

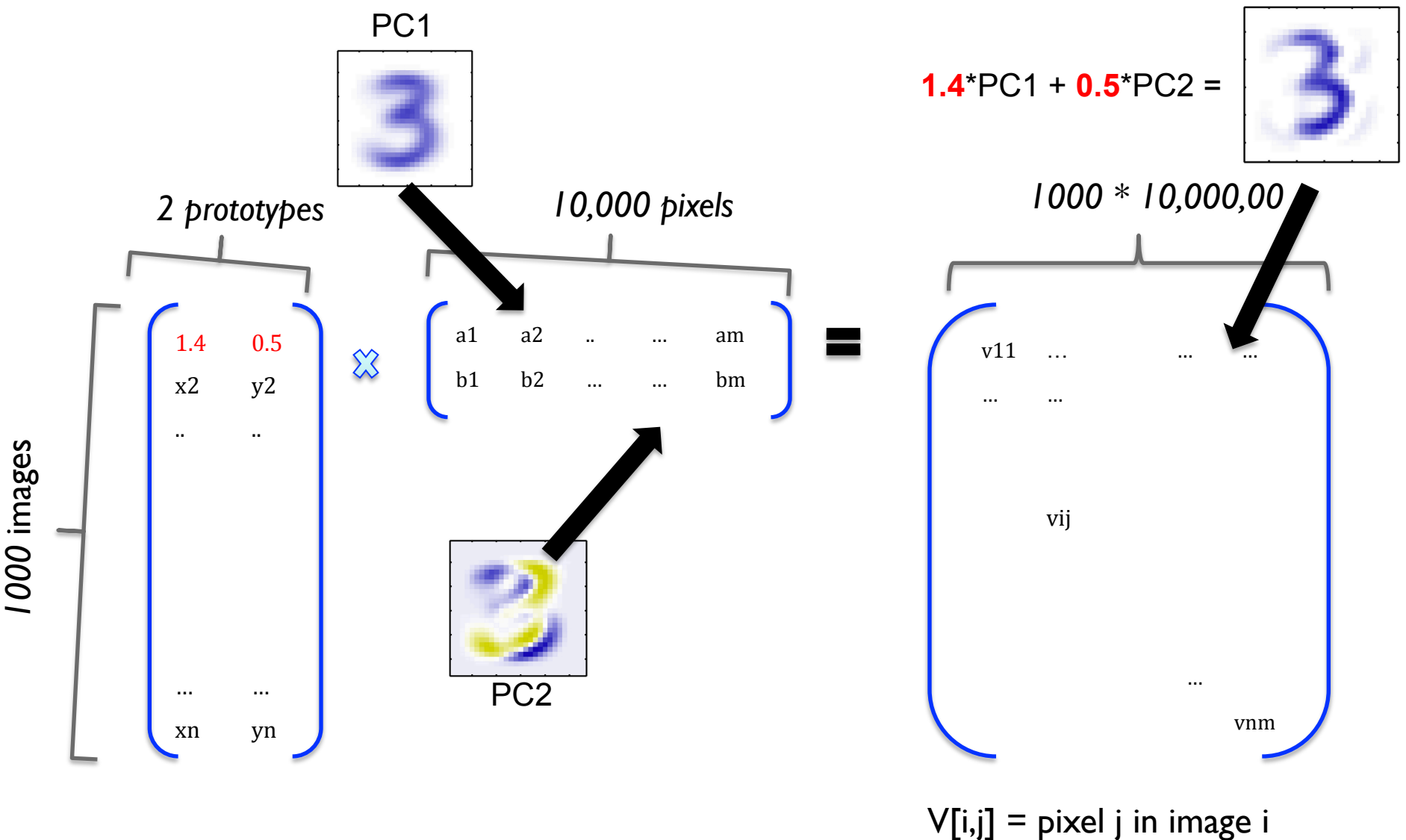
Techniques for Dimensionality Reduction

Topic Models/LDA
Latent Dirichlet Allocation

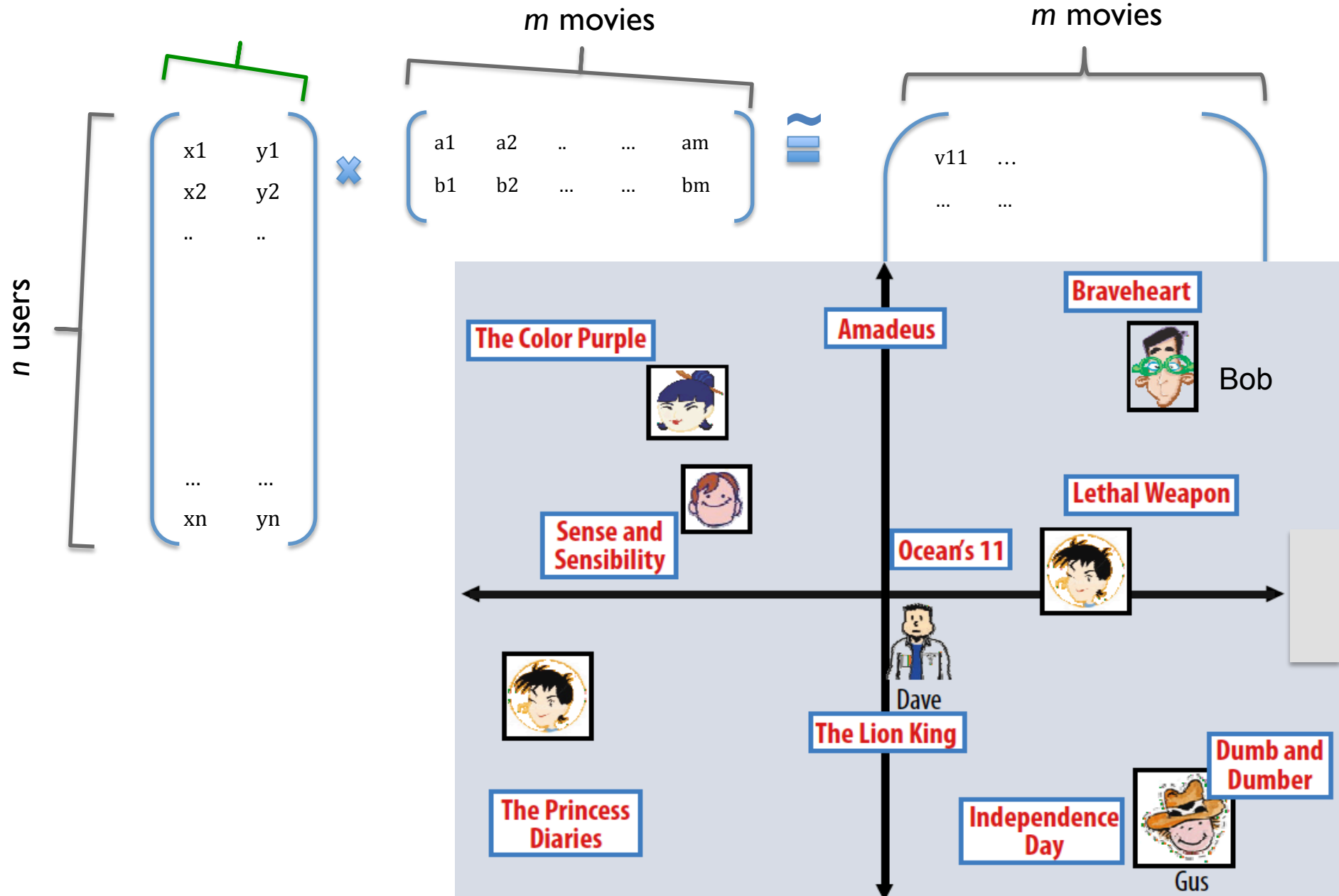
Outline

- Recap: PCA
- Directed models for text
 - Naïve Bayes
 - unsupervised Naïve Bayes (mixtures of multinomials)
 - probabilistic latent semantic indexing pLSI
 - and connection to PCA, LSI, MF
 - latent Dirichlet allocation (LDA)
 - inference with Gibbs sampling
- LDA-like models for graphs & etc

Quick Recap....

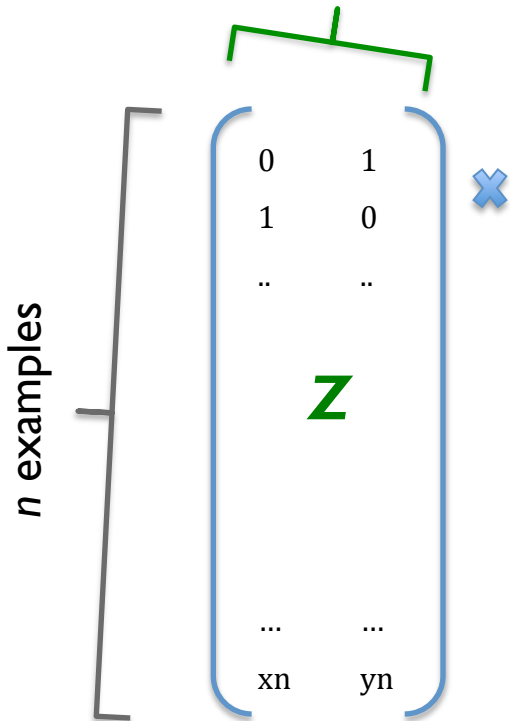


PCA for movie recommendation...

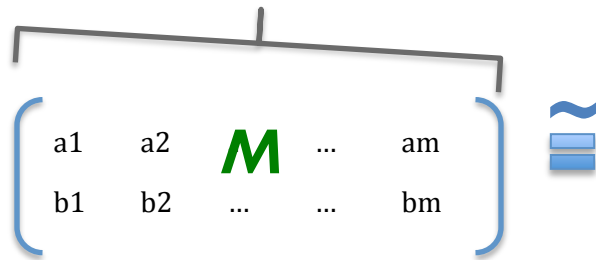


..... vs k-means

indicators for r
clusters

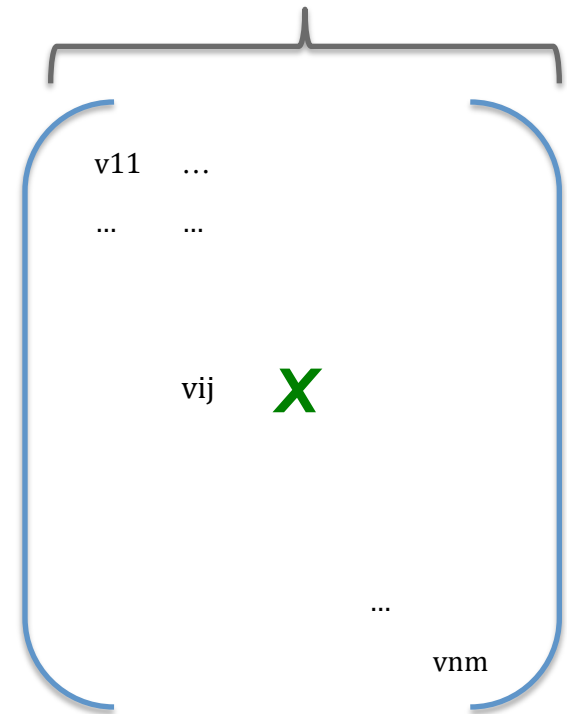


cluster means



Question: what other
generative models look like MF?

