to decide which words to remove
 - Keep only nouns? Example: saw vs. saw



HMM

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



Combining Syntax and Semantics: HMM-LDA

- All words (both syntactic and semantic) exhibit short range dependencies.
- Only content (semantic) 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 from only those words which are deemed semantic



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**.

Definitions

<u>Words</u> $\mathbf{w} = \{w_1, \ldots, w_n\}$ form document d where each word w_i is one of **W** words

<u>Topic assignments</u> $\mathbf{z} = \{z_1, \dots, z_n\}$ for each word, where each z_i taking one of **T** topics

<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



How to Incorporate Context into LDA?





Dirichlet Distribution



- The Dirichlet is a distribution on the simplex, positive vectors that sum to 1.
- It assumes that components are nearly independent.
- In real data, an article about *fossil fuels* is more likely to also be about *geology* than about *genetics*.



Generative Process 2

 $\theta^{(d)} \sim \text{Dirichlet}(\alpha)$ $\phi^{(z)} \sim \text{Dirichlet}(\beta)$ $\pi^{(c)} \sim \text{Dirichlet}(\gamma)$ Where $\pi^{(c)}$ is the row of the transition matrix indicated by c. $\phi^{(c)} \sim \text{Dirichlet}(\delta)$ For document d Draw topic • 1. Sample $\theta^{(d)}$ from a Dirichlet(α) prior distribution 2. For each word w_i in document d Draw a topic • (a) Draw z_i from $\theta^{(d)}$ for word i (b) Draw c_i from $\pi^{(c_{i-1})}$ (c) If $c_i = 1$, then draw w_i from $\phi^{(z_i)}$, else draw w_i from $\phi^{(c_i)}$ Draw a class for word i from Semantic class transition matrix Draw a semantic word OR Draw a syntactic word VERSITÄT ZU LÜBECK **IM FOCUS DAS LEBEN**

Simplified Example



network used for images image obtained with kernel output described with objects neural network trained with svm images

Preposition class

Verb class



IM FOCUS DAS LEBEN

LDA-HMM: Summary

- HMM-LDA is a composite topic model
 - Long range semantic dependencies
 - Short-range syntactic dependencies
- Quite competitive with traditional HMM POS tagger
- Outperforms LDA when stop-words 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.

1789



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... Inaugural addresses



2009

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**.

Recap: Smoothed LDA Model



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





Topics drift through time



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

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

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

• Lets us make inferences about sequences of documents



Logistic Normal Distribution



- The **logistic normal** is a distribution on the simplex that can model dependence between components (Aitchison, 1980).
- The log of the parameters of the multinomial are drawn from a multivariate Gaussian distribution,

$$X \sim \mathcal{N}_{\mathcal{K}}(\mu, \Sigma)$$

 $\theta_i \propto \exp\{x_i\}.$



Logit Normal Distribution

Normal Distribution $f(x) = rac{1}{\sigma\sqrt{2\pi}} e^{-rac{1}{2}\left(rac{x-\mu}{\sigma}
ight)^2}$

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

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

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

....

Genome sequencing projects reveal the genetic makeup of an organism by reading off the sequence of the human genome, the largest amount of any center so far (3). We DNA bases, which encodes all of the information 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 obtaining the base-by-base sequence of DNA easier. By application of an electric field across a gel matrix, these sequencers separate 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 entire 3 gigabases (Gb) of the human genome in 3 years. The specifications of the ABI 3700 machine say that with less than 1 hour of human labor per day, it can se-quence 768 samples per day. Assuming that each sample gives an average of 400 base pairs (bp) of usable sequence data (its read length) and any section from the entire 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 error for unexpected developments.

At the Sanger Centre, we have finished 146 Mb of genomic sequence from a vari-crotter plates of DNA samples are located.

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

are aiming to sequence 1 Gb of human sequence in rough-draft form by 2001, with a finished version by 2003. Our sequencing equipment includes 44 ABI 373XL, 61 ABI 377XL, and 31 ABI 377XL-96 slab gel sequencers from Perkin-Elmer plus 6

throughput of 32,000 samples per day. Two ABI 3700 capillary sequencers-delivered of the CCD detector. We have evaluated these ma-



Fig. 1. Comparison of read-length histograms for sequences collected with the ABI 3700 capillary machine and the ABI 377XL-96 slab gel machine. The capillary machine under-performs the slab gel machine by about 200 bases. Both sets of reads are from runs with ABI Big Dye Terminator chemistries. Read length is computed as the number of bases per read where the predicted error rate is less than or equal to 1.0% ($Q \ge 20$). The "phred" Q value was recalible for a single part of the single part of brated for each type of read.

