

Probabilistic Topic Models and Applications

Eric Xing

Machine Learning Department
Language Technology Institute
Computer Science Department
Carnegie Mellon University

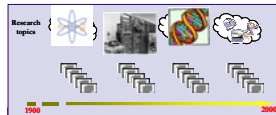


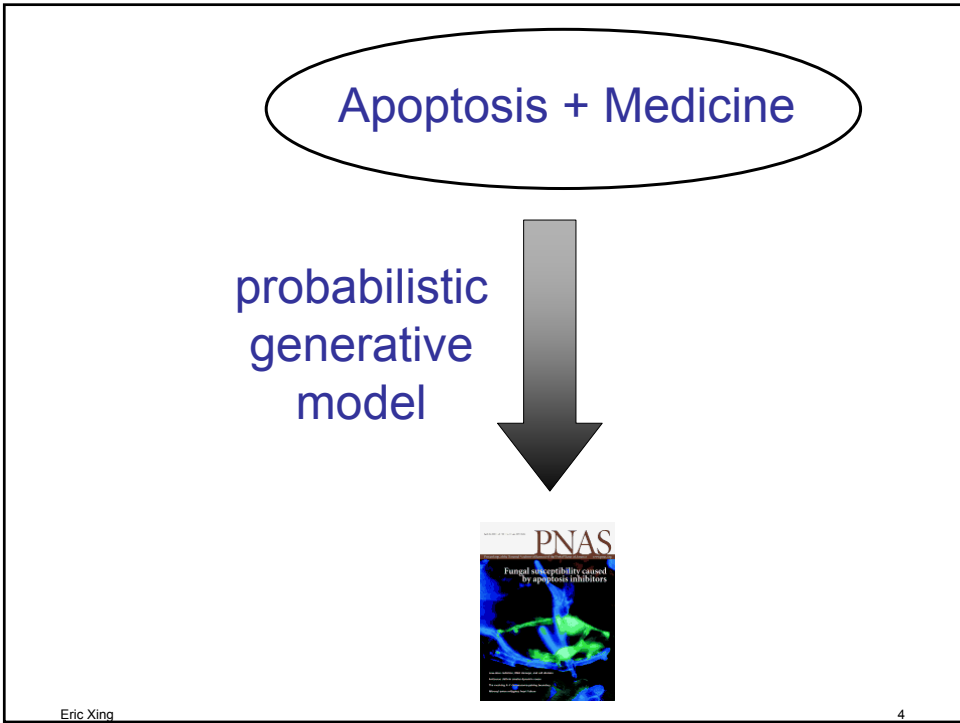
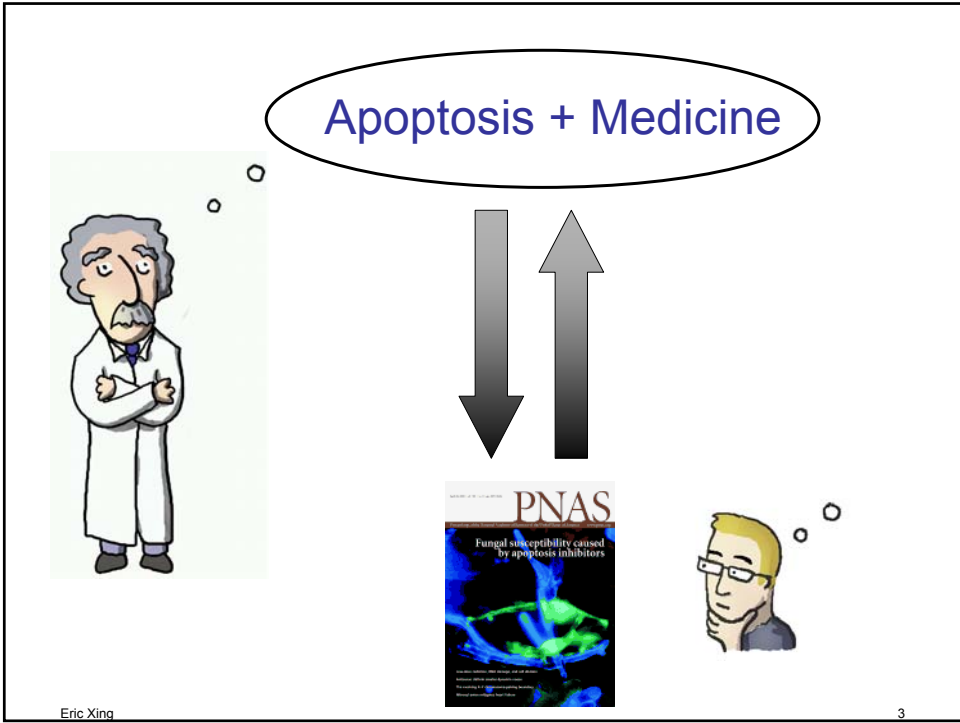
NLP and Data Mining

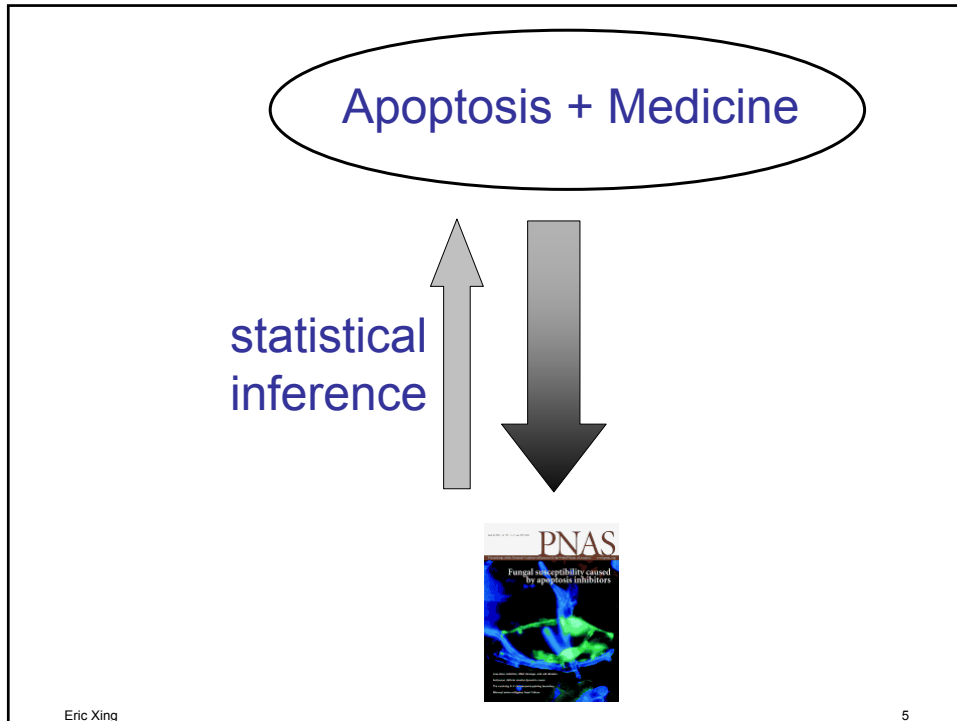


We want:

- Semantic-based search
- infer topics and categorize documents
- Multimedia inference
- Automatic translation
- Predict how topics evolve
- ...







This Talk

- A recap of graphical model
- Two families of probabilistic topics models and approximate inference
 - Bayesian admixture models
 - Random models
- Four applications
 - Topic evolution
 - Machine translation
 - Image topics
 - Multimedia inference

Eric Xing

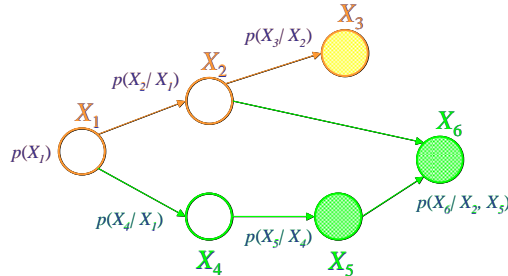
6

A decorative graphic in the top right corner of the slide, consisting of a grid of colored dots in shades of purple, blue, and green.