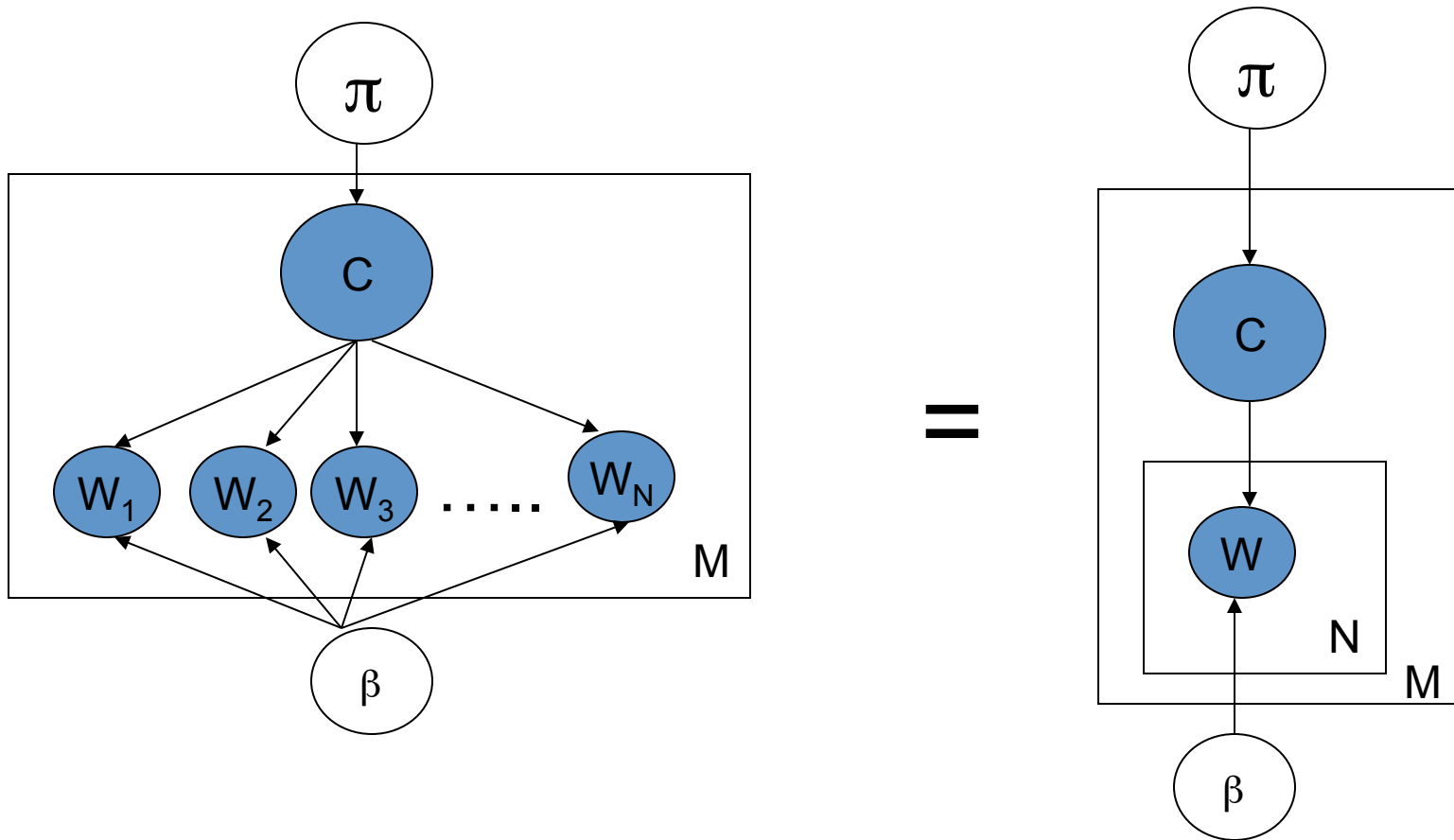
original data set



LDA AND OTHER DIRECTED MODELS FOR MODELING TEXT

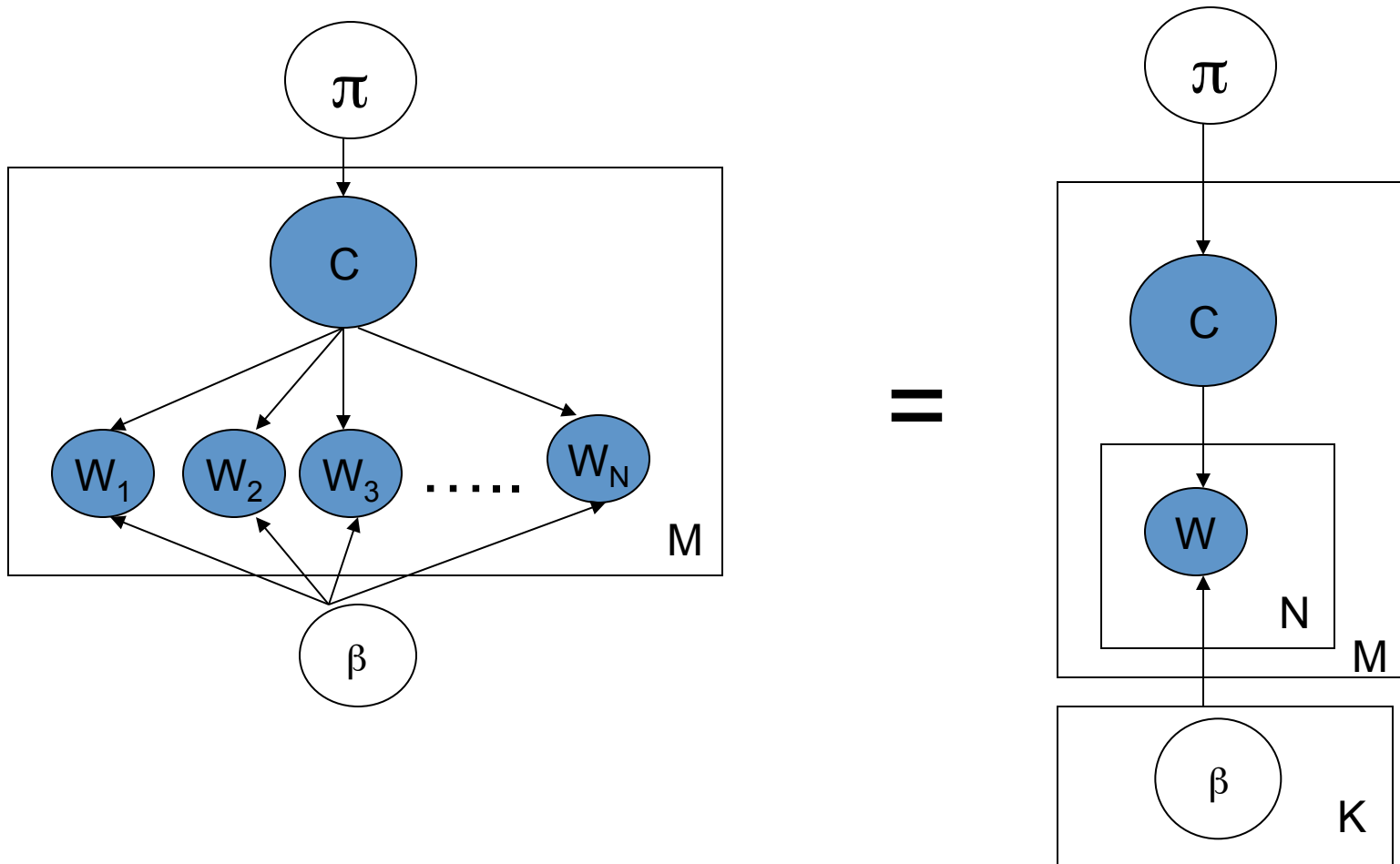
Supervised Multinomial Naïve Bayes

- Naïve Bayes Model: Compact notation



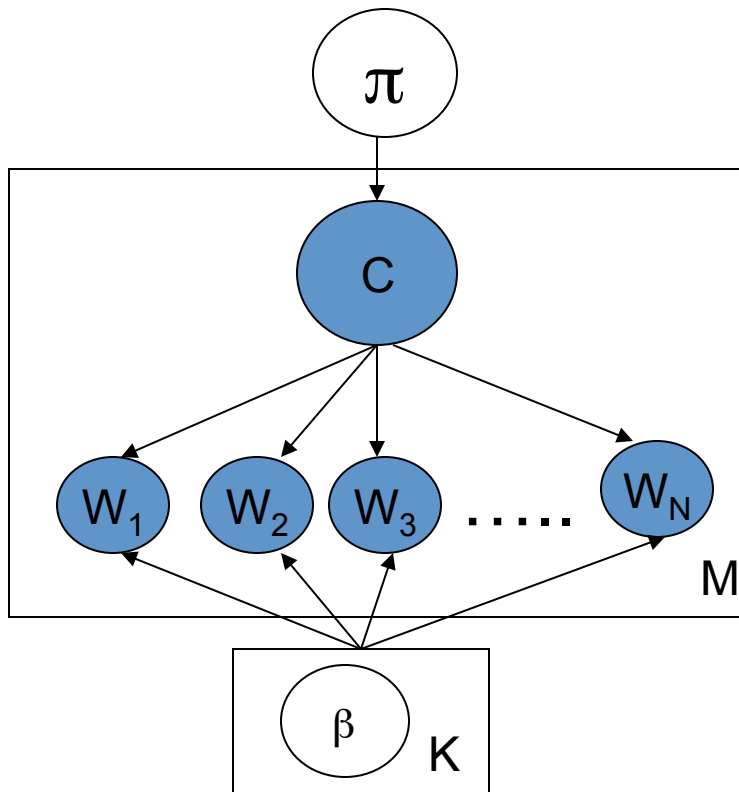
Supervised Multinomial Naïve Bayes

- Naïve Bayes Model: Compact representation



Supervised Naïve Bayes

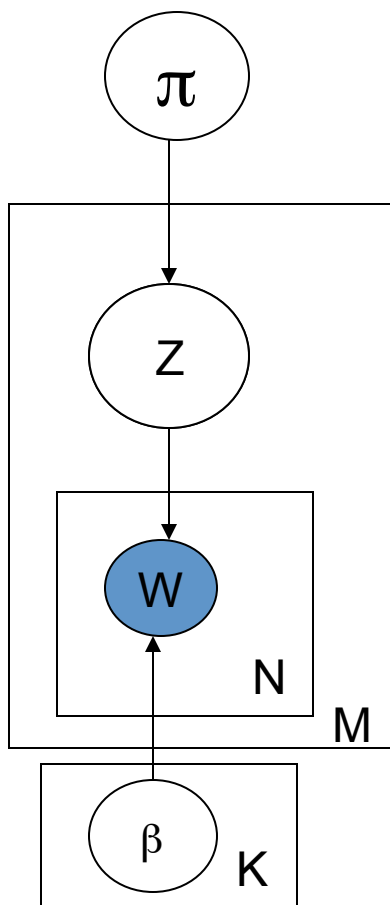
- Multinomial Naïve Bayes



- For each class $1..K$
 - Construct a multinomial β_i
- For each document $d = 1, \dots, M$
 - Generate $C_d \sim \text{Mult}(\cdot | \pi)$
 - For each position $n = 1, \dots, N_d$
 - Generate $w_n \sim \text{Mult}(\cdot | \beta, C_d) \dots$
or if you prefer $w_n \sim \text{Pr}(w | \beta_{C_d})$

