

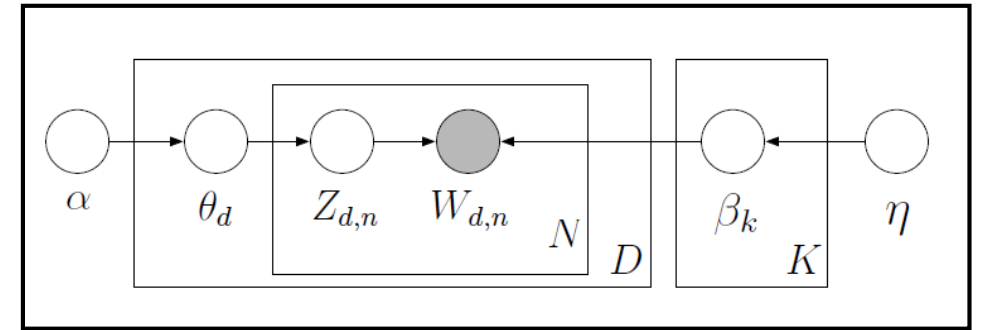
Correlated Topic Models

Authors: Blei and Lafferty, 2006

Reviewer: Casey Hanson

Recap Latent Dirichlet Allocation

- $D \equiv$ set of documents.
- $K =$ set of topics.
- $V =$ set of all words. $|N|$ words in each doc.
- $\theta_d \equiv$ Multi over topics for a document $d \in D$. $\theta_d \sim Dir(\alpha)$
- $\beta_k \equiv$ Multi over words in a topic, $k \in K$. $\beta_k \sim Dir(\eta)$
- $Z_{d,n} \equiv$ topic selected for word n in document d . $Z_{d,n} \sim Multi(\theta_d)$
- $W_{d,n} \equiv n_{th}$ word in document d . $W_{d,n} \sim Multi(B_{Z_{d,n}})$



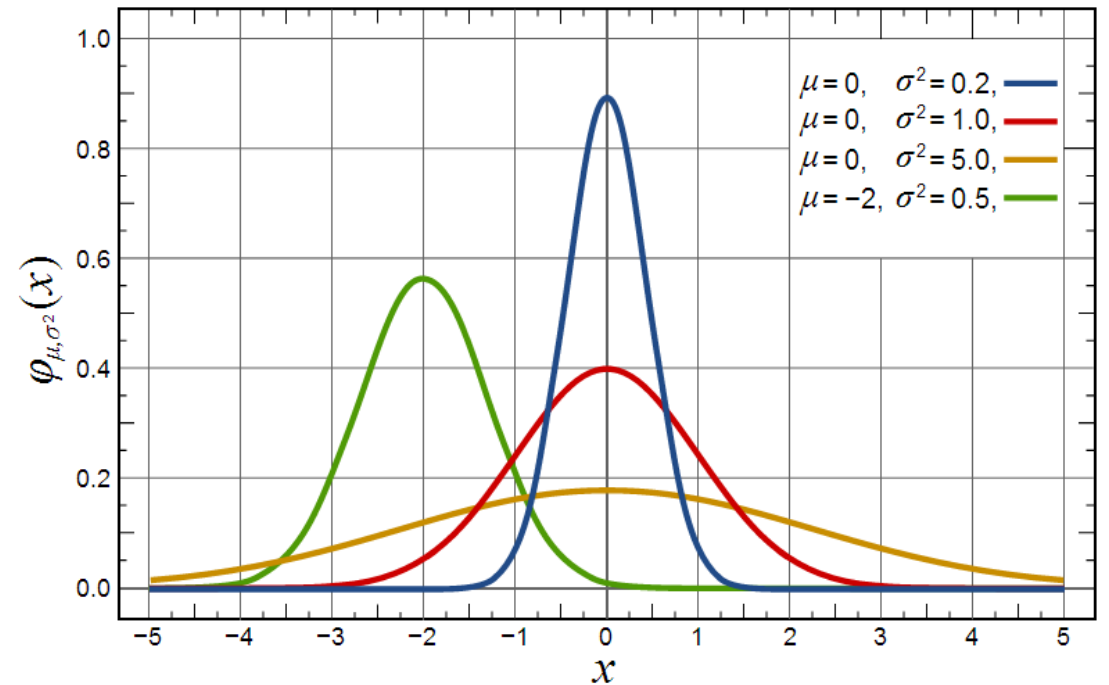
Latent Dirichlet Allocation

- Need to calculate posterior: $P(\theta_{1:D}, Z_{1:D,1:N}, \beta_{1:K} | W_{1:D,1:N}, \alpha, \eta)$
 - $\propto p(\theta_{1:D}, Z_{1:D,1:N}, \beta_{1:K}, W_{1:D,1:N}, \alpha, \eta)$
 - Normalization factor, $\int_{\beta} \int_{\theta} \sum_Z p(\dots)$, is intractable
 - Need to use approximate inference.
 - Gibbs Sampling
- Drawback
 - No intuitive relationship between topics.
- Challenge
 - Develop method similar to LDA with relationships between topics.