to the Sanger Centre in December 1998ent capacity to reach our goal. The ABI 3700 DNA sequencer is built into a floor-standing cabinet, which con-

tains in its base all the reagents required for its operation. The reagent containers are readily accessible for replenishment, which is required every day under high-throughput operation. At bench height within the The operator places the prepared plates into position, closes the front of the machine and programs it by using a personal com-puter. A robotic arm transfers DNA sam-

www.sciencemag.org SCIENCE VOL 283 19 MARCH 1999

an important parameter when evaluating new sequencing technologies.

protocols for Perkin-Elmer Big Dye Terminator chemistry

taking approximately 16 hours before oper-ator intervention is required. This rate falls ator intervention is required. This rate fails short of the design specification of four 96-well plates in 12 hours. The main innovation of the ABI 3700 is 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 capillary within a fused silica cuvette. A laminar fluid flows over the ends of the capillaries, drawing the DNA fragments as they emerge from the capillaries through a fixed laser beam that simultaneously intersects with all of the samples. The emitted fluorescence is Molecular Dynamics MegaBACE 1000 capillary sequencers, allowing a maximum detected with a spectral CCD (charge-cou-pled device) detector. This arrangement means that there are no moving parts in the detection system, other than a shutter in front

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,

chines for their performance, op-eration, ease of use, and reliability in comparison to the more commonly used slab gel seauencing machines. In automataencers, 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

trix into a capillary (internal diameter <0.2 mm). Most sequenc ing facilities use the slab gel method, because multicapillary sequencers have only recently 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 are desirable. In fact, a system that could are in our Research and Development de-partment for evaluation. Thus, the ABI 3700 will ultimately be added to our psystems cost the same. This is because assembling relatively fewer long-se-quenced fragments is easier than assem-bling many short ones. So, read length is

> 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 DNA samples. These samples were subcloned into plasmid or m13 phage and prepared and sequenced with our standard

> > 1867





Original article

TECHVIEW: DNA SEQUENCING

Sequencing the Genome, Fast

James C. Mullikin and Amanda A. McMurray

Genome sequencing projects reveal the genetic makeup of an organism by reading off the sequence of the DNA bases, which encodes all of the information 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 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 separate 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 generat-ed because Craig Venter of Celera Genomics Corporation anticipates that ~230 of these machines (1) will enable the comnany to produce raw sequence for the entire 3 gigabases (Gb) of the human genome in 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. Assuming that each sample gives an average of 400 base pairs (bp) of usable sequence data (its read length) and any section from the entire 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, UK.E-mail: jcm@sanger.ac.uk

ety of genomes, including 81 Mb of se-quence from the human genome, the largest amount of any center so far (3). We are aiming to sequence 1 Gb of human sequence 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 gel sequencers from Perkin-Elmer plus 6 Molecular Dynamics MegaBACE 1000 capillary sequencers, allowing a maximum throughput of 32,000 samples per day. Two ABI 3700 capillary sequencers-delivered

Fig. 1. Comparison of read-length histograms for se s collected with the ABI 3700 capillary machine and the ABI 377XL-96 slab gel machine. The capillary machine under-performs the slab gel machine by about 200 bases Both sets of reads are from runs with ABI Big Dye Terminator chemistries. Read length is computed as the number of bases per read where the predicted error rate is less than or equal to 1.0% ($Q \ge 20$). The "phred" Q value was recalibrated for each type of read.

to the Sanger Centre in December 1998are in our Research and Development de-partment for evaluation. Thus, the ABI 3700 will ultimately be added to our presit capacity to reach our goal. The ABI 3700 DNA sequencer is built

into a floor-standing cabinet, which con-tains in its base all the reagents required for its operation. The reagent containers are readily accessible for replenishment, which is required every day under high-throughput operation. At bench height within the cabinet is a four-position bed, on which mi-crotiter plates of DNA samples are located. The operator places the prepared plates in-to position, closes the front of the machine

and programs it by using a personal com-puter. A robotic arm transfers DNA samwww.sciencemag.org SCIENCE VOL 283 19 MARCH 1999

taking approximately 16 hours before oper-ator 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 the use of a sheath flow fluorescence detec-tion system (4). Detection of the DNA fragments 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 emerge from the capillaries through a fixed lase beam that simultaneously intersects with all of the samples. The emitted fluorescence is detected with a spectral CCD (charge-cou-

sequencing operation is fully automatic The machine can currently process four

96-well plates of DNA samples unattended.