Unsupervised Naïve Bayes

- Mixture model: unsupervised naïve Bayes



- Joint probability of words and classes:

$$\prod_{d=1}^M P(w_1, \dots, w_{N_d}, z_d | \beta, \pi) = \prod_{d=1}^M \left\{ \pi_{z_d} \prod_{n=1}^{N_d} \beta_{z_d, w_n} \right\}$$



$$\prod_{d=1}^M P(w_1, \dots, w_{N_d} | \pi, \beta) = \prod_{d=1}^{N_d} \left\{ \sum_{k=1}^K \left(\pi_k \prod_{n=1}^{N_d} \beta_{k, w_n} \right) \right\}$$

Solved using EM

Unsupervised Naïve Bayes

$$\log \left(\sum_{k=1}^K \left[\pi_k \prod_{n=1}^{N_d} \beta_{k,w_n} \right] \right) \geq \sum_{k=1}^K \left\{ \gamma_k \log \left(\pi_k \prod_{n=1}^{N_d} \beta_{k,w_n} \right) \right\} + H(\gamma)$$

- Mixture model: EM solution

E-step:

$$\gamma_{dk}^{(t+1)} = \frac{\pi_k^{(t)} \prod_{n=1}^{N_d} \beta_{k,w_n}^{(t)}}{\sum_k \pi_k^{(t)} \prod_{n=1}^{N_d} \beta_{k,w_n}^{(t)}}$$

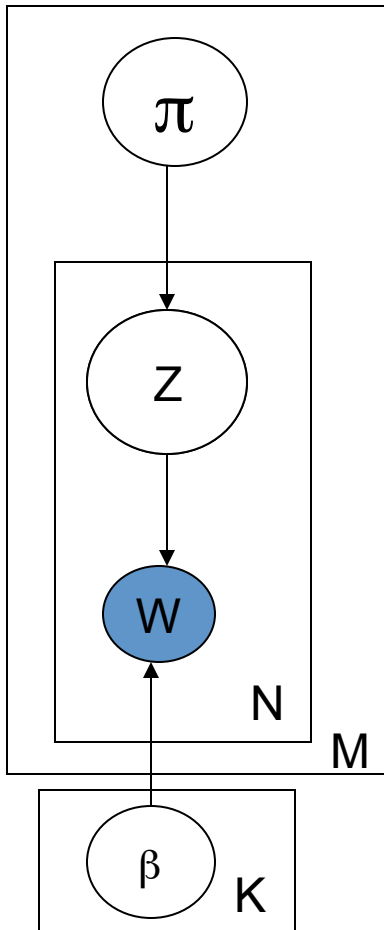
M-step:

$$\left\{ \begin{aligned} \pi_k^{(t+1)} &= \frac{\sum_{d=1}^M \gamma_{dk}^{(t)}}{M} \\ \beta_{k,w}^{(t+1)} &= \frac{\sum_{d=1}^M \gamma_{dk}^{(t)} n(d, w)}{\sum_{d=1}^M \gamma_{nk}^{(t)} \sum_w n(d, w)} \end{aligned} \right.$$

Key capability: estimate distribution of **latent variables** given **observed variables**

Beyond Naïve Bayes - Probabilistic Latent Semantic Indexing (PLSI)

- Every document is a mixture of topics

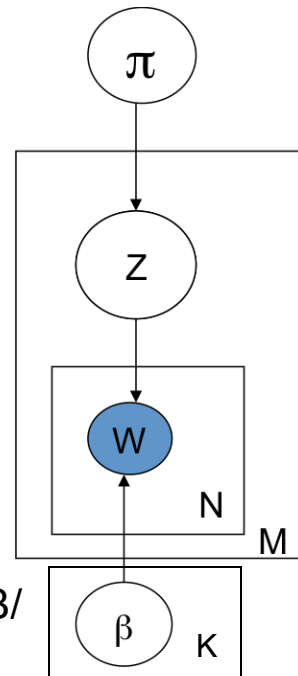


- For $i=1\dots K$:
 - Let β_i be a multinomial over words
- For each document d :
 - Let π_d be a distribution over $\{1,\dots,K\}$
 - For each word position in d :
 - Pick a topic z from π_d
 - Pick a word w from β_i

Latent variable for each word not for each document

Can still learn with EM

compare to unsupervised NB/
mixture of multinomials

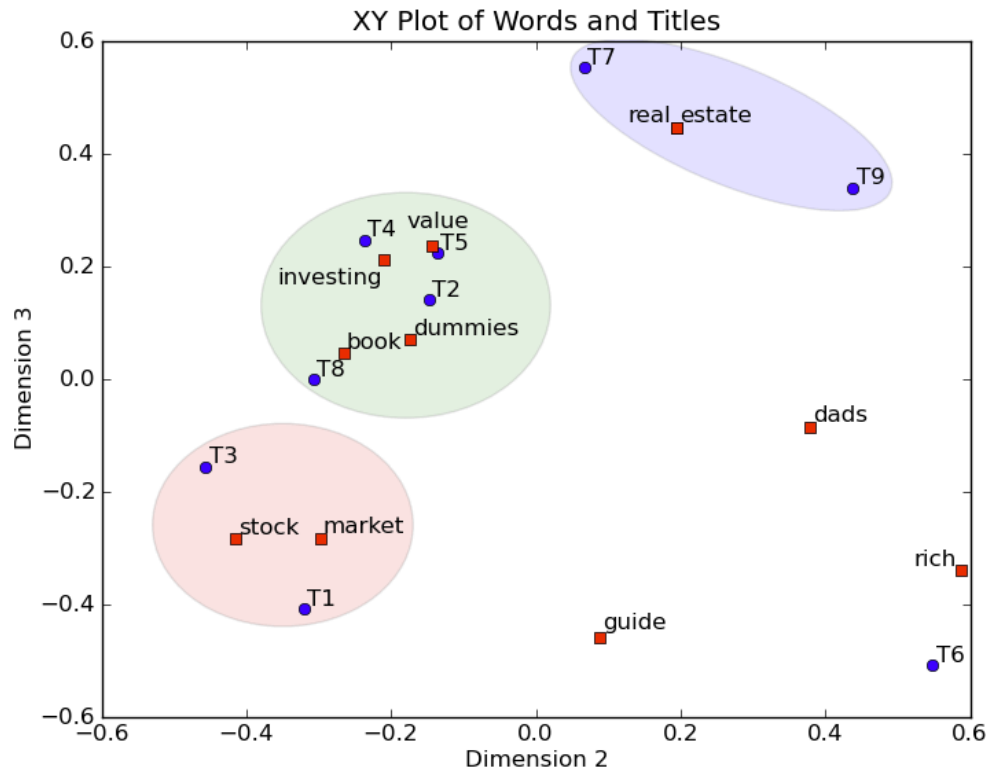
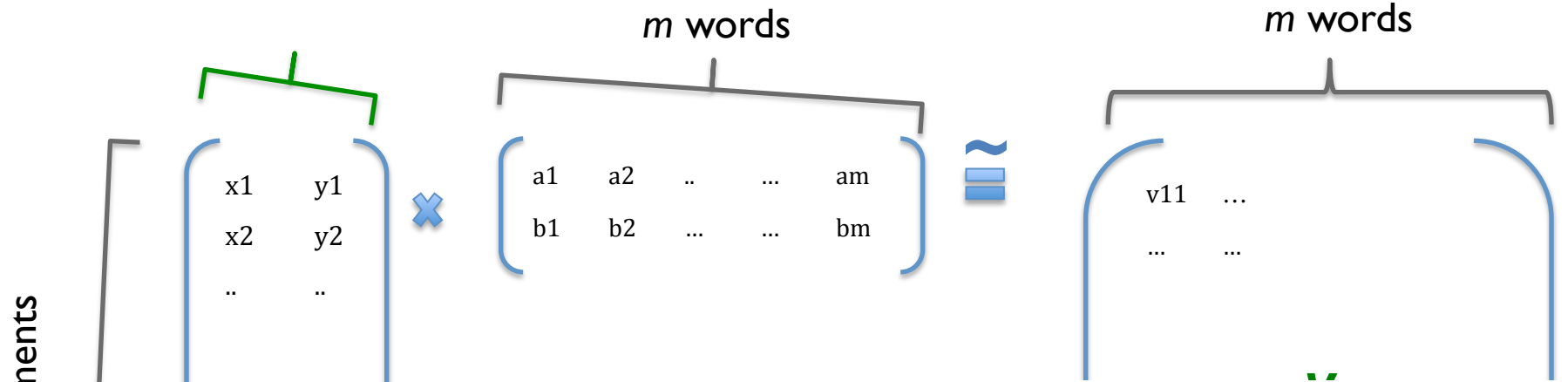