Normal or Gaussian Distribution

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- Continuous distribution
 - Symmetrical and defined for $-\infty < x < \infty$
- Parameters: $\mathcal{N}(\mu, \sigma^2)$
 - $\mu \equiv$ mean
 - $\sigma^2 \equiv$ variance
 - $\sigma \equiv$ standard deviation
- Estimation from Data: $X = \{x_1 \dots x_n\}$
 - $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$
 - $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$



Multivariate Gaussian Distribution: k dimensions

$$f(\mathbf{X}) = f_x(X_1 \dots X_k) = \frac{1}{(2\pi)^{k/2} \sqrt{\det \Sigma}} e^{-\frac{1}{2}(\mathbf{X}-\boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{X}-\boldsymbol{\mu})}$$

- $\mathbf{X} = [X_1 \dots X_k]^T \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$
 - $\boldsymbol{\mu} \equiv k \times 1$ vector of means for each dimension
 - $\Sigma \equiv k \times k$ covariance matrix.

Example: 2D Case

- $\boldsymbol{\mu} = E[\mathbf{X}] = \begin{bmatrix} E[x_1] \\ E[x_2] \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$
- $\Sigma = \begin{bmatrix} E[(x_1 - \mu_1)^2] & E[(x_1 - \mu_1)(x_2 - \mu_2)] \\ E[(x_1 - \mu_1)(x_2 - \mu_2)] & E[(x_2 - \mu_2)^2] \end{bmatrix}$

2D Multivariate Gaussian:

$$\bullet \Sigma = \begin{bmatrix} \sigma_{X_1}^2 & \rho_{X_1, X_2} \sigma_{X_1} \sigma_{X_2} \\ \rho_{X_1, X_2} \sigma_{X_1} \sigma_{X_2} & \sigma_{X_2}^2 \end{bmatrix}$$

- Topic Correlations on Off Diagonal

$$\bullet \rho_{X_1, X_2} \sigma_{X_1} \sigma_{X_2} = E[(x_1 - \mu_1)(x_2 - \mu_2)] = \sum_{i=1}^n \frac{(X_{i,1} - \mu_1)(X_{i,2} - \mu_2)}{n}$$

- Covariance matrix is diagonal!

Matlab Demo

...Back to Topic Models

- How can we adapt LDA to have correlations between topics.
- In LDA, we assume two things:
 - **Assumption 1:** Topics in a document are independent. $\theta_d \sim \text{Dir}(\alpha)$
 - **Assumption 2:** Distribution of words in a topic is stationary. $B_k \sim (\eta)$
- To sample topic distributions for topics that are correlated, we need to correct assumption 1.

Exponential Family of Distributions

- Family of distributions that can be placed in the following form:

$$f(x|\theta) = h(x) \cdot e^{\eta(\theta) \cdot T(x) - A(\theta)}$$

- Ex: Binomial distribution: $\theta = p$

$$f(x|\theta) = \binom{n}{x} p^x (1-p)^{n-x}, x \in \{0, 1, 2, \dots, n\}$$

- $\eta(\theta) = \log \frac{p}{1-p}$

$$h(x) = \binom{n}{x}, \quad A(\theta) = n \log 1-p, \quad T(x) = x$$

$$f(x) = \binom{n}{x} e^{x \cdot \log\left(\frac{p}{1-p}\right) + n \cdot \log(1-p)}$$

Natural Parameterization

Categorical Distribution

- Multinomial $n=1$:
 - $f(x_1) = \theta_1$; $f(Z_1) = \theta^T \cdot Z_1$
 - where $Z_1 = [1 \ 0 \ 0 \dots 0]^T$ (**Iverson Bracket or Indicator Vector**)
 - $z_i = 1$
- Parameters: θ
 - $\theta = [p_1 \ p_2 \ p_3]$, where $\sum_i p_i = 1$
 - $\theta' = \begin{bmatrix} p_1 & p_2 & 1 \\ p_k & p_k & \end{bmatrix}$
 - $\log \theta' = \begin{bmatrix} \log \frac{p_1}{p_k} & \log \frac{p_2}{p_k} & 1 \end{bmatrix}$