Probabilistic Graphical Models



- Graph-theoretic representations of probabilistic distributions



$$p(X_1, X_2, X_3, X_4, X_5, X_6) = p(X_1) p(X_2|X_1) p(X_3|X_2) p(X_4|X_1) p(X_5|X_4) p(X_6|X_2, X_5)$$

- Bayesian philosophy $\theta \rightarrow \text{grey} \Rightarrow \alpha \rightarrow \theta \rightarrow \text{grey}$

- Modular combination of heterogeneous parts -- **divide and conquer**

Eric Xing

7

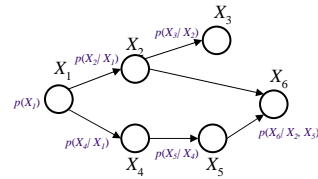
Probabilistic Inference



- Many modern problems in data mining/NLP can be formulated as probabilistic inference problems

$P(\text{query variable} \mid \text{query data \& KB})$

- Is this text document relevant to my query?
- Which category is this image in?
- What movies would I probably like?
- Create a caption for this image.
- Modeling document collections



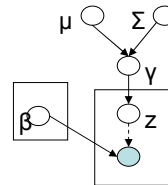
- General purpose algorithms exist to fully automate such computation
 - Computational cost depends on the topology of the network
 - Exact inference:
 - The junction tree algorithm
 - Approximate inference;
 - Loopy belief propagation, variational inference, Monte Carlo sampling

Eric Xing

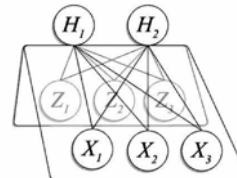
8

Two types of GMs

- **Directed edges** give **causality** relationships (**Bayesian Network** or **Directed Graphical Model**):



- **Undirected edges** simply give **correlations** between variables (**Markov Random Field** or **Undirected Graphical model**):



Eric Xing

9

This Talk

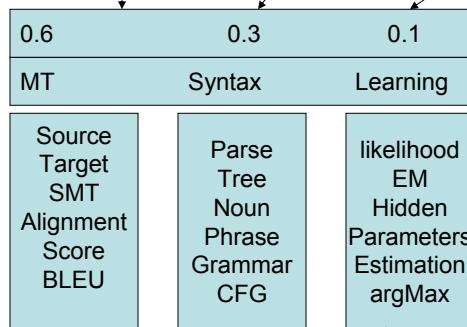
- A graphical model primer
- Two families of probabilistic topics models and approximate inference
 - Bayesian admixture models
 - Random models
- Three applications
 - Topic evolution
 - Machine translation
 - Image topics
 - Multimedia inference

Eric Xing

10

How to Model Semantic?

- Q: What is it about?
- A: Mainly MT, with syntax, some learning



Mixing Proportion

Topics

Unigram over vocabulary

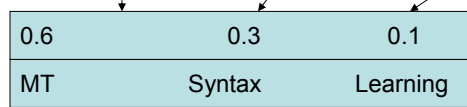
Topic Models

A Hierarchical Phrase-Based Model for Statistical Machine Translation

We present a statistical phrase-based Translation model that uses *hierarchical phrases*—phrases that contain sub-phrases. The model is formally a synchronous context-free grammar but is learned from a bitext without any syntactic information. Thus it can be seen as a shift to the *formal* machinery of syntax based translation systems without any *linguistic* commitment. In our experiments using BLEU as a metric, the hierarchical Phrase based model achieves a relative Improvement of 7.5% over Pharaoh, a state-of-the-art phrase-based system.

Why this is Useful?

- Q: What is it about?
- A: Mainly MT, with syntax, some learning



Mixing Proportion

- Q: give me similar document?
 - Structured way of browsing the collection
- Other tasks
 - Dimensionality reduction
 - TF-IDF vs. topic mixing proportion
 - Classification, clustering, and more ...

A Hierarchical Phrase-Based Model for Statistical Machine Translation

We present a statistical phrase-based Translation model that uses *hierarchical phrases*—phrases that contain sub-phrases. The model is formally a synchronous context-free grammar but is learned from a bitext without any syntactic information. Thus it can be seen as a shift to the *formal* machinery of syntax based translation systems without any *linguistic* commitment. In our experiments using BLEU as a metric, the hierarchical Phrase based model achieves a relative Improvement of 7.5% over Pharaoh, a state-of-the-art phrase-based system.

Words in Contexts



- “It was a nice **shot**. ”



Eric Xing

13

Words in Contexts (con'd)



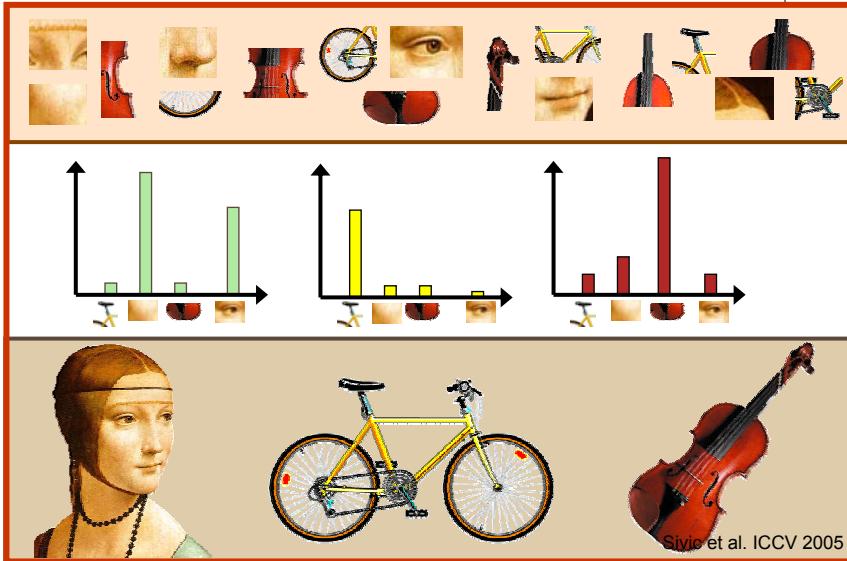
- the opposition Labor **Party** fared even worse, with a predicted 35 **seats**, seven less than last **election**.



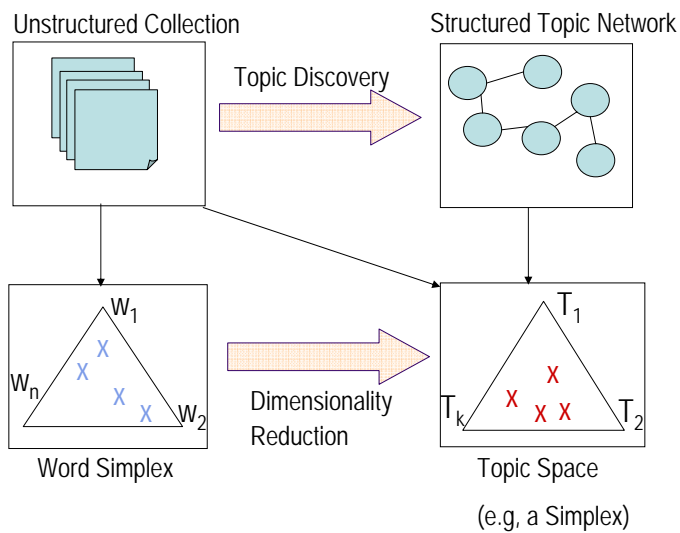
Eric Xing

14

"Words" in Contexts (con'd)



Topic Models: The Big Picture



Method One:



- **Hierarchical Bayesian Admixture**

A. Ahmed and E.P. Xing
AISTAT 2007

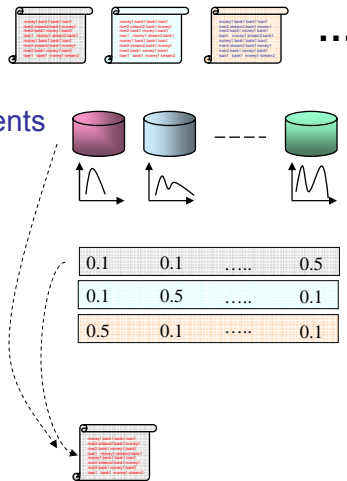
Eric Xing

17

Admixture Models



- Objects are **bags** of elements
- Mixtures are **distributions** over elements
- Objects have **mixing vector** θ
 - Represents each mixtures' contributions
- Object is **generated** as follows:
 - Pick a mixture component from θ
 - Pick an element from that component



Eric Xing

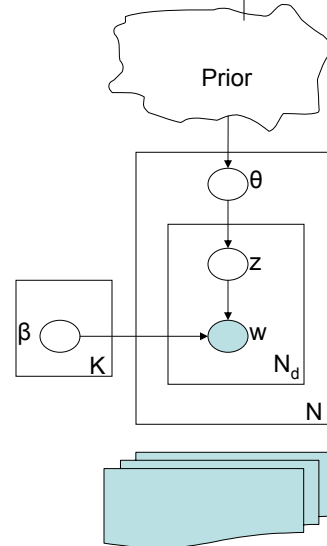
18

Topic Models = Admixture Models

Generating a document

- Draw θ from the prior
- For each word n
 - Draw z_n from *multinomial* $l(\theta)$
 - Draw w_n | $z_n, \{\beta_{1,k}\}$ from *multinomial* $l(\beta_{z_n})$

Which prior to use?