**LATENT SEMANTIC INDEXING
VS
SINGULAR VALUE DECOMPOSITION
VS
PRINCIPLE COMPONENTS ANALYSIS**

Background: PCA vs LSI

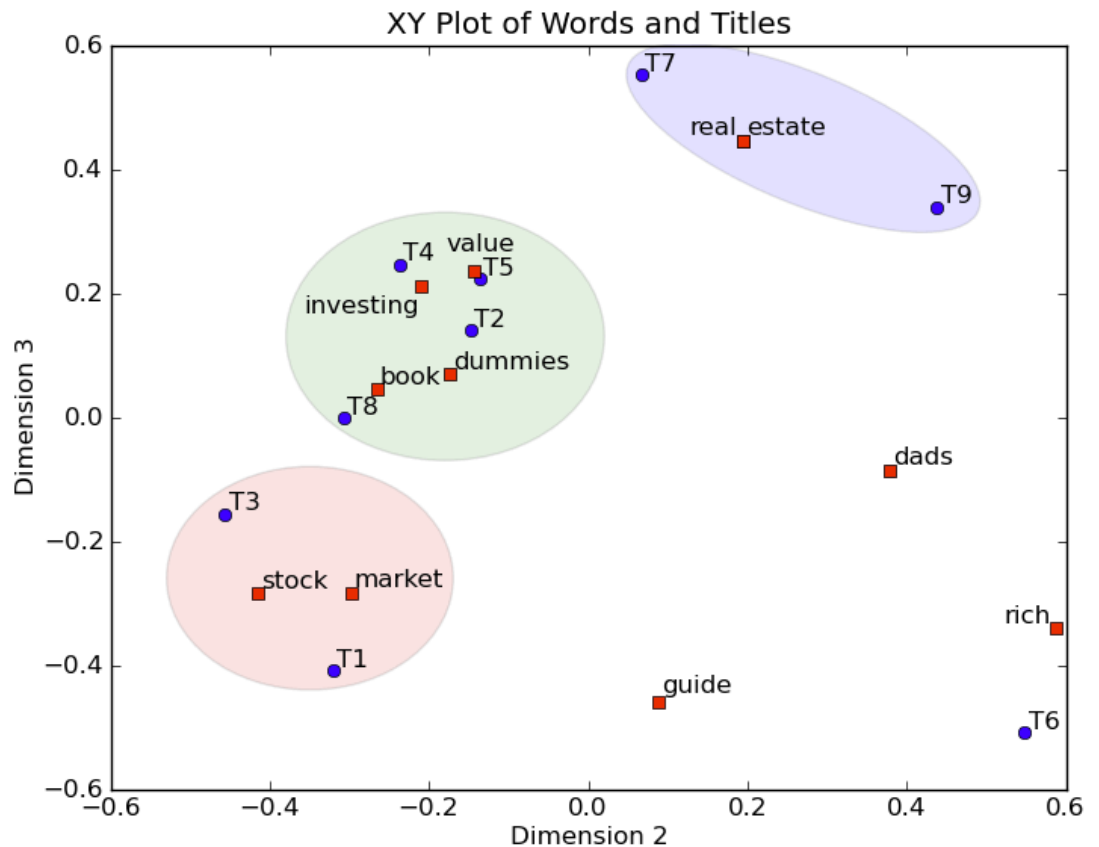
- Some connections
 - PCA is closely related to *singular value decomposition* (SVD): see handout
 - we factored data matrix X into $Z*V$, V is eigenvectors of C_X
 - *one more step*: factor Z into $U*\Sigma$ where Σ is diagonal and $\Sigma(i,i)=\text{sqrt}(\lambda_i)$ and variables in U have unit variance
 - then $X = U*\Sigma*V$ and this is called SVD
- When X is a term-document matrix this is called *latent semantic indexing* (LSI)

LSI



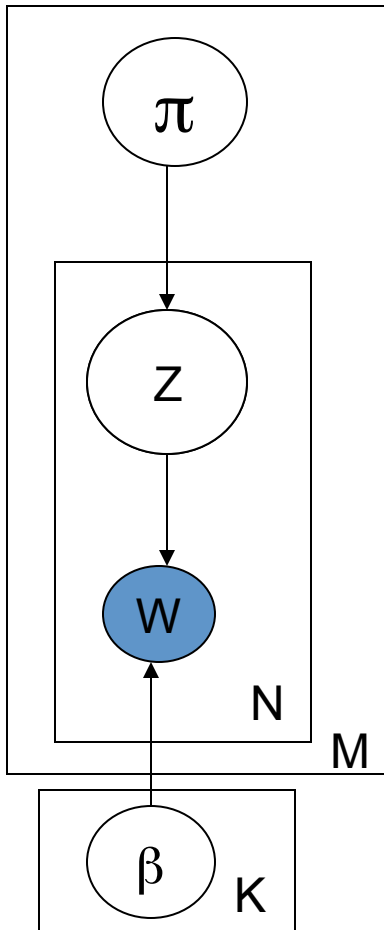
LSI

- i. The Neatest Little Guide to Stock Market Investing
- ii. Investing For Dummies, 4th Edition
- iii. The Little Book of Common Sense Investing: The Only Way to Guarantee Your Fair Share of Stock M
- iv. The Little Book of Value Investing
- v. Value Investing: From Graham to Bufl
- vi. Rich Dad's Guide to Investing: What tl and the Middle Class Do Not!
- vii. Investing in Real Estate, 5th Edition
- viii. Stock Investing For Dummies
- ix. Rich Dad's Advisors: The ABC's of Re Finding Hidden Profits Most Investor:



Beyond Naïve Bayes - Probabilistic Latent Semantic Indexing (PLSI)

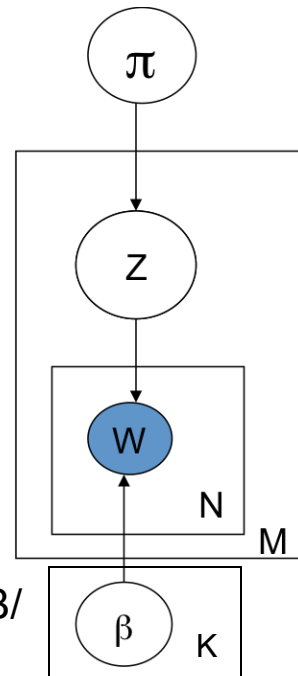
- Every document is a mixture of topics



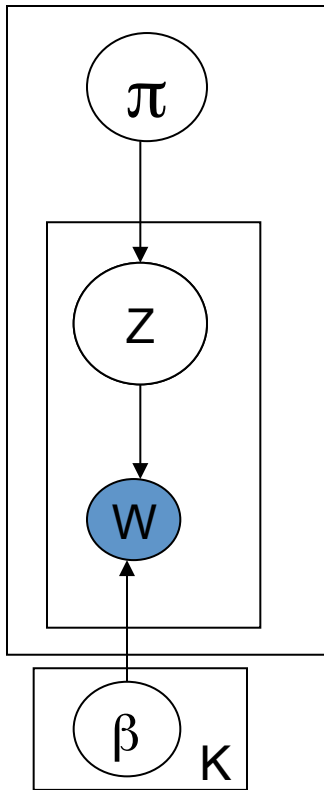
- For $i=1\dots K$:
 - Let β_i be a multinomial over words
- For each document d :
 - Let π_d be a distribution over $\{1,\dots,K\}$
 - For each word position in d :
 - Pick a topic z from π_d
 - Pick a word w from β_i

Latent variable for each word not for each document

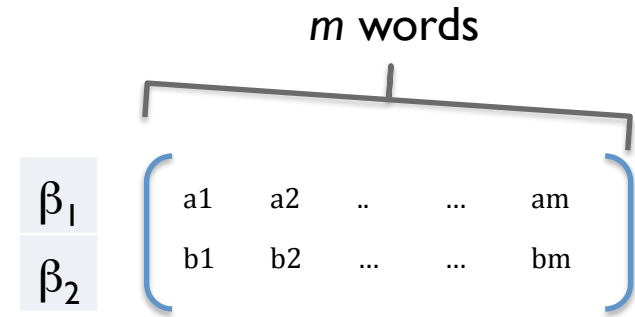
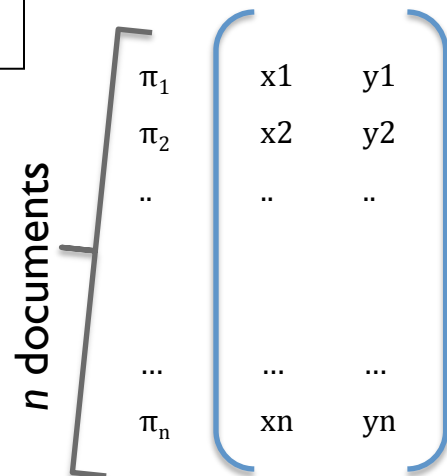
compare to unsupervised NB/
mixture of multinomials



Beyond Naïve Bayes - Probabilistic Latent Semantic Indexing (PLSI)

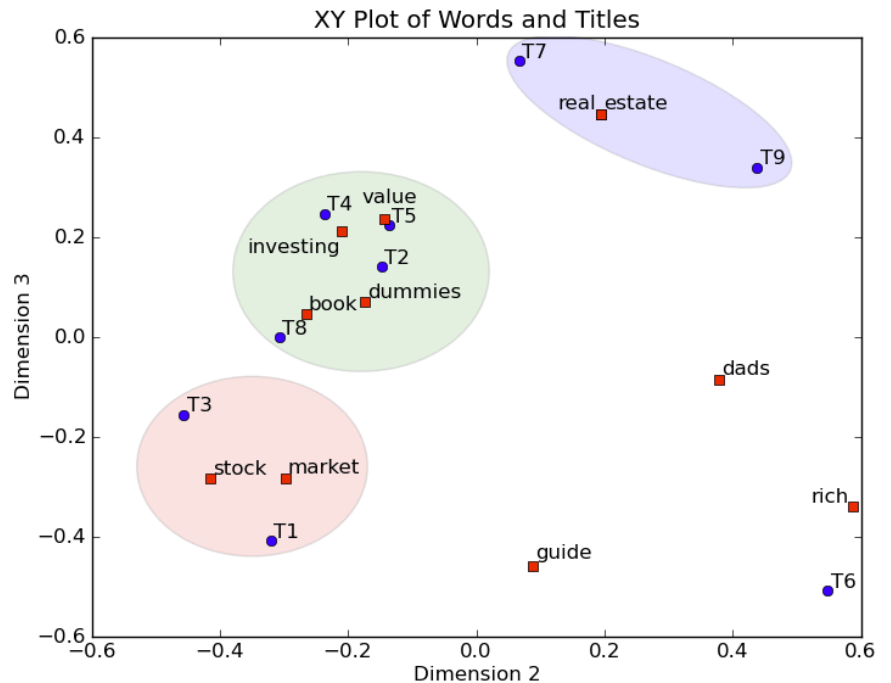
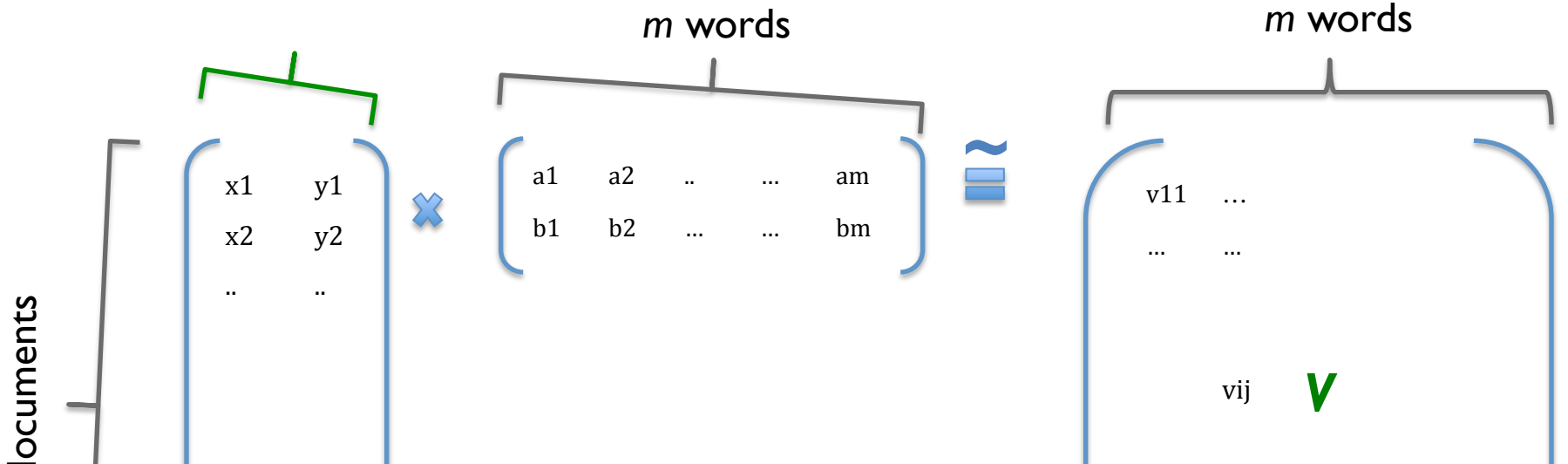