Exponential Family Multinomial With N=1

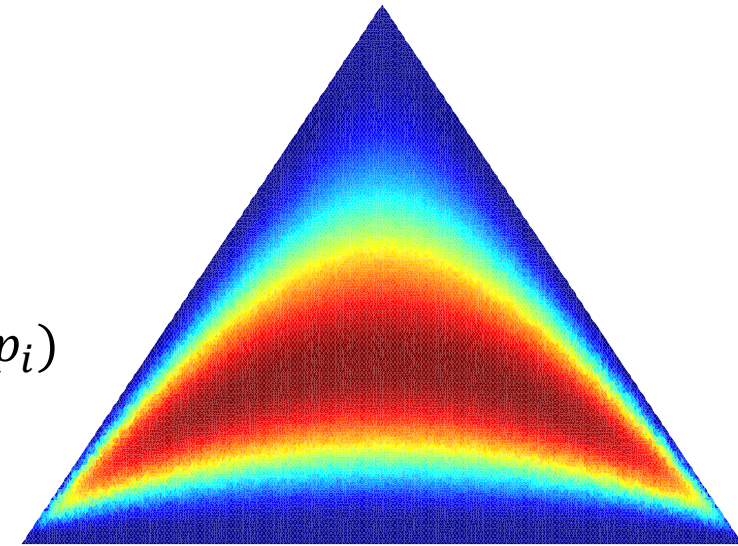
- **Recall:** $f(Z_i|\theta) = \theta^T \cdot Z_i$
- We want: $f(x|\theta) = h(x) \cdot e^{\eta(\theta) \cdot T(x) - A(\theta)}$
- $f(Z_i|\eta) = e^{\eta^T Z_i - \log \sum_{i=1} e^{\eta_i}} = \frac{e^{\eta^T \cdot Z_i}}{\sum_{i=1} e^{\eta_i}}$
- *Note: k-1 independent dimensions in Multinomial*
- $\eta' = [\log \frac{p_1}{p_k} \log \frac{p_2}{p_k} \dots 0]$, $\eta'_i = \log \frac{p_i}{p_k}$
- $f(Z_i|\eta') = \frac{e^{\eta'^T \cdot Z_i}}{1 + \sum_{i=1}^{k-1} e^{\eta'_i}}$

Verify: Classroom participation

- Given: $\eta = [\log \frac{p_1}{p_k} \log \frac{p_2}{p_k} \dots 0]$
- Show: $f(Z_i | \theta) = \theta^T \cdot Z_i = e^{\eta^T Z_i - \log \sum_{i=1} e^{\eta_i}}$

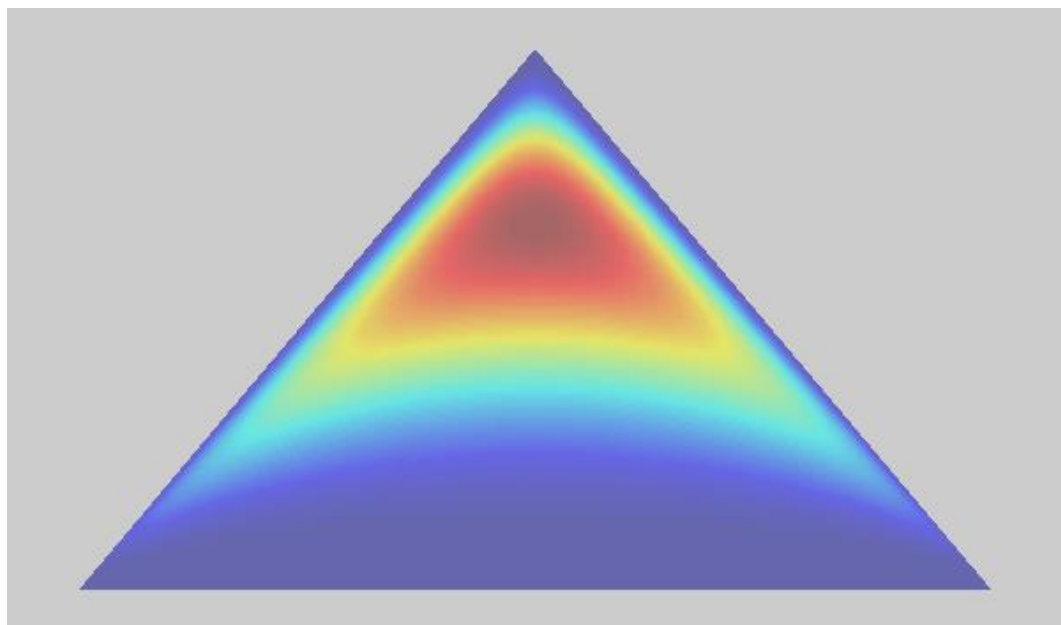
Intuition and Demo

- Can sample η from any number of places.
 - Choose normal (allows for correlation between topic dimensions)
- Get a topic distribution for each document by sampling:
 $\eta \sim \mathcal{N}_{k-1}(\mu, \sigma)$
 - What is the μ
 - Expected deviation from last topic: $\log\left(\frac{p_i}{p_k}\right)$
 - Negative means push density towards last topic ($\eta_i < 0, p_k > p_i$)
 - What about the covariance
 - Shows variability in deviation from last topic between topics.