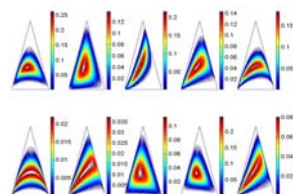
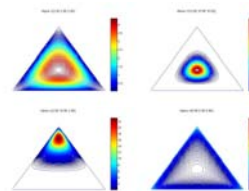


Eric Xing

19

Prior Comparison

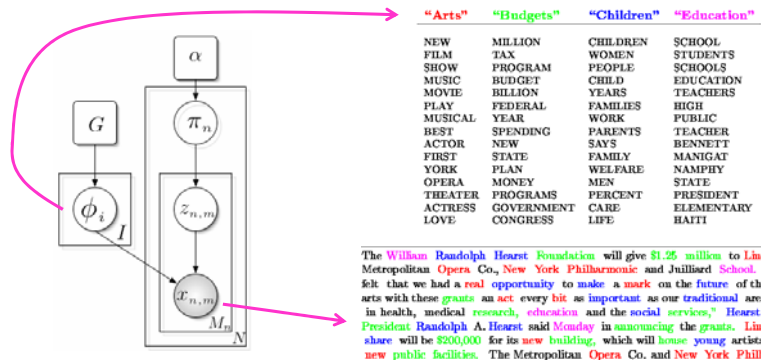
- Dirichlet (LDA) (Blei et al. 2003)
 - Conjugate prior means efficient inference
 - Can **only** capture variations in each topic's intensity **independently**
- Logistic Normal (CTM=LoNTAM) (Blei & Lafferty 2005, Ahmed & Xing 2006)
 - Capture the intuition that some topics are highly correlated and can rise up in intensity together
 - **Not** a conjugate prior implies **hard** inference



Eric Xing

20

Inference Tasks



The **William Randolph Hearst Foundation** will give **\$1.25 million** to **Lincoln Center**, **Metropolitan Opera Co.**, **New York Philharmonic** and **Julliard School**. "Our board felt that we had a **real opportunity** to make a **mark** on the future of the **performing arts** with these **grants** an **act** every bit as **important** as our **traditional areas of support** in **health**, **medical research**, **education** and the **social services**," **Hearst Foundation** President **Randolph A. Hearst** said **Monday** in **announcing** the **grants**. **Lincoln Center's** share will be **\$200,000** for its **new building**, which will **house young artists** and **provide new public facilities**. The **Metropolitan Opera Co.** and **New York Philharmonic** will receive **\$400,000** each. The **Julliard School**, where **music** and the **performing arts** are **taught**, will get **\$250,000**. The **Hearst Foundation**, a **leading supporter** of the **Lincoln Center Consolidated Corporate Fund**, will make its usual **annual \$100,000** donation, too.

Approximate Inference

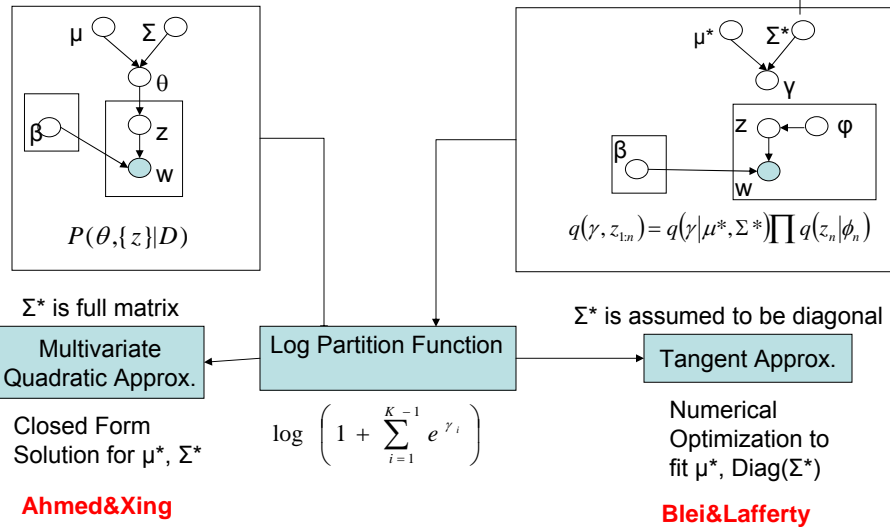


- Variational Inference
 - Mean field approximation (Blei et al)
 - Expectation propagation (Minka et al)
 - Variational 2nd-order Taylor approximation (Xing)

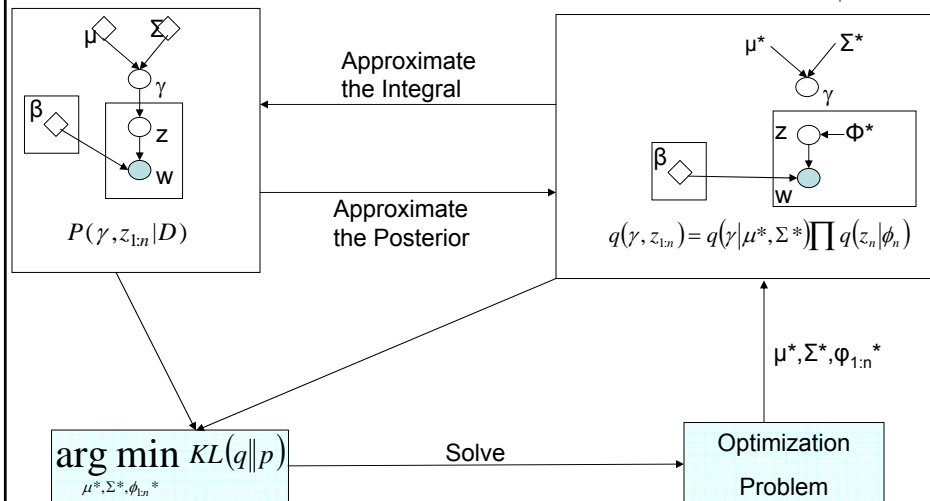
- Markov Chain Monte Carlo
 - Gibbs sampling (Griffiths et al)

Approximate Inference

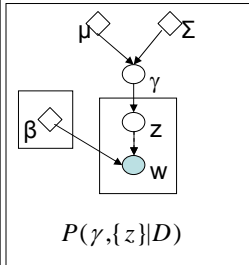
(e.g., MF, Jordan et al 1999, GMF, Xing et al 2004)



Variational Inference



Variational Inference With no Tears



Iterate until Convergence

- Pretend you know $E[Z_{1:n}]$
 - $P(\gamma | E[Z_{1:n}], \mu, \Sigma)$
- Now you know $E[\gamma]$
 - $P(z_{1:n} | E[\gamma], w_{1:n}, \beta_{1:k})$

- More Formally: $q^*(X_C) = P\left(X_C \left| \langle S_Y \rangle_{q_Y} : \forall y \in X_{MB}\right.\right)$

Message Passing Scheme (GMF)

Equivalent to previous method (Xing et. al.2003)

LoNTAM Variations Inference



- Fully Factored Distribution

$$q(\gamma, z_{1:n}) = q(\gamma) \prod q(z_n)$$

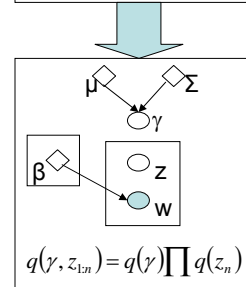
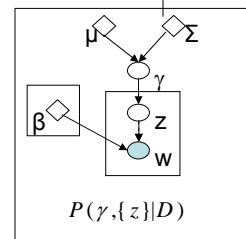
- Two clusters: λ and $Z_{1:n}$

$$q^*(X_C) = P\left(X_C \left| \langle S_Y \rangle_{q_Y} : \forall y \in X_{MB}\right.\right)$$