- For $i=1 \dots K$:
 - Let β_i be a multinomial over words
- For each document d :
 - Let π_d be a distribution over $\{1, \dots, K\}$
 - For each word position in d :
 - Pick a topic z from π_d
 - Pick a word w from β_i



$$\Pr[X(i,j)] = \sum_k \Pr(Z = k | \pi[i,k]) * \Pr(W = j | \beta[k,j])$$

LSI or pLSI



Background: LSI vs NMF

- <http://papers.nips.cc/paper/1861-algorithms-for-non-negative-matrix-factorization.pdf>

PLSI and MF and Nonnegative MF



Available online at www.sciencedirect.com



Computational Statistics and Data Analysis 52 (2008) 3913–3927

**COMPUTATIONAL
STATISTICS
& DATA ANALYSIS**

www.elsevier.com/locate/csa

On the equivalence between Non-negative Matrix Factorization and Probabilistic Latent Semantic Indexing

Chris Ding^a, Tao Li^{b,*}, Wei Peng^b

PLSI and MF and Nonnegative MF

The general form of NMF is

$$F \cong CH^T, \quad \text{We considered minimized L2 reconstruction error.}$$

Another choice: constrain C, H to be non-negative and minimize

J_{NMF} :

$$J_{\text{NMF}} = \sum_{i=1}^m \sum_{j=1}^n F_{ij} \log \frac{F_{ij}}{(CH^T)_{ij}} - F_{ij} + (CH^T)_{ij}.$$

PLSI maximizes the likelihood

$$\max J_{\text{PLSI}}, \quad J_{\text{PLSI}} = \sum_{i=1}^m \sum_{j=1}^n F_{ij} \log P(w_i, d_j)$$

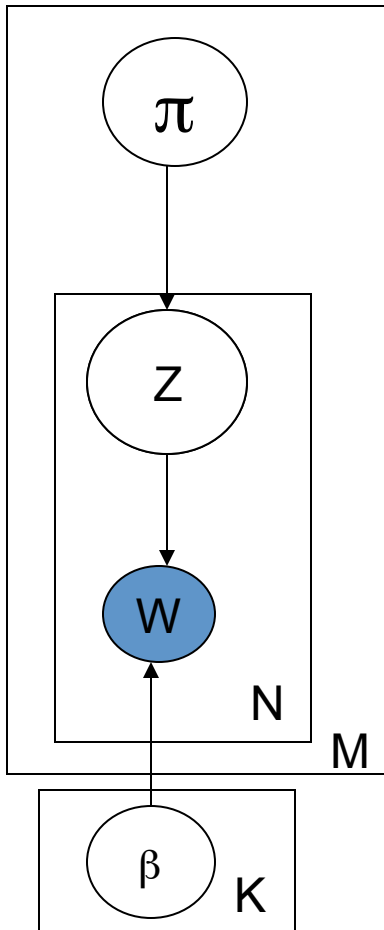
$$\text{where } P(w_i, d_j) = \sum_k P(w_i, d_j | z_k) P(z_k)$$

$$= \sum_k P(w_i | z_k) P(d_j | z_k) P(z_k),$$

Ding et al: the objective functions for NMF and PLSI are equivalent

Beyond Naïve Bayes - Probabilistic Latent Semantic Indexing (PLSI)

- Every document is a mixture of topics



- For $i=1\dots K$:
 - Let β_i be a multinomial over words
- For each document d :
 - Let π_d be a distribution over $\{1,\dots,K\}$
 - For each word position in d :
 - Pick a topic z from π_d
 - Pick a word w from β_i
- **Turns out to be hard to fit:**
 - **Lots of parameters!**
 - **Also: only applies to the training data**

LATENT DIRICHLET ANALYSIS (LDA)

The LDA Topic Model

Latent Dirichlet Allocation

David M. Blei

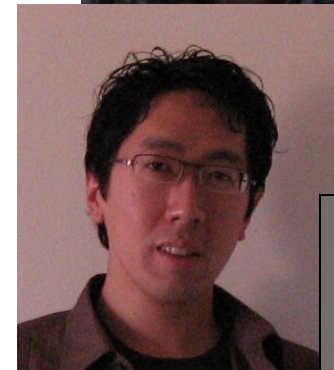
*Computer Science Division
University of California
Berkeley, CA 94720, USA*

Andrew Y. Ng

*Computer Science Department
Stanford University
Stanford, CA 94305, USA*

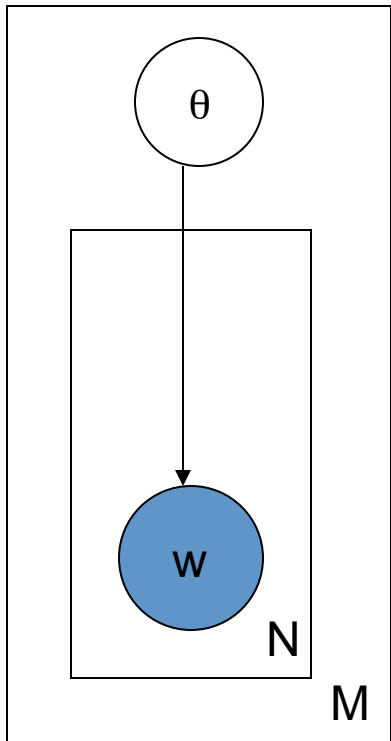
Michael I. Jordan

*Computer Science Division and Department of Statistics
University of California
Berkeley, CA 94720, USA*



LDA

- Motivation



Assumptions: 1) documents are i.i.d 2) *within* a document, words are i.i.d. (bag of words)

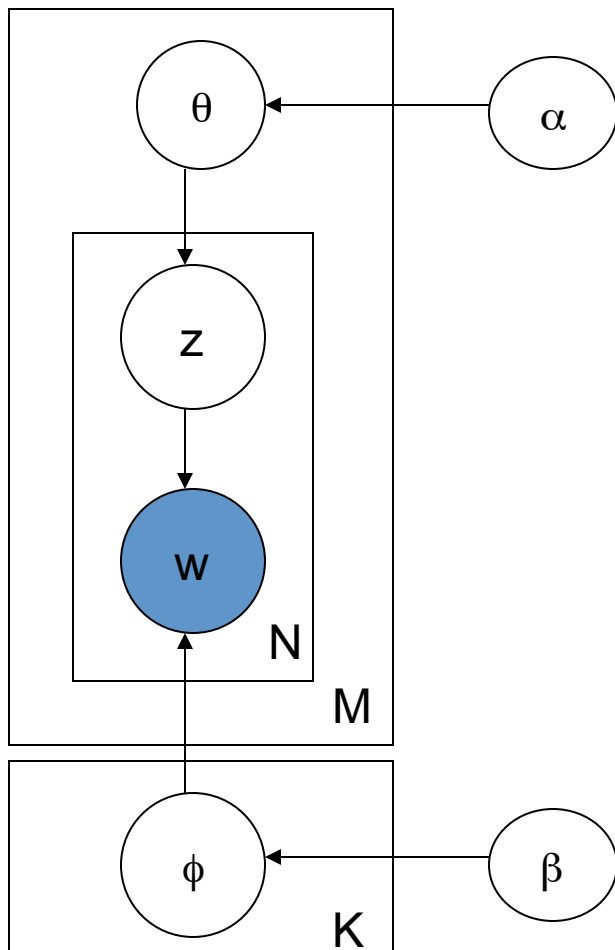
- For each document $d = 1, \dots, M$
 - Generate $\theta_d \sim D_1(\dots)$
 - For each word $n = 1, \dots, N_d$
 - generate $w_n \sim D_2(\cdot | \vartheta_{d_n})$

Now pick your favorite distributions for D_1, D_2

$$\Pr(z = j | n_1, n_2, \dots, n_k, \alpha) = \frac{n_j + \alpha_j}{n_1 + \alpha_1 + \dots + n_k + \alpha_k}$$

“Mixed membership”

- Latent Dirichlet Allocation



- For each document $d = 1, \dots, M$
 - Generate $\theta_d \sim \text{Dir}(\cdot | \alpha)$
 - For each position $n = 1, \dots, N_d$
 - generate $z_n \sim \text{Mult}(\cdot | \theta_d)$
 - generate $w_n \sim \text{Mult}(\cdot | \beta_{z_n})$

$$\prod_{d=1}^{N_d} P(w_1, \dots, w_{N_d} | \beta, \alpha)$$

$$= \prod_{d=1}^{N_d} \int_{\theta_d} P(\theta_d | \alpha) \left\{ \prod_{n=1}^{N_d} \left(\sum_k \theta_{dk} \beta_{kw_n} \right) \right\} d\theta_d$$

- LDA' s view of a document

The William Randolph Hearst Foundation will give \$1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard 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 Juilliard 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.