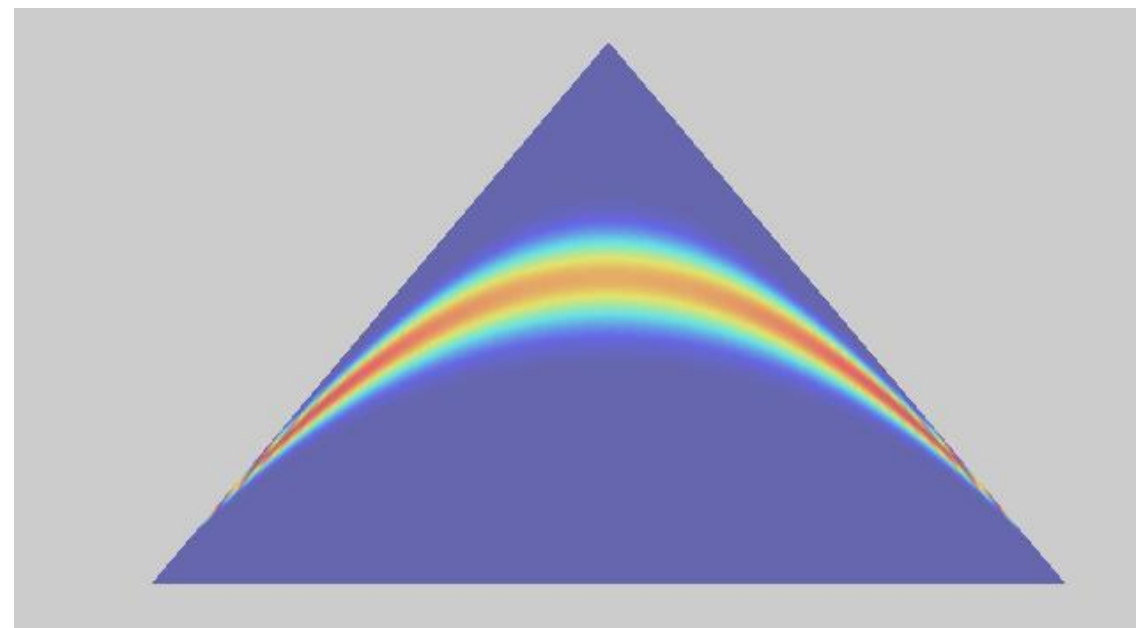


$$\mu = [0 \ 0]^T, \sigma = [1 \ 0; 0 \ 1]$$

Favoring Topic 3

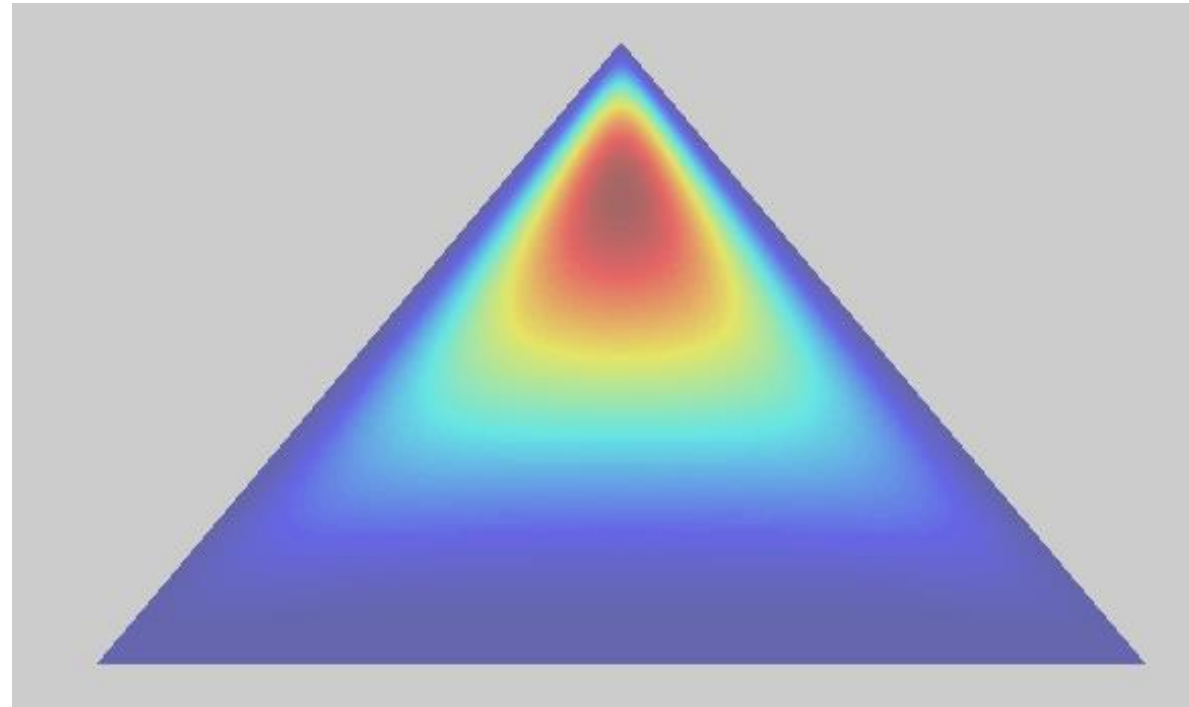


$$\mu = [-0.9, -0.9], \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