- Fixed Point Equations

$$q_Y^*(\gamma) = P\left(\gamma \left| \langle S_z \rangle_{q_z}, \mu, \Sigma\right.\right)$$

$$q_z^*(z) = P\left(z \left| \langle S_\gamma \rangle_{q_\gamma}, \beta_{1:k}\right.\right)$$



Variational γ

$$q_{\lambda}^*(\gamma) = P(\gamma | \langle S_z \rangle_{q_z}, \mu, \Sigma)$$

$$\propto P(\gamma | \mu, \Sigma) P(\langle S_z \rangle_{q_z} | \gamma)$$

Now what is $\langle S_z \rangle_{q_z}$?

$$S_z = m = \left[\sum_n I(z_n = 1), \dots, \sum_n I(z_n = k) \right]$$

$$\propto N(\gamma | \mu, \Sigma) \exp\{\langle m \rangle_{q_z} \gamma - N \times C(\gamma)\}$$

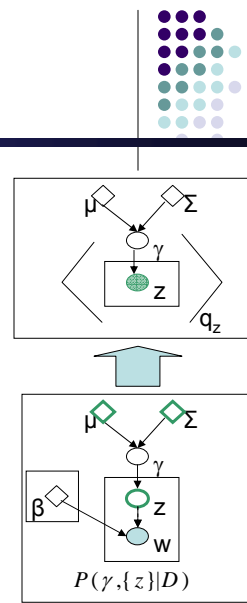
$$\propto \exp\left\{ -\frac{1}{2} \gamma' \Sigma^{-1} \gamma + \gamma \Sigma^{-1} \mu + \langle m \rangle_{q_z} \gamma - N \times C(\gamma) \right\}$$

$$C(\gamma) = C(\gamma_{\wedge}) + g'_{\lambda} (\gamma - \gamma_{\wedge}) + .5 (\lambda - \gamma_{\wedge})' H (\gamma - \gamma_{\wedge})$$

$$q_{\lambda}^*(\gamma) = N(\mu_{\gamma}, \Sigma_{\gamma})$$

$$\Sigma_{\gamma} = \text{inv}(\Sigma^{-1} + NH)$$

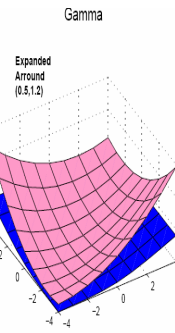
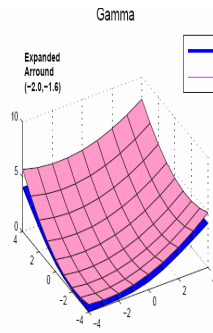
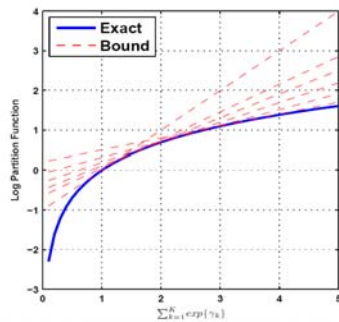
$$\mu_{\gamma} = \Sigma_{\gamma} (\Sigma^{-1} \mu + NH \gamma_{\wedge} + \langle m \rangle - Ng)$$



Eric Xing

27

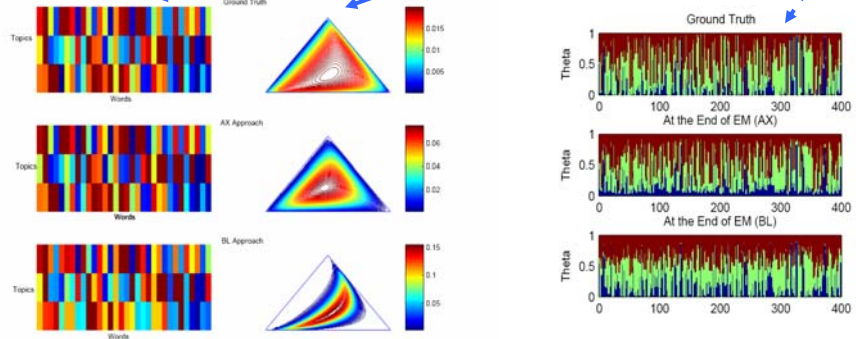
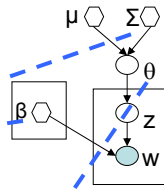
Tangent Approximation



Eric Xing

28

Test on Synthetic Text



Eric Xing

29

Comparison: accuracy and speed

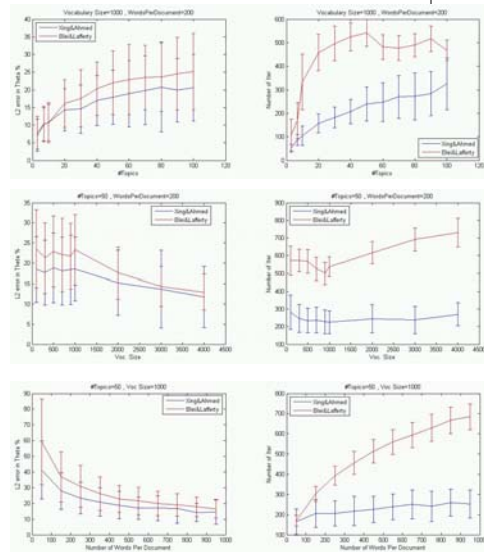


L2 error in topic vector est. and # of iterations

- Varying Num. of Topics

- Varying Voc. Size

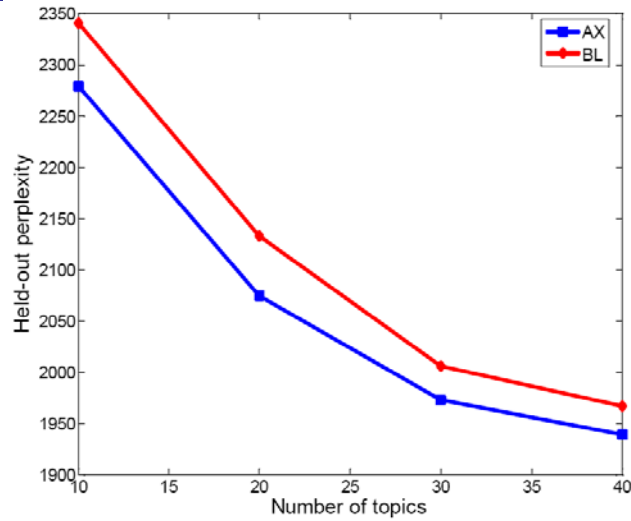
- Varying Num. Words Per Document



Eric Xing

30

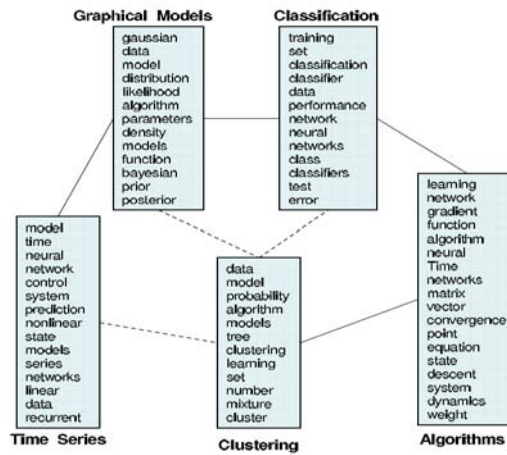
Comparison: perplexity



Eric Xing

31

Topics and topic graphs



Eric Xing

32

Result on PNAS collection



- PNAS abstracts from 1997-2002
 - 2500 documents
 - Average of 170 words per document
- Fitted 40-topics model using both approaches
- Use low dimensional representation to predict the abstract category
 - Use SVM classifier
 - 85% for training and 15% for testing

Classification Accuracy

Category	Doc	BL	AX
Genetics	21	61.9	61.9
Biochemistry	86	65.1	77.9
Immunology	24	70.8	66.6
Biophysics	15	53.3	66.6
Total	146	64.3	72.6

Eric Xing

33

Method Two:



- Layered Boltzmann machines

E.P. Xing, R. Yan and A. G. Hauptmann,
UAI 2006

Eric Xing

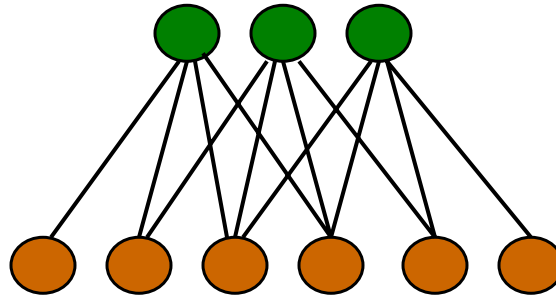
34

The Harmonium



hidden units

visible units



Boltzmann machines:

$$p(x, h | \theta) = \exp \left\{ \sum_i \theta_i \phi_i(x_i) + \sum_j \theta_j \phi_j(h_j) + \sum_{i,j} \theta_{i,j} \phi_{i,j}(x_i, h_j) - A(\theta) \right\}$$

Eric Xing

35

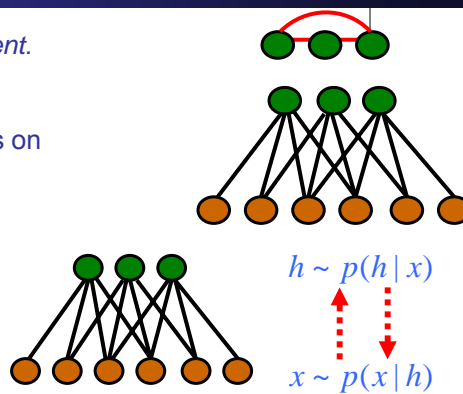
Properties of Harmoniums



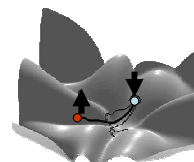
- Factors are marginally *dependent*.
- Factors are conditionally *independent* given observations on the visible nodes.

$$P(\ell | \mathbf{w}) = \prod_i P(\ell_i | \mathbf{w})$$

- Iterative Gibbs sampling.



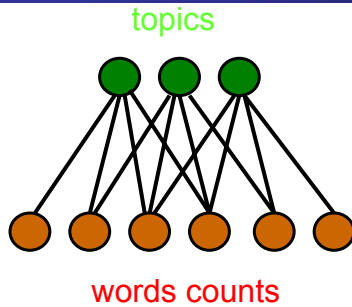
- Learning with contrastive divergence



Eric Xing

36

A Binomial Word-count Model



$h_j = 3$: topic j has strength 3

$$h_j \in \mathbf{R}, \quad \langle h_j \rangle = \sum_i W_{i,j} x_i$$

$x_i = n$: word i has count n

$$x_i \in \mathbf{I}$$

$$p(\mathbf{h} | \mathbf{x}) = \prod_j \text{Normal}_{h_j} \left[\sum_i \vec{W}_{ij} \vec{x}_i, 1 \right]$$

$$p(\mathbf{x} | \mathbf{h}) = \prod_i \text{Bi}_{x_i} \left[N, \frac{\exp(\alpha_j + \sum_j W_{ij} h_j)}{1 + \exp(\alpha_j + \sum_j W_{ij} h_j)} \right]$$

$$\text{Bi}_{x_i} [N, p] = C_{x_i}^N p^{x_i} (1-p)^{N-x_i} = C_{x_i}^N \left(\frac{p}{1-p}\right)^{x_i} (1-p)^N$$

$$\text{Let } p = \frac{\exp(\alpha_j + \sum_j W_{ij} h_j)}{1 + \exp(\alpha_j + \sum_j W_{ij} h_j)}$$

$$\text{Bi}_{x_i} [N, p] = C_{x_i}^N \frac{(\exp(\alpha_j + \sum_j W_{ij} h_j))^{x_i}}{(1 + \exp(\alpha_j + \sum_j W_{ij} h_j))^N}$$

$$\propto C_{x_i}^N \exp\left\{(\alpha_j + \sum_j W_{ij} h_j)x_i + A_j\right\}$$

Reduce to softmax when $N=1$!

$$\Rightarrow p(\mathbf{x}) \propto \exp\left\{\left(\sum_i \alpha_i x_i - \log \Gamma(x_i) - \log \Gamma(N - x_i)\right) + \frac{1}{2} \sum_j \left(\sum_i W_{i,j} x_i\right)^2\right\}$$

Eric Xing

37

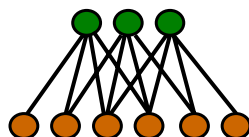
The Computational Trade-off



Undirected model: Learning is hard, inference is easy.

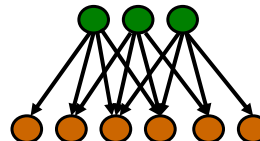
Directed Model: Learning is "easier", inference is hard.

Example: Document Retrieval.



topics

words



Retrieval is based on comparing (posterior) topic distributions of documents.

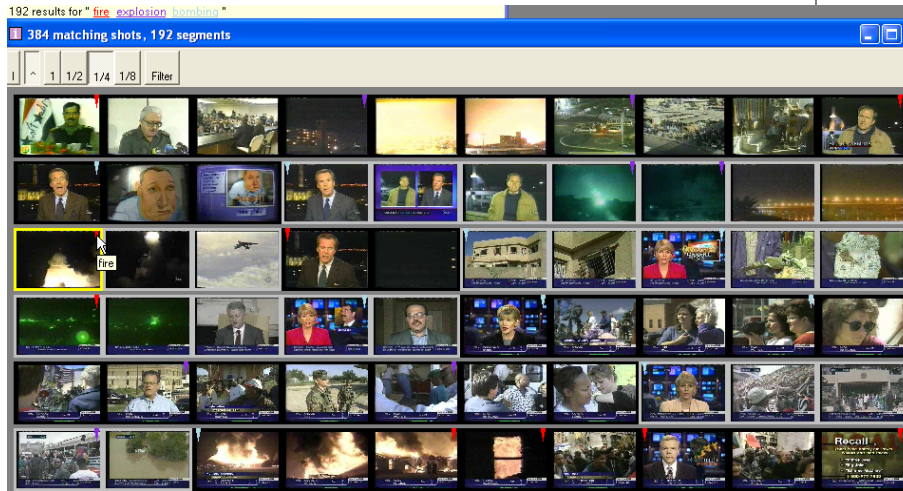
- **directed models:** inference is slow. Learning is relatively "easy".

- **undirected model:** inference is fast. Learning is slow but can be done offline.

Eric Xing

38

Multi-Source Data

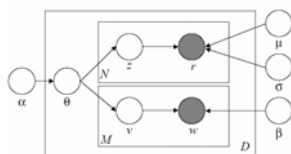


TRECVID 2004 Example Images

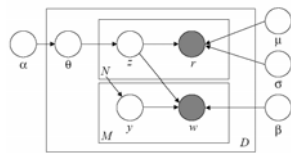
Eric Xing

39

Inter-Source Associations

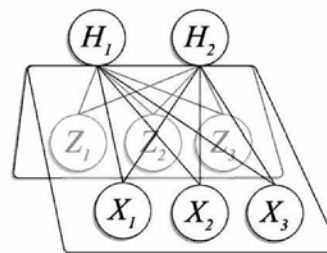


GM-LDA



Co-LDA

DWH

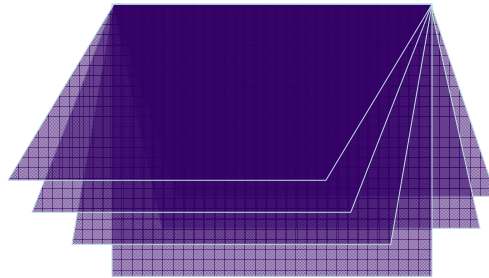


Z and X are marginally dependent (same as GM-LDA)

Eric Xing

40

Multi-wing Harmoniums



Learning and Inference



- Maximal likelihood learning based on gradient ascent.

$$\delta\theta_i \propto \langle f_i(x_i) \rangle_{\text{data}} - \langle f_i(x_i) \rangle_p$$

- gradient computation requires model distribution $p(\cdot)$
- $p(\cdot)$ is intractable
- Contrastive Divergence
 - approximate $p(\cdot)$ with Gibbs sampling
- Variational approximation
 - GMF approximation

$$q(\mathbf{x}, \mathbf{z}, \mathbf{h}) = \prod_i q(x_i | v_i) \prod_k q(z_k | \mu_k, \sigma_k) \prod_j q(h_j | \gamma_j)$$