“Arts”

“Budgets”

“Children”

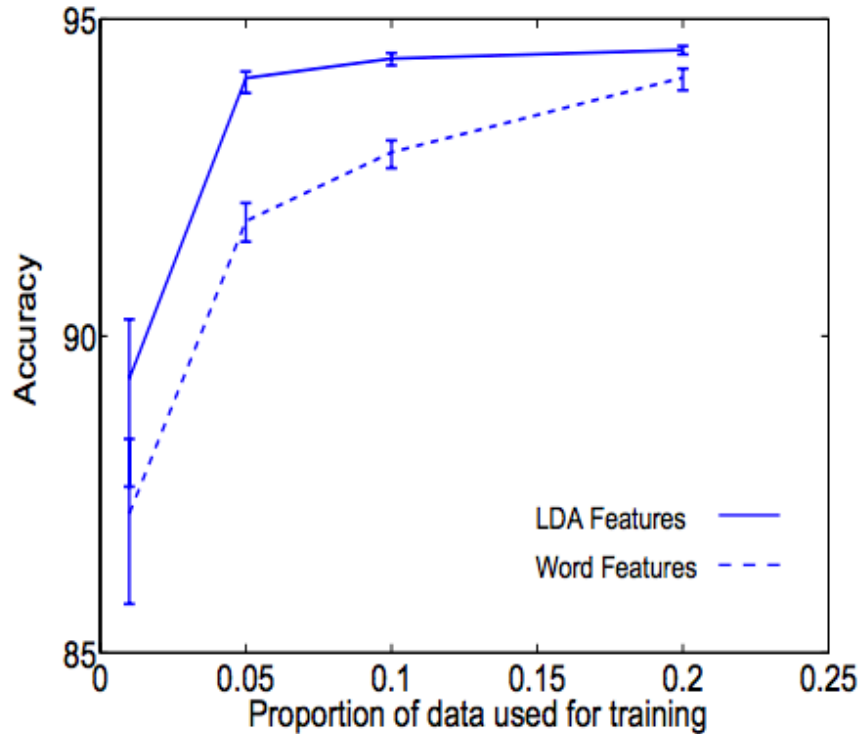
“Education”

- LDA topics

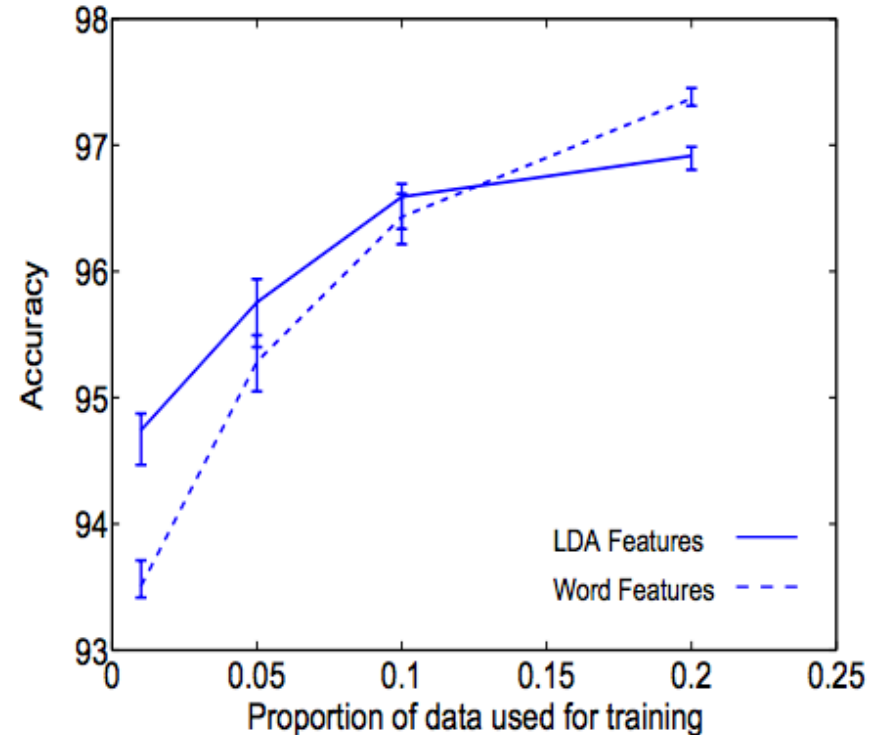
“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

LDA used as dimension reduction for classification

50 topics vs all words, SVM



(a)



(b)

Figure 10: Classification results on two binary classification problems from the Reuters-21578 dataset for different proportions of training data. Graph (a) is EARN vs. NOT EARN. Graph (b) is GRAIN vs. NOT GRAIN.

LDA used for CF

Users rating 100+ movies only

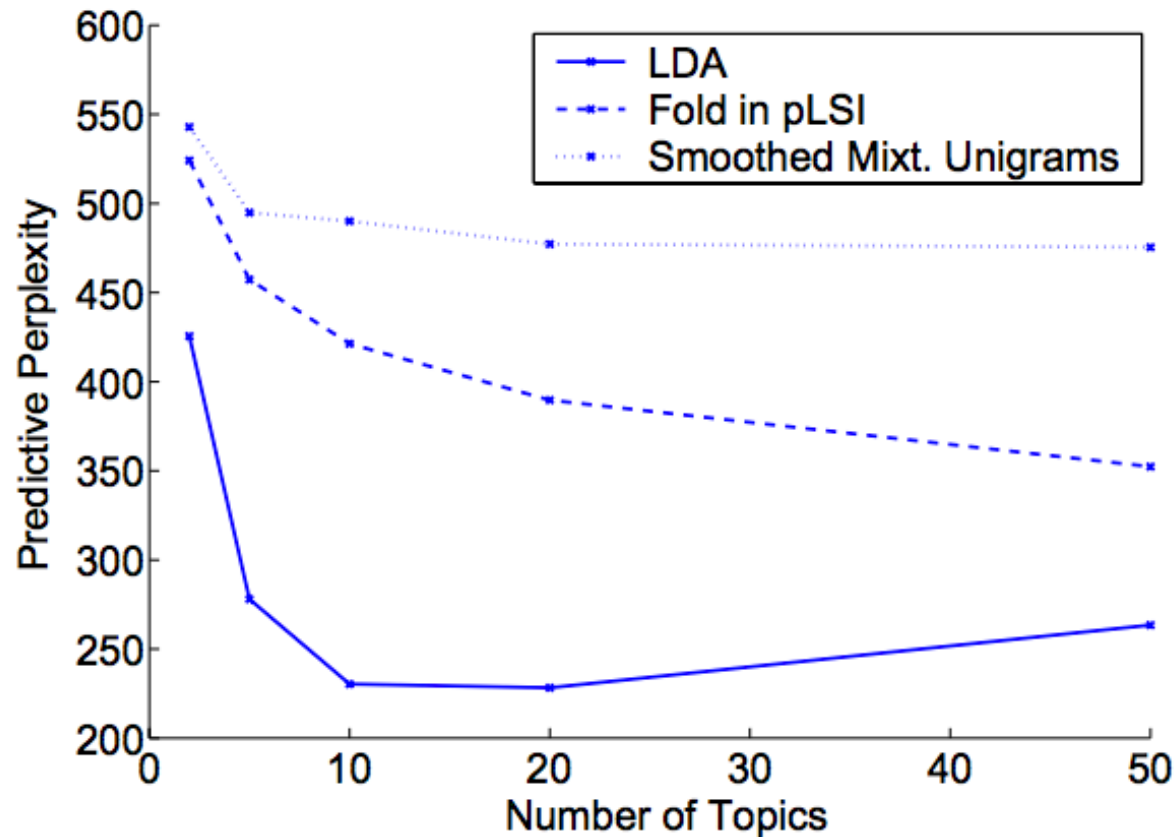


Figure 11: Results for collaborative filtering on the EachMovie data.

LDA

- Latent Dirichlet Allocation
 - Parameter learning:
 - Variational EM
 - Numerical approximation using lower-bounds
 - Results in biased solutions
 - Convergence has numerical guarantees
 - Gibbs Sampling
 - Stochastic simulation
 - unbiased solutions
 - Stochastic convergence

LDA

- Gibbs sampling – works for *any* directed model!
 - Applicable when joint distribution is hard to evaluate but conditional distribution is known
 - Sequence of samples comprises a Markov Chain
 - Stationary distribution of the chain is the joint distribution

1. Initialise $x_{0,1:n}$.

2. For $i = 0$ to $N - 1$

– Sample $x_1^{(i+1)} \sim p(x_1 | x_2^{(i)}, x_3^{(i)}, \dots, x_n^{(i)})$.

– Sample $x_2^{(i+1)} \sim p(x_2 | x_1^{(i+1)}, x_3^{(i)}, \dots, x_n^{(i)})$.

⋮

– Sample $x_j^{(i+1)} \sim p(x_j | x_1^{(i+1)}, \dots, x_{j-1}^{(i+1)}, x_{j+1}^{(i)}, \dots, x_n^{(i)})$.

⋮

– Sample $x_n^{(i+1)} \sim p(x_n | x_1^{(i+1)}, x_2^{(i+1)}, \dots, x_{n-1}^{(i+1)})$.

Key capability: estimate distribution of **one** latent variables given **the other latent variables** and observed variables.

Why does Gibbs sampling work?