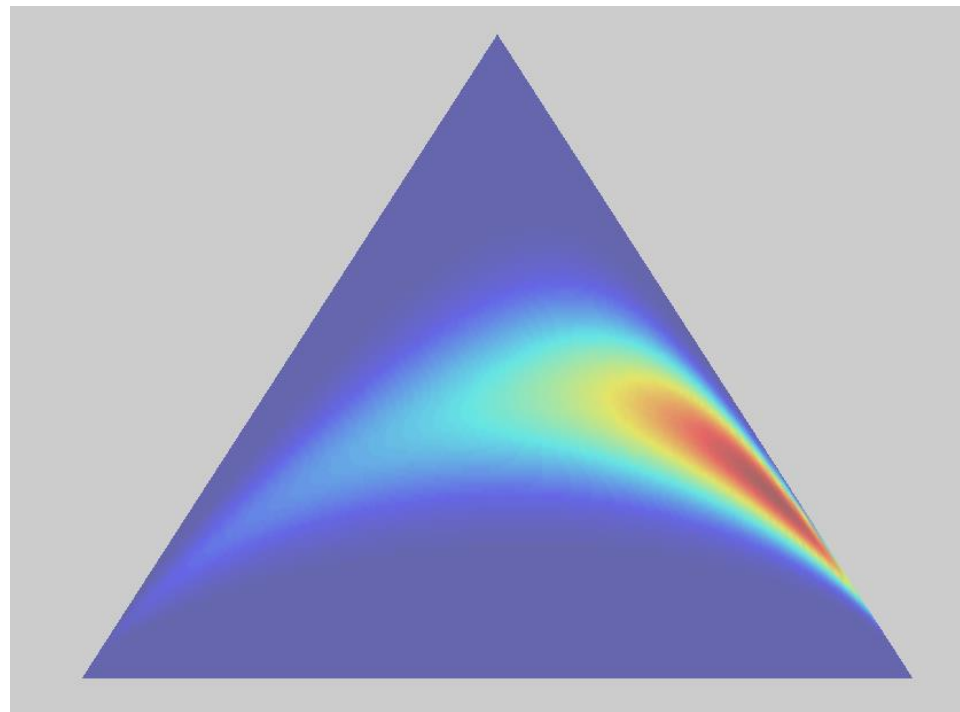
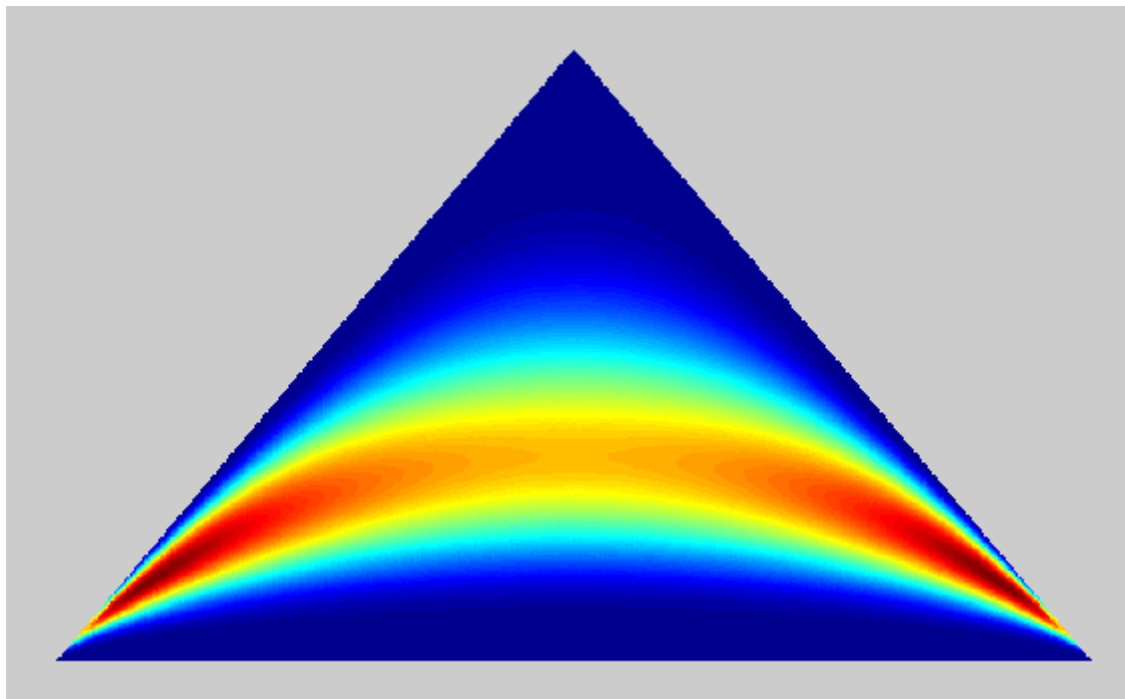
$$\mu = [-0.9, -0.9], \quad \Sigma = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}$$

Favoring Topic 3:



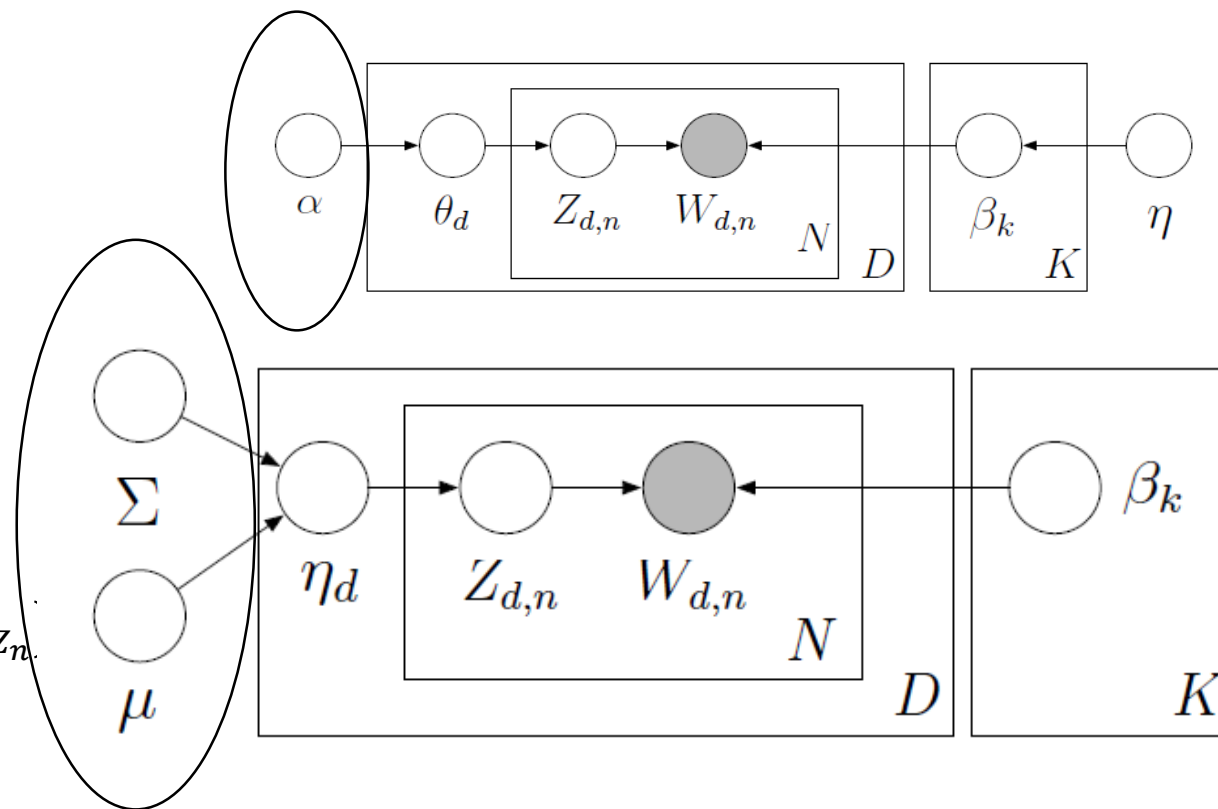
$$\mu = [-0.9, -0.9], \quad \Sigma = [1 \ 0.4; 0.4 \ 1]$$

Exercises



Correlated Topic Model

- Algorithm:
- $\forall d \in D$
 - Draw $\eta_d | \{\mu, \Sigma\} \sim \mathcal{N}(\mu, \Sigma)$
 - $\forall n \in \{1 \dots N\}$:
 - Draw topic assignment
 - $Z_{d,n} | \eta_d \sim \text{Categorical}(f(\eta_d))$
 - Draw word
 - $W_{d,n} | \{Z_{d,n}, \beta_{1:K}\} \sim \text{Categorical}(\beta_{Z_{d,n}})$
- Parameter Estimation:
 - Intractable
 - User variational inference (later)