Inter-source Inference



- GMF approximation to DWH

$$q(\mathbf{x}, \mathbf{z}, \mathbf{h}) = \prod_i q(x_i | N, v_i) \prod_k q(z_k | \mu_k, \sigma_k) \prod_j q(h_j | \gamma_j)$$

- Expected mean value of topic strength:

$$\gamma_j = \sum_i W_{i,j} v_i + \sum_k U_{k,j} \mu_k$$

- Expected mean value of image-feature :






$$\mu_k = \sigma_k^2 \left(\beta_k + \sum_j U_{k,j} \gamma_j \right)$$

- Expected mean count

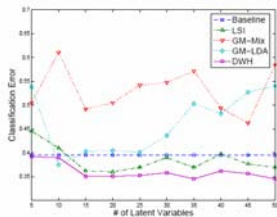
$$N v_i = N \frac{\exp(\alpha_j + \sum_j W_{ij} \gamma_j)}{1 + \exp(\alpha_j + \sum_j W_{ij} \gamma_j)}$$

Examples of Latent Topics

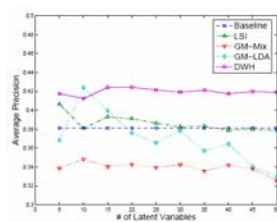


T_1	storms gulf hawaii low forecast southeast showers 
T_2	rebounds 14 shouting tests guard cut hawks 
T_3	engine flying craft asteroid say hour aerodynamic 
T_4	safe cross red sure dry providing services 
T_5	losing jersey sixth antonio david york orlando 

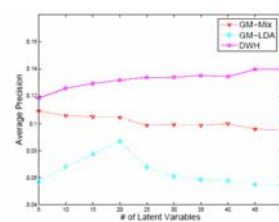
Performance



Classification



Retrieval



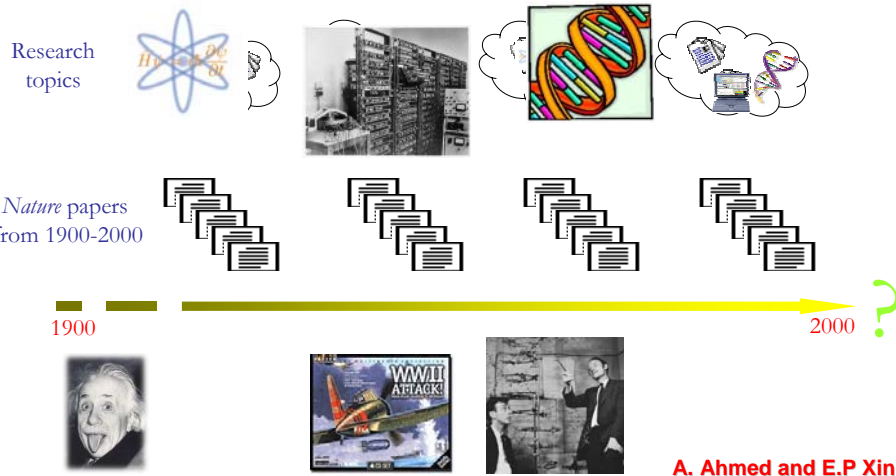
Annotation

This Talk



- A graphical model primer
- Two families of probabilistic topics models and approximate inference
 - Bayesian admixture models
 - Random models
- Three applications
 - topic evolution
 - Machine translation
 - Image topics
 - Multimedia inference

Application 1: topic evolution?