- What's the fixed point?
 - Stationary distribution of the chain is the joint distribution
- When will it converge (in the limit)?
 - If graph defined by the chain is connected
- How long will it take to converge?
 - Depends on second eigenvector of that graph

□ initialisation

zero all count variables, $n_m^{(k)}, n_m, n_k^{(t)}, n_k$

for all documents $m \in [1, M]$ **do**

for all words $n \in [1, N_m]$ in document m **do**

 sample topic index $z_{m,n}=k \sim \text{Mult}(1/K)$

 increment document–topic count: $n_m^{(k)} + 1$

 increment document–topic sum: $n_m + 1$

 increment topic–term count: $n_k^{(t)} + 1$

 increment topic–term sum: $n_k + 1$

end for

end for

□ Gibbs sampling over burn-in period and sampling period

while not finished **do**

for all documents $m \in [1, M]$ **do**

for all words $n \in [1, N_m]$ in document m **do**

 □ for the current assignment of k to a term t for word $w_{m,n}$:

 decrement counts and sums: $n_m^{(k)} - 1; n_m - 1; n_k^{(t)} - 1; n_k - 1$

sample topic index $\tilde{k} \sim p(z_i | \vec{z}_{-i}, \vec{w})$ ep):

 □ use the new assignment of $z_{m,n}$ to the term t for word $w_{m,n}$ to:

 increment counts and sums: $n_m^{(\tilde{k})} + 1; n_m + 1; n_{\tilde{k}}^{(t)} + 1; n_{\tilde{k}} + 1$

end for

end for

□ check convergence and read out parameters

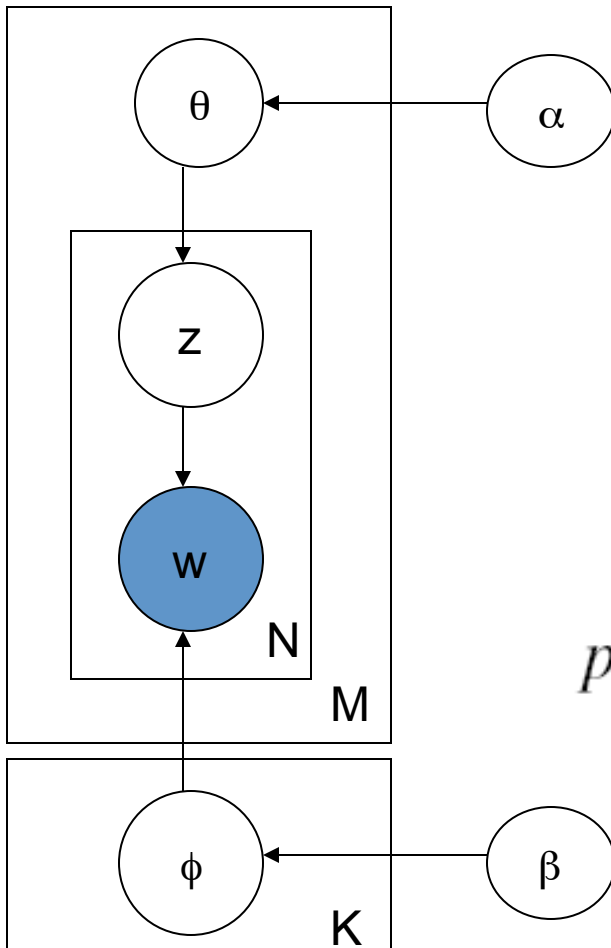
$$\begin{aligned}
p(z_i=k|\vec{z}_{\neg i}, \vec{w}) &= \frac{p(\vec{w}, \vec{z})}{p(\vec{w}, \vec{z}_{\neg i})} = \frac{p(\vec{w}|\vec{z})}{p(\vec{w}_{\neg i}|\vec{z}_{\neg i})p(w_i)} \cdot \frac{p(\vec{z})}{p(\vec{z}_{\neg i})} \\
&\propto \frac{\Delta(\vec{n}_z + \vec{\beta})}{\Delta(\vec{n}_{z,\neg i} + \vec{\beta})} \cdot \frac{\Delta(\vec{n}_m + \vec{\alpha})}{\Delta(\vec{n}_{m,\neg i} + \vec{\alpha})} \\
&\propto \frac{\Gamma(n_k^{(t)} + \beta_t) \Gamma(\sum_{t=1}^V n_{k,\neg i}^{(t)} + \beta_t)}{\Gamma(n_{k,\neg i}^{(t)} + \beta_t) \Gamma(\sum_{t=1}^V n_k^{(t)} + \beta_t)} \cdot \frac{\Gamma(n_m^{(k)} + \alpha_k) \Gamma(\sum_{k=1}^K n_{m,\neg i}^{(k)} + \alpha_k)}{\Gamma(n_{m,\neg i}^{(k)} + \alpha_k) \Gamma(\sum_{k=1}^K n_m^{(k)} + \alpha_k)} \\
&\propto \frac{n_{k,\neg i}^{(t)} + \beta_t}{\sum_{t=1}^V n_{k,\neg i}^{(t)} + \beta_t} \cdot \frac{n_{m,\neg i}^{(k)} + \alpha_k}{[\sum_{k=1}^K n_{m,\neg i}^{(k)} + \alpha_k] - 1}
\end{aligned}$$

Called “collapsed Gibbs sampling” since you’ve marginalized away some variables

LDA

- Latent Dirichlet Allocation

“Mixed membership”



- Randomly initialize each $z_{m,n}$
- Repeat for $t=1, \dots$
 - For each doc m , word n
 - Find $\Pr(z_{mn}=k | \text{other } z' \text{ s})$
 - Sample z_{mn} according to that distr.

$$p(z_i=k | \vec{z}_{-i}, \vec{w}) =$$

$$\propto \frac{n_{k,-i}^{(t)} + \beta_t}{\sum_{t=1}^V n_{k,-i}^{(t)} + \beta_t} \cdot \frac{n_{m,-i}^{(t)} + \alpha_k}{[\sum_{k=1}^K n_m^{(k)} + \alpha_k] - 1}$$

EVEN MORE DETAIL ON LDA...

Way way more detail

```
# topic k, docId d, and wordId w are integer indices
#
# x[d][j] = w, index of j-th word in doc d
# z[d][j] = k, index of latent topic of j-th word in doc d
# vocab[w] = string for the word with index w
#
# totalTopicCount[k] = number of words in topic k
# docTopicCount[d][k] = number of words in topic k for document d
# wordTopicCount[w][k] = number of occurrences of word w in topic k
# totalWords = number of words in the corpus
```

```

# topic k, docId d, and wordId w are integer indices
#
# x[d][j] = w, index of j-th word in doc d
# z[d][j] = k, index of latent topic of j-th word in doc d
# vocab[w] = string for the word with index w
#
# totalTopicCount[k] = number of words in topic k
# docTopicCount[d][k] = number of words in topic k for document d
# wordTopicCount[w][k] = number of occurrences of word w in topic k
# totalWords = number of words in the corpus

```

```

def initGibbs(self):
    print '.initializing latent vars'
    self.totalTopicCount = self.topicCounter()
    self.docTopicCount = [self.topicCounter() for d in xrange(len(self.x))]
    self.wordTopicCount = [self.topicCounter() for w in xrange(len(self.vocab))]
    self.z = [[-1 for j in xrange(len(self.x[d]))] for d in xrange(len(self.x))]
    for d in xrange(len(self.x)):
        if (d+1)%self.dstep==0: print '..',d+1,'of',len(self.x)
        for j in xrange(len(self.x[d])):
            w = self.x[d][j]
            k = random.randint(0, self.numTopics-1)
            self.z[d][j] = k
            self.docTopicCount[d].add(k, 1)
            self.wordTopicCount[w].add(k, 1)
            self.totalTopicCount.add(k, 1)
    #reasonable parameters
    self.alpha = 1.0/self.numTopics
    self.beta = 1.0/len(self.vocab)
    print "alpha:", self.alpha, "beta:", self.beta

```

```

def runGibbs(self, maxT):
    for t in xrange(maxT):
        print '. iteration', t+1, 'of', maxT
        for d in xrange(len(self.x)):
            if (d+1)%self.dstep==0: print '.. doc', d+1, 'of', len(self.x)
            for j in xrange(len(self.x[d])):
                k = self.resample(d, j)
                self.flip(d, j, self.z[d][j], k)

def flip(self, d, j, k_old, k_new):
    """update counts to reflect a changed value of z[d][j]"""
    if k_old != k_new:
        w = self.x[d][j]
        self.docTopicCount[d].add(k_old, -1)
        self.wordTopicCount[w].add(k_old, -1)
        self.wordTopicCount[w].add(k_new, +1)
        self.totalTopicCount.add(k_old, -1)
        self.totalTopicCount.add(k_new, +1)
        self.z[d][j] = k_new

```

```

def runGibbs(self, maxT):
    for t in xrange(maxT):
        print '. iteration', t+1, 'of', maxT
        for d in xrange(len(self.x)):
            if (d+1)%self.dstep==0: print '.. doc', d+1, 'of', len(self.x)
            for j in xrange(len(self.x[d])):
                k = self.resample(d, j)
                self.flip(d, j, self.z[d][j], k)

def resample(self, d, j):
    """sample a new value of z[d][j]"""
    p = []
    norm = 0.0
    #compute pk = Pr(z_dj=k | everything else)
    for k in xrange(self.numTopics):
        w = self.x[d][j]
        z_dj_equals_k = 1 if self.z[d][j]==k else 0
        #unnormalized chance of picking topic k in doc d
        ak = (self.docTopicCount[d][k] - z_dj_equals_k + self.alpha)
        #unnormalized chance of picking topic k for word w
        bk = ((self.wordTopicCount[w][k] - z_dj_equals_k + self.beta)
              / (self.totalTopicCount[k] - z_dj_equals_k + self.totalWords * self.beta))
        pk = ak*bk
        p.append(pk)
        norm += pk
    #pick randomly from the normalized pk
    sum_p_up_to_k = 0.0
    r = random.random()
    for k in xrange(self.numTopics):
        sum_p_up_to_k += p[k]/norm
        if r < sum_p_up_to_k:
            return k

```

What gets learned.....

```
def phi(self, w, k):  
    """weight of word w under topic k"""  
    num = (self.wordTopicCount[w][k] + self.beta)  
    denom = (self.totalTopicCount[k] + self.totalWords * self.beta)  
    return num/denom  
  
def theta(self, d, k):  
    """weight of doc unde  
    num = (self.docTopicC  
    denom = (sum(self.doc'  
    return num/denom
```

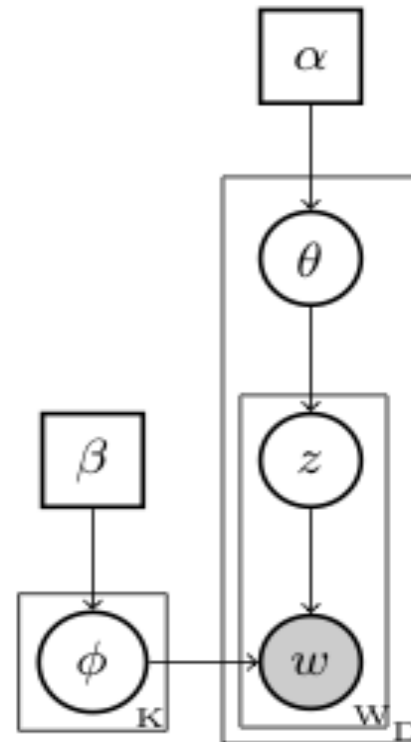


Figure 1: Graphical model for LDA.