Evaluation I: CTM on Test Data

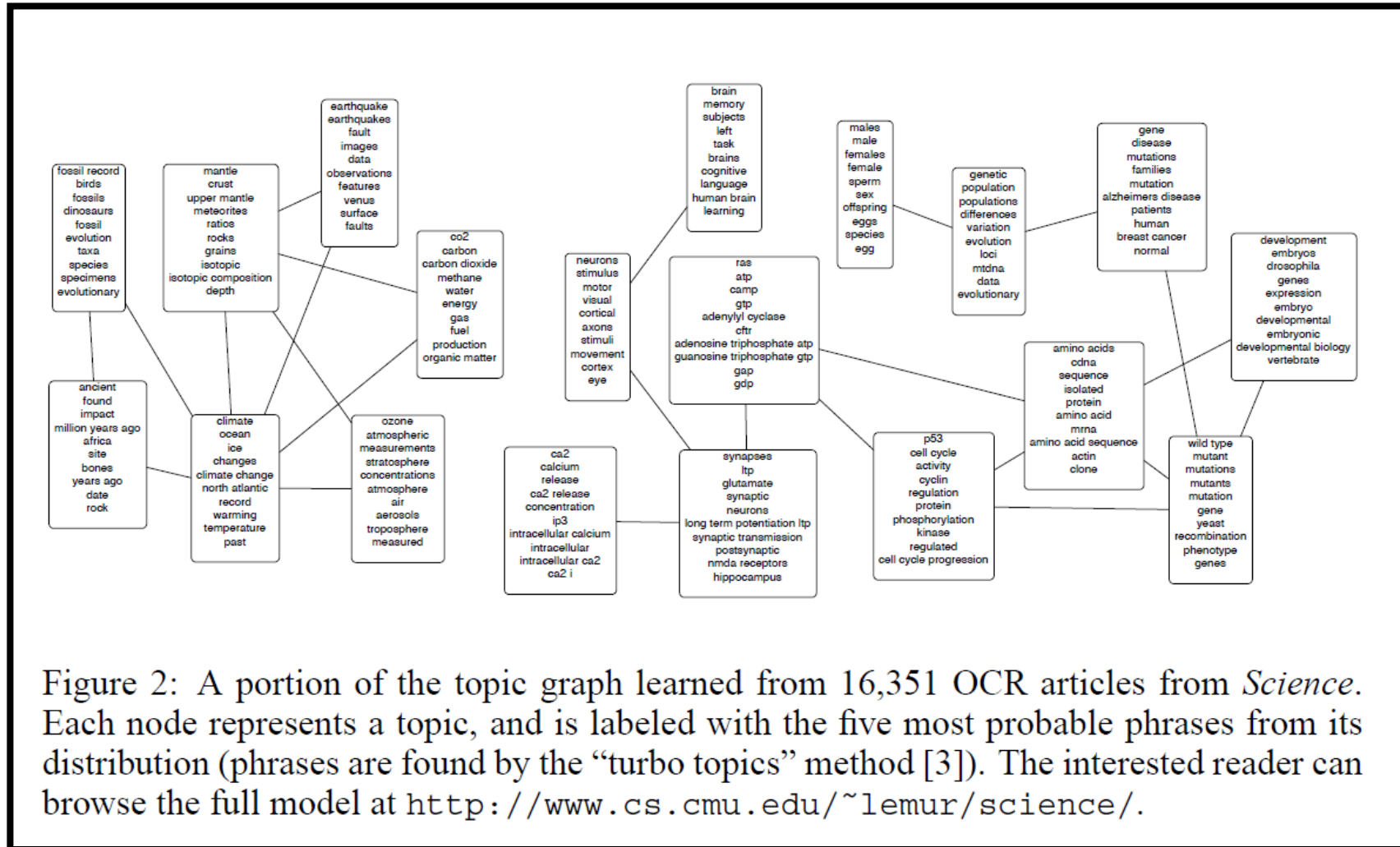
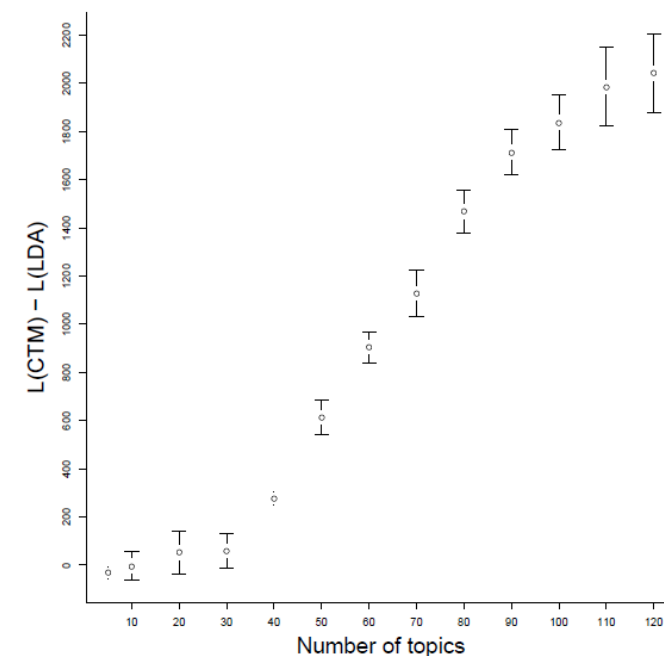
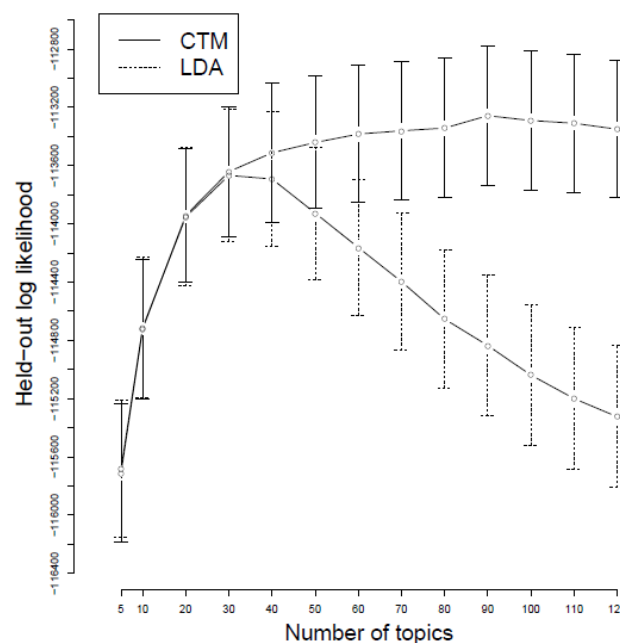