• TECH.SIGHT ples from the plates into wells that open in-to the capillaries. This and the rest of the

> 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 We have evaluated these ma

> > 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 trix into a capillary (internal diameter <0.2 mm). Most sequence

ing facilities use the slab gel method, because multicapillary sequencers have only recently become commercially available. 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

are desirable. In fact, a system that could read twice as many bases but at half the speed of another system is preferable, if both systems cost the same. This is because assembling relatively fewer long-se-quenced fragments is easier than assembling many short ones. So, read length is 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 DNA samples. These samples were sub cloned into plasmid or m13 phage and pre-pared and sequenced with our standard protocols for Perkin-Elmer Big Dye Ter

chines for their performance, op-eration, ease of use, and reliabili-

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

Most likely words from top topics

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







IM FOCUS DAS LEBEN 70

FORCE **OXYGEN ASER** 0-0-0 NERVE 00 RELATIVITY NEURON 1880 1900 1920 1940 1980 1900 1980 1960 2000 1880 1920 1940 1960 2000

"Theoretical Physics"

"Neuroscience"



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

IM FOCUS DAS LEBEN 71

Probabilistic Topic Models, David Blei, 2013

- Understand developments
- Distributions of topics over time
- Discretization of time might be a problem
 - Runtime increases dramatically
 - Continuous dynamic topic models
- Many applications
 - E.g., comparison of science areas, analysis of scientific work
- How can we compare distributions?


Recap: Huffman code example





Μ	code 1	ength	prob	Exp. len
А	000	3	0,125	0,375
В	001	3	0,125	0,375
С	01	2	0,250	0,500
D	1	1	0,500	0,500
average message length				1,750

If we need to send many messages (A,B,C or D) and they have this probability distribution and we use this code, then over time, the average bits/message should approach 1.75

Recap: Information Theory Background

- Assume that you need to send messages from a repertoire of n messages
- If there are n equally probable possible messages, then the probability p of each is 1/n or n = 1/p
- Information (number of bits) conveyed by a message is log(n) = log(1/p)= -log(p)
- E.g., if there are 16 messages, then log(16) = 4 and we need 4 bits to identify/send each message.
- In general, if we are given a probability distribution

 $P = (p_1, p_2, .., p_n)$

• Information conveyed by distribution (aka entropy of P) is:

 $I(P) = -(p_1^* log(p_1) + p_2^* log(p_2) + ... + p_n^* log(p_n))$ = - \Sigma_i p_i^* log(p_i) = \Sigma_i p_i^* log(1/p_i)



The KL Divergence

 The cross-entropy, or Kullback-Leibler divergence, between two distributions p and q measures the expected information gain (reduction in average number of bits per event) due to replacing the "wrong" distribution q with the "right" distribution p:

$$D^{KL}(\mathbf{p},\mathbf{q}) \equiv \sum_{i} p_i (\ln(1/q_i) - \ln(1/p_i)) = \mathbf{E}_{\mathbf{p}} [\ln(\mathbf{p}/\mathbf{q})]$$

• Not symmetric



Hellinger Distance

 The Hellinger distance is a symmetric measure of distance between two distributions that is popular in machine learning applications:

$$D^{HEL}(\mathbf{p}, \mathbf{q}) \equiv \|\sqrt{\mathbf{p}} - \sqrt{\mathbf{q}}\|_2 = \left(\sum_{j=1}^n \left(\sqrt{p_j} - \sqrt{q_j}\right)^2\right)$$
$$\in [0, \sqrt{2}]$$

• Sometimes value should be in [0, 1]

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]

S DAS LEBEN

Dynamic Topic Models

- 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 of documents based only on topic proportions



Dynamic Topic Models

The Brain of the Orang (1880)

SCIENCE.

F. Barker, Professor O. C. Marsh

THE BRAIN OF THE ORANG.*

BY HENRY C. CHAPMAN, M.D.

The brain of the Orang has been figured by Tiede The brain of the Orag has been figured by Tiede-man, Sandiffort, Schweder an der Neck and Vrolik, Gratisiel, Rolleston, etc. On account, however, of the few illustrations extant, and of the importance of the subject, I avail myself of the copportunity of presenting several views of my Orang's brain (Gray, 1 to g), which was removed from the skull only a few hours afther death. The membranes were in a bigh state of congen-tion, and a little of the surface of the left hemisphere had been discontinger the bicases, obteneits the bach were to en disorganized by disease, otherwise the brain was in usorganized by disease, otherwise the orain was in condition. It weighed exactly ten ounces. The of the Orang in its general contour resembled that an more than those of either of the Chimpanzees I examined. In these the brain was more elong-The general character of the folds and fissures in



and man are th hree. The fissure of Sulvius nior branch pro--occipital fis-the first oc-al side of the ng the parietal from the occipal lobe

lemy of Natural oces, Phila., 1880.

is these cases, which were submitted to the the Orang, the parteto-occidital fastere does not react the the dist December laws for corrections on the second secon Gorilla, and se is to be i In the female Chir

F16. # proc. a. proc. b. proc. b. proc. a. portions. The precureus, or the space on the mesial side of the parietal lobe between the parieto-occipital





Dynamic Topic Models

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





Dynamic Topic Models: Summary

- Can model changes of topics (= word distributions) in corpora over time
- Uses HMM as a technique for modeling temporal influences
- As a by-product we have discussed techniques for comparing distributions