Eric Xing

47

How to Model Topic Evolution



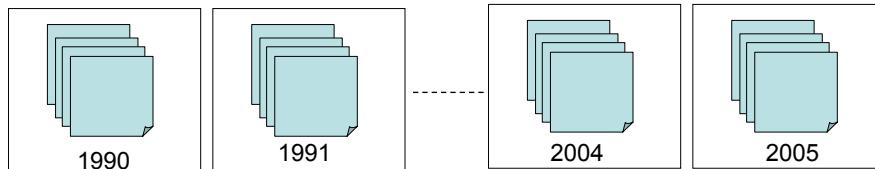
Topic Trends

Topic Keywords

Topic correlations

~~Number of topics~~

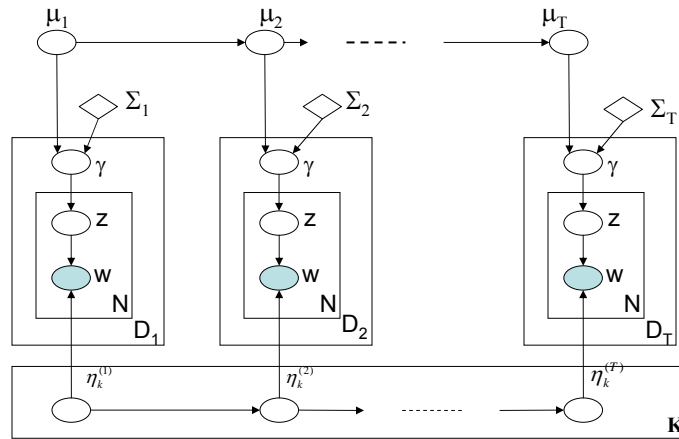
The Dynamic Correlated Topic model



Eric Xing

48

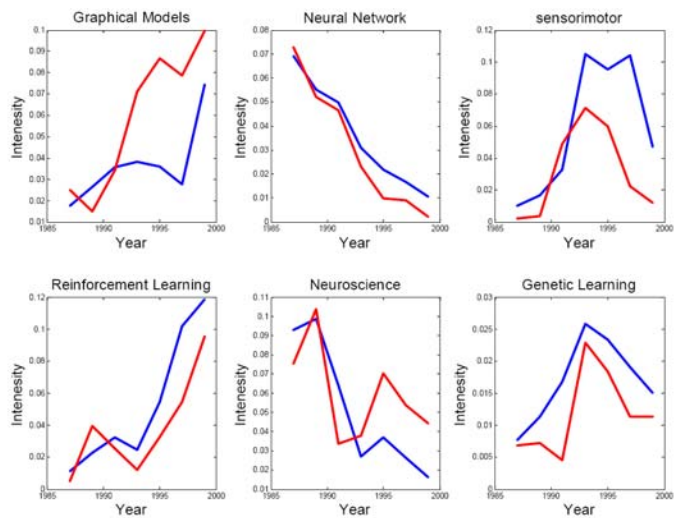
The Dynamic CTM



Eric Xing

49

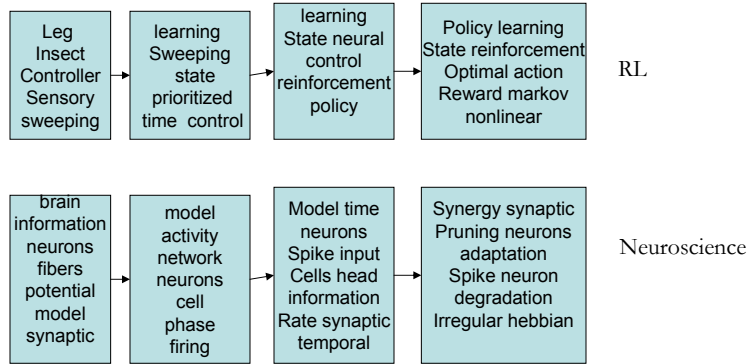
Topic Trends



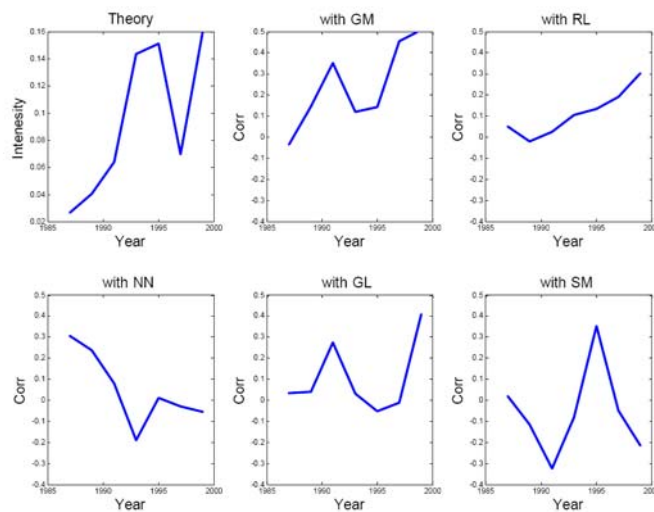
Eric Xing

50

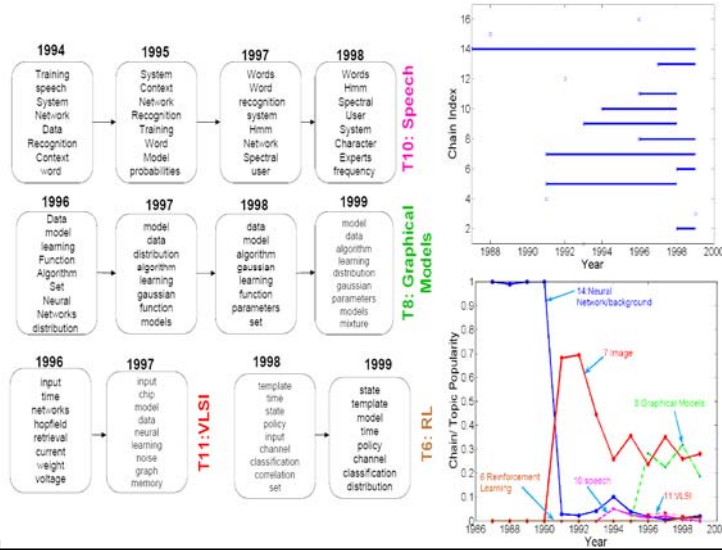
Topic Words over Time



Topic Correlations Over Time



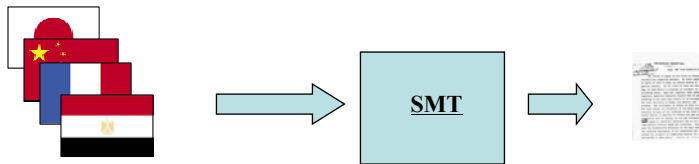
Birth-Death of Topics via Dynamic Dirichlet Process



Eric Xing

53

Application 2: Machine translation



Eric Xing

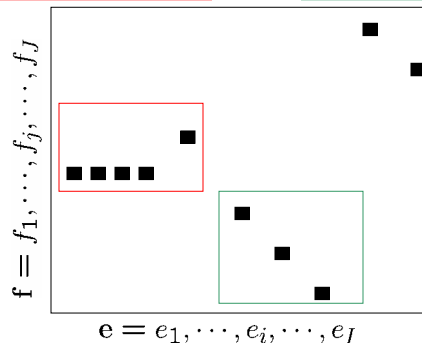
B. Zhao and E.P Xing, ACL 2006, NIPS 2007

54

Word Alignment

天津 与 俄罗斯 经贸 关系 稳步 发展

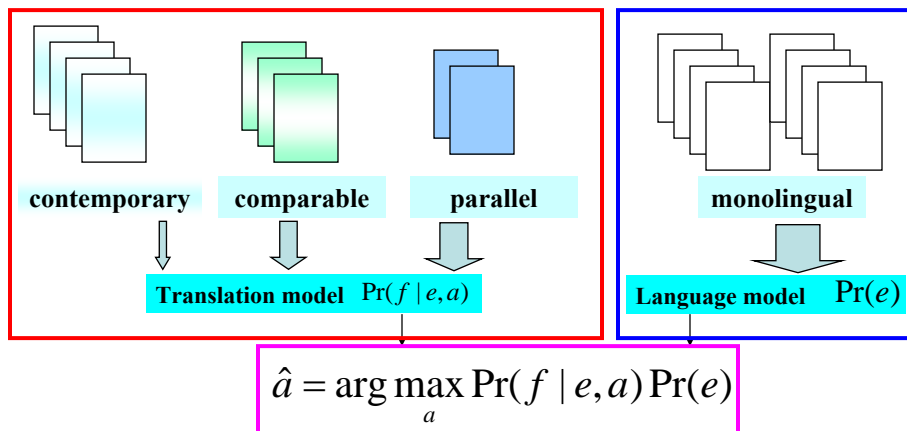
The economy and trade relations between russia and tianjin develop steadily



Eric Xing

55

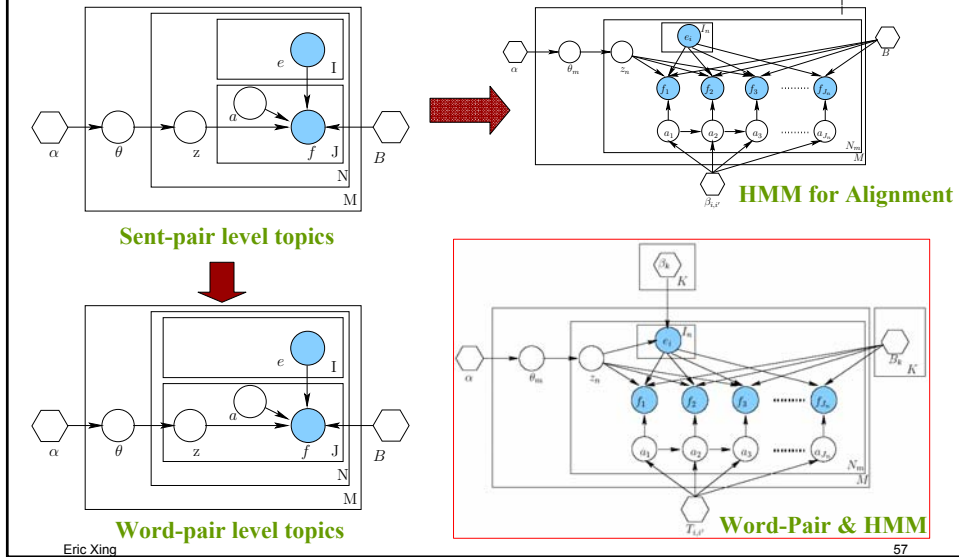
The Statistical Formulation



Eric Xing

56

An upgrade path for BiTAMs



Experiments

- Training data
 - Small: Treebank 316 doc-pairs (133K English words)
 - Large: FBIS-Beijing, Sinorama, XinHuaNews, (15M English words).

Train	#Doc.	#Sent.	#Tokens	
			English	Chinese
Treebank	316	4172	133K	105K
FBIS.BJ	6,111	105K	4.18M	3.54M
Sinorama	2,373	103K	3.81M	3.60M
XinHua	19,140	115K	3.85M	3.93M
FOUO	15,478	368K	13.14M	11.93M
Test	95	627	25,500	19,726

- Word Alignment Accuracy & Translation Quality
 - F-measure
 - BLEU

Topics



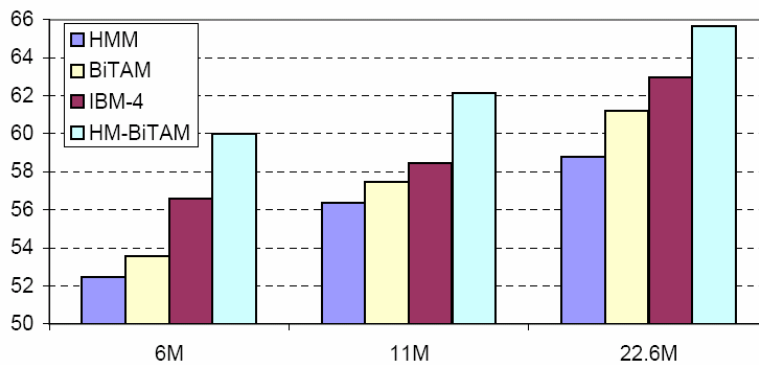
T1	Teams, sports, disabled, games members, people, cause, water, national, handicapped
T2	Shenzhen, singapore, hongkong, stock, national, investment, yuan, options, million, dollar
T3	Chongqing, company, takeover, shenzhen, tianjin, city, national, government, project, companies
T4	Hongkong, trade, export, import, foreign, tech., high, 1998, year, technology
T5	House, construction, government, employee, living, provinces, macau, anhui, yuan
T6	Gas, company, energy, usa, russia, france, chongqing, resource, china, economy, oil

T1	人, 残疾, 体育, 事业, 水, 世界, 区, 新华社, 队员, 记者
T2	深圳, 深, 新, 元, 有, 股, 香港, 国有, 外资, 新华社
T3	国家, 重庆, 市, 区, 厂, 天津, 政府, 项目, 国, 深圳
T4	香港, 贸易, 出口, 外资, 合作, 今年, 项目, 利用, 新, 技术
T5	住房, 房, 九江, 建设, 澳门, 元, 职工, 目前, 国家, 占, 省
T6	公司, 天然气, 两, 国, 美国, 记者, 关系, 俄, 法, 重庆