Figure 2: A portion of the topic graph learned from 16,351 OCR articles from *Science*. Each node represents a topic, and is labeled with the five most probable phrases from its distribution (phrases are found by the “turbo topics” method [3]). The interested reader can browse the full model at <http://www.cs.cmu.edu/~lemur/science/>.

Evaluation II: 10-Fold Cross Validation LDA vs CTM

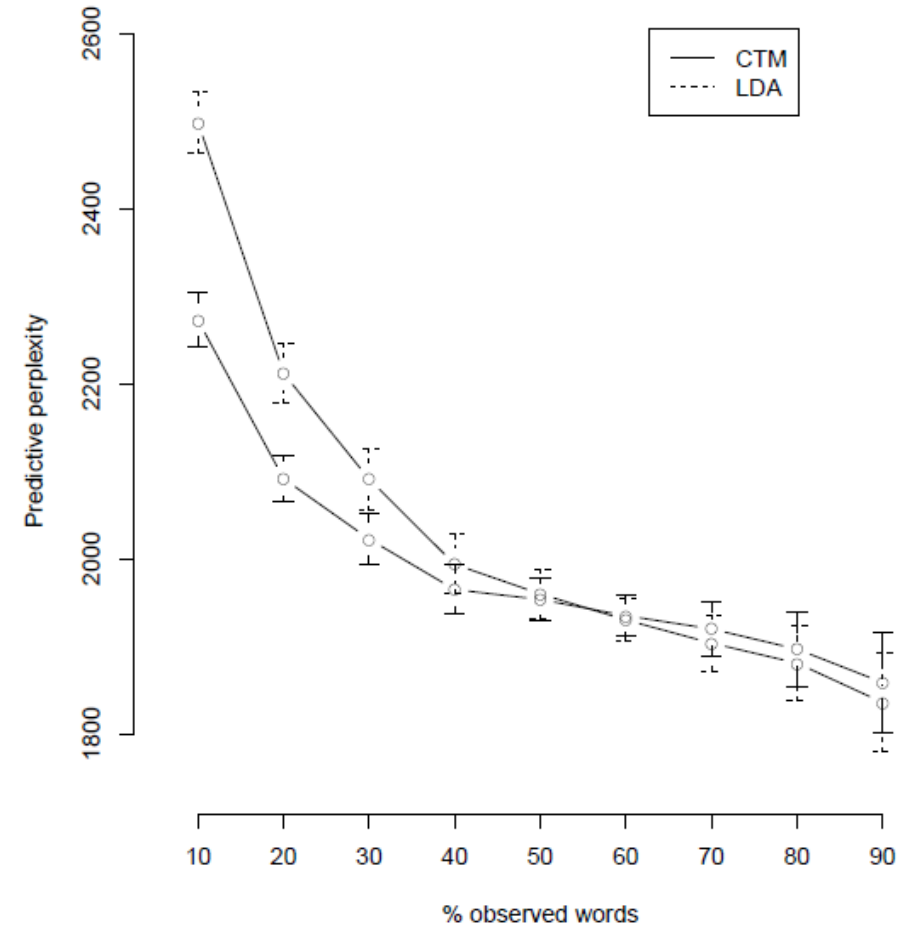
- ~1500 documents in corpus.
- ~5600 unique words
 - After pruning
- Methodology:
 - Partition data into 10 sets
 - 10 fold cross validation
 - Calculate the log likelihood of a set, given you trained on the previous 9 sets, for both LDA and CTM.
- Right($L(\text{CTM}) - L(\text{LDA})$)
- Left($L(\text{CTM}) - L(\text{LDA})$)



CTM shows a much higher log likelihood as the number of topics increases.

Evaluation II: Predictive Perplexity

- Perplexity measure \equiv expected number of equally likely words
 - Lower perplexity means higher word resolution.
- Suppose you see a percentage of words in a document, how likely is the rest of the words in the document according to your model?
- CTM does better with lower #'s of observed words.
 - Able to infer certain words given topic probabilities.



Conclusions

- CTM changes the distribution from which hyper parameters are drawn, from a Dirichlet to a logistic normal function.
 - Very similar to LDA
- Able to model correlations between topics.
- For larger topic sizes, CTM performs better than LDA.
- With known topics, CTM is able to infer words associations better than LDA.