Eric Xing

59

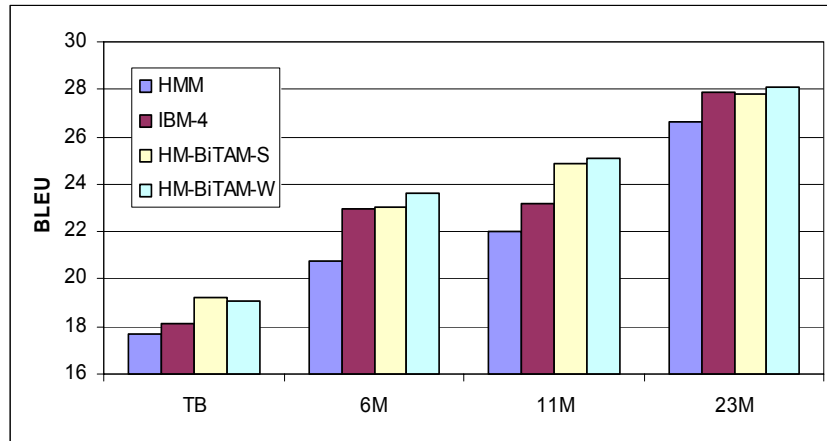
HM-BiTAM versus others



Eric Xing

60

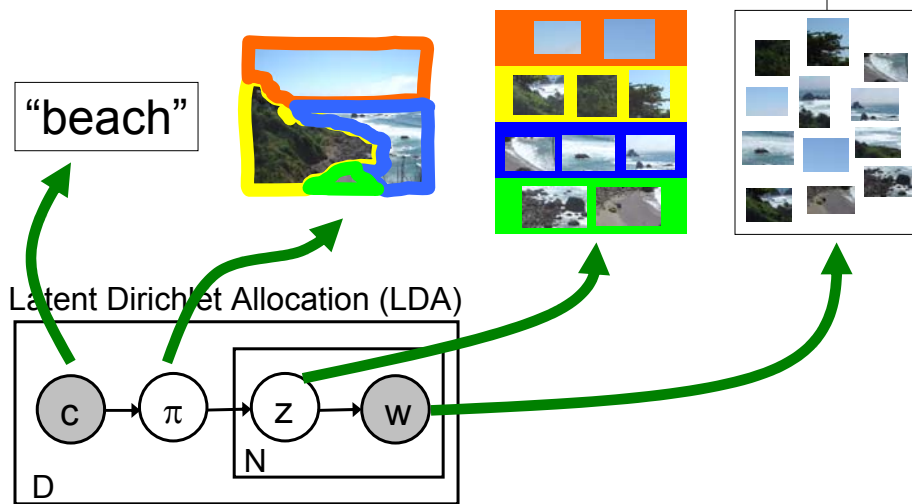
Translation Evaluations



Eric Xing

61

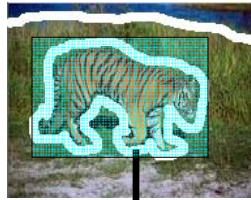
Application 3: Topic Models for Images



Eric Xing

Fei-Fei et al. ICCV 2005

Image Representation



cat, grass, tiger, water

$$[r_{11} \cdots r_{1d}], [w_1 \cdots w_{|V|}]$$

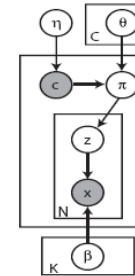
representation vector : **annotation vector**
 (real, 1 per image segment) : (binary, same for each segment)

$$[r_{n1} \cdots r_{nd}], [w_1 \cdots w_{|V|}]$$

Eric Xing

63

To Generate an Image ...



$$p(x, z, \pi, c | \theta, \eta, \beta) = p(c | \eta) p(\pi | c, \theta) \cdot \prod_{n=1}^N p(z_n | \pi) p(x_n | z_n, \beta)$$

$$p(c | \eta) = \text{Mult}(c | \eta)$$

$$p(\pi | c, \theta) = \prod_{j=1}^C \text{Dir}(\pi | \theta_{j \cdot})^{\delta(c, j)}$$

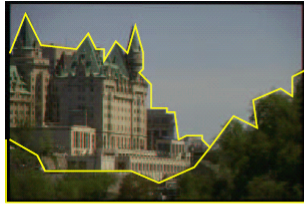
$$p(z_n | \pi) = \text{Mult}(z_n | \pi)$$

$$p(x_n | z_n, \beta) = \prod_{k=1}^K p(x_n | \beta_{k \cdot})^{\delta(z_n, k)}$$

eric xing

64

Annotated images



{9.32, 2.44, 0.02, 3.23}
{4.35, 3.12, -0.23, 9.41}
{6.65, 2.11, 1.02, 2.31}

This cozy place is nestled in the heart of the Mission. Easy access to bars, restaurants, and BART.

This, cozy, place, is, nestled, in, the, heart, of, the, Mission. Easy, access, to, bars, restaurants, and, BART

- Forsyth et. al. (2001): images as documents where region-specific feature vectors are like visual words.
- A captioned image can be thought of as annotated data: two documents, one of which describes the other.

Eric Xing

65

Application 4: video representation/classification



- Video: a complex, multi-modal data type for representation and classification
 - Image, text (closed-captions, speech transcript), audio
- Goal: classify video segments called video shots into semantic categories



anchor



building



meeting



speech

J. Yang, Y. Liu, E. P. Xing and A. Hauptmann,
SDM 2007, **BEST PAPER Award**

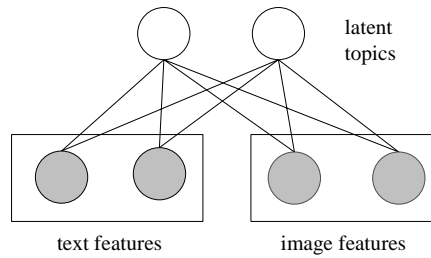
Eric Xing

66

Harmoniums for Multi-modal Data



- Dual-wing harmoniums (DWH) [Xing et al. 05]
 - modeling bi-modal data: captioned images, video
 - learning hidden topics from two "wings" of observed features



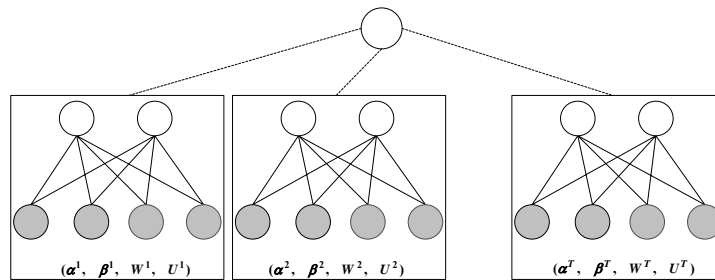
Eric Xing

67

Mixture-of-Harmoniums (MoH)



- A family of category-specific dual-wing harmoniums



- classification by finding the "best-fitting" harmonium

Eric Xing

68

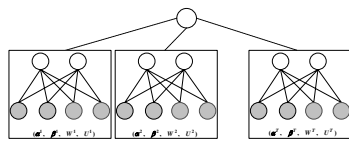
$H_1 \dots$

$\dots X_N$

Semantic Topics by FoH



- Revealing “sub-topics” of each category
- Co-clusters of both text and image features



Topic 1 life, call, way, fire, know, thousands, rain, farmers, control

Topic 2 space, flight, thousands, fifteen, Florida, radar, track, amount

Topic 3 asteroid, scientists, destroy, miss, destruction, actually, come, course

Topic 4 rain, control, area, forest, years, fires, large, burning, state, nature

Topic 5 panic, sized, type, headaches, freedom, love, turning, beautiful

Eric Xing

69

Conclusion



- GM-based topic models are cool
 - Flexible
 - Modular
 - Interactive
- There are many ways of implementing topic models
 - Directed
 - Undirected
- Efficient Inference/learning algorithms
 - GMF, with Laplace approx. for non-conjugate dist.
 - MCMC
- Many applications
 - ...
 - Word-sense disambiguation
 - Word-net
 - Network inference

$$H_1 \cdots H_k$$

$$X_1 \cdots X_N \quad Z_1 \cdots Z_N$$

category T

Eric Xing

70