In A Math-ier Notation

```
# topic k, docId d, and wordId w are integer indices
#
# x[d][j] = w, index of j-th word in doc d
# z[d][j] = k, index of latent topic of j-th word in doc d
# vocab[w] = string for the word with index w
#
# totalTopicCount[k] = number of words in topic k       $N[*,k]$ 
# docTopicCount[d][k] = number of words in topic k for document d       $N[d,k]$ 
# wordTopicCount[w][k] = number of occurrences of word w in topic k
# totalWords = number of words in the corpus       $N[*,*]=V$        $M[w,k]$ 
```

for each document d and word position j in d

- $z[d,j] = k$, a random topic
- $N[d,k]++$
- $W[w,k]++$ where $w = \text{id of } j\text{-th word in } d$

```
def initGibbs(self):
    print '.initializing latent vars'
    self.totalTopicCount = self.topicCounter()
    self.docTopicCount = [self.topicCounter() for d in xrange(len(self.x))]
    self.wordTopicCount = [self.topicCounter() for w in xrange(len(self.vocab))]
    self.z = [[-1 for j in xrange(len(self.x[d]))] for d in xrange(len(self.x))]
    for d in xrange(len(self.x)):
        if (d+1)%self.dstep==0: print '..doc', d+1, 'of', len(self.x)
        for j in xrange(len(self.x[d])):
            w = self.x[d][j]
            k = random.randint(0, self.numTopics-1)
            self.z[d][j] = k
            self.docTopicCount[d].add(k, 1)
            self.wordTopicCount[w].add(k, 1)
            self.totalTopicCount.add(k, 1)
    #reasonable parameters
    self.alpha = 1.0/self.numTopics
    self.beta = 1.0/len(self.vocab)
    print "alpha:", self.alpha, "beta:", self.beta
```



```

def runGibbs(self, maxT):
    for t in xrange(maxT):
        print '. iteration', t+1, 'of', maxT
        for d in xrange(len(self.x)):
            if (d+1)%self.dstep==0: print '.. doc', d+1, 'of', len(self.x)
            for j in xrange(len(self.x[d])):
                k = self.resample(d, j)
                self.flip(d, j, self.z[d][j], k)

```

for each pass $t=1,2,\dots$

for each document d and word position j in d

- $z[d,j] = k$, a new random topic

- update N, W to reflect the new assignment of z :

- $N[d,k]++; N[d,k']--$ where k' is old $z[d,j]$
- $W[w,k]++; W[w,k']--$ where w is $w[d,j]$

```

def flip(self, d, j, k_old, k_new):
    """update counts to reflect a changed value of z[d][j]"""
    if k_old != k_new:
        w = self.x[d][j]
        self.docTopicCount[d].add(k_old, -1)
        self.wordTopicCount[w].add(k_old, -1)
        self.wordTopicCount[w].add(k_new, +1)
        self.totalTopicCount.add(k_old, -1)
        self.totalTopicCount.add(k_new, +1)
        self.z[d][j] = k_new

```

$$p(Z_{d,j} = k | ..) \propto \Pr(Z_{d,j} = k | "d") * \Pr(W_{d,k} = w | Z_{d,j} = k, ...)$$

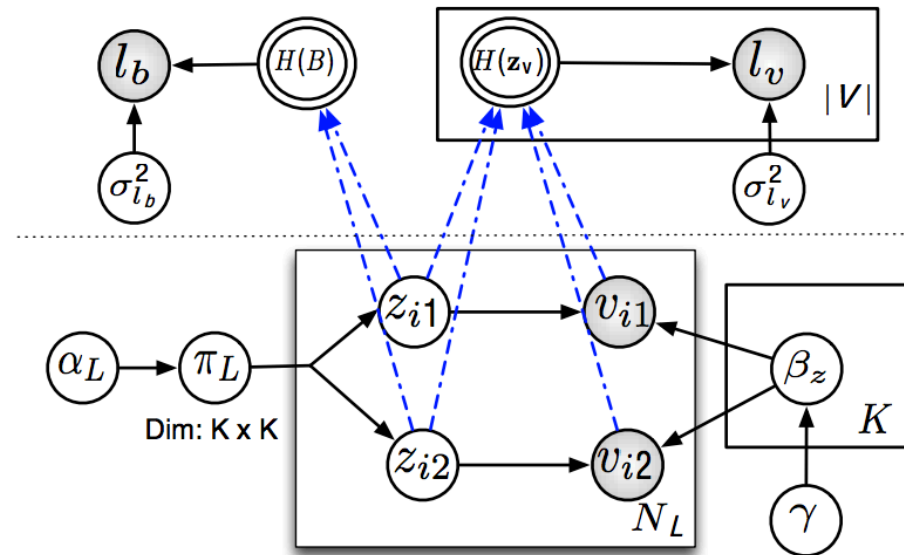
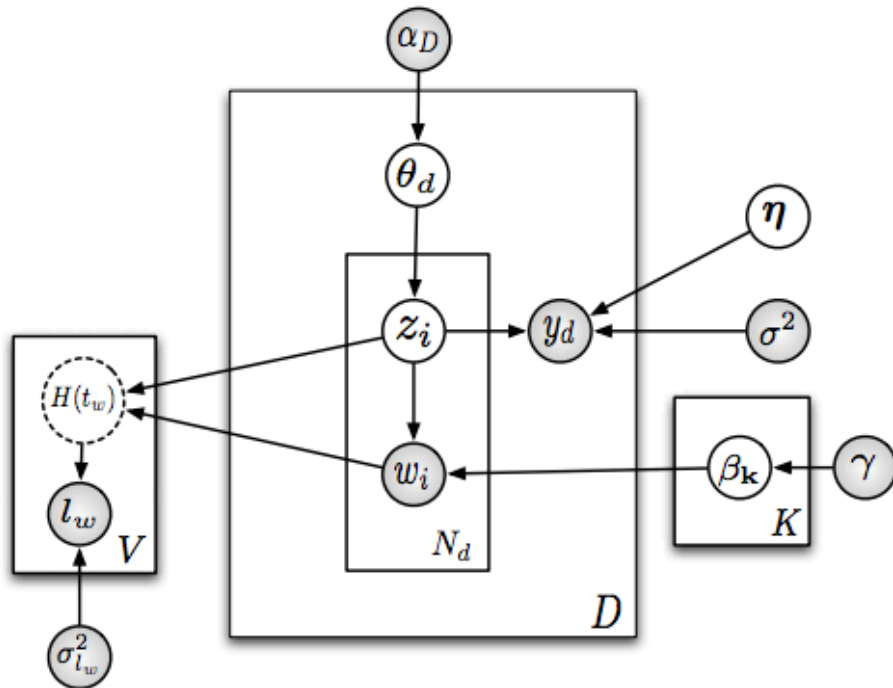
$$= \frac{N[k, d] - C_{d,j,k} + \alpha}{Z} \cdot \frac{W[w, k] - C_{d,j,k} + \beta}{(W[* , k] - C_{d,j,k}) + \beta N[* , *]}$$

```
def resample(self, d, j):
    """sample a new value of z[d][j]"""
    p = []
    norm = 0.0
    #compute pk = Pr(z_dj=k | everything else)
    for k in xrange(self.numTopics):
        w = self.x[d][j]
        z_dj_equals_k = 1 if self.z[d][j]==k else 0
        #unnormalized chance of picking topic k in doc d
        ak = (self.docTopicCount[d][k] - z_dj_equals_k + self.alpha)
        #unnormalized chance of picking topic k for word w
        bk = ((self.wordTopicCount[w][k] - z_dj_equals_k + self.beta)
              / (self.totalTopicCount[k] - z_dj_equals_k + self.totalWords * self.beta))
        pk = ak*bk
        p.append(pk)
        norm += pk
    #pick randomly from the normalized pk
    sum_p_up_to_k = 0.0
    r = random.random()
    for k in xrange(self.numTopics):
        sum_p_up_to_k += p[k]/norm
        if r < sum_p_up_to_k:
            return k
```

$$C_{d,j,k} = \begin{cases} 1 & Z_{d,j} = k \\ 0 & \text{else} \end{cases}$$

Some comments on LDA

- Very widely used model
- Also a component of many other models

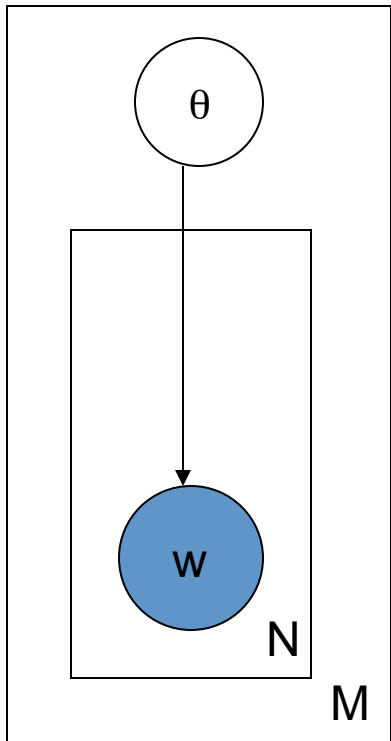


Techniques for Modeling Graphs

Latent Dirichlet Allocation++

Review - LDA

- Motivation



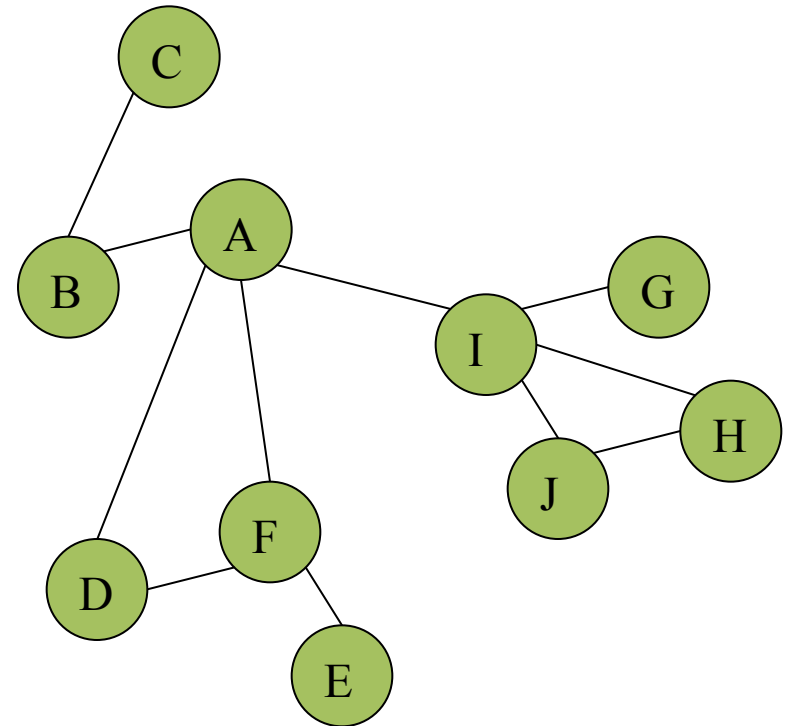
Assumptions: 1) documents are i.i.d 2) *within* a document, words are i.i.d. (bag of words)

- For each document $d = 1, \dots, M$
 - Generate $\theta_d \sim D_1(\dots)$
 - For each word $n = 1, \dots, N_d$
 - generate $w_n \sim D_2(\phi \mid \vartheta_{d_n})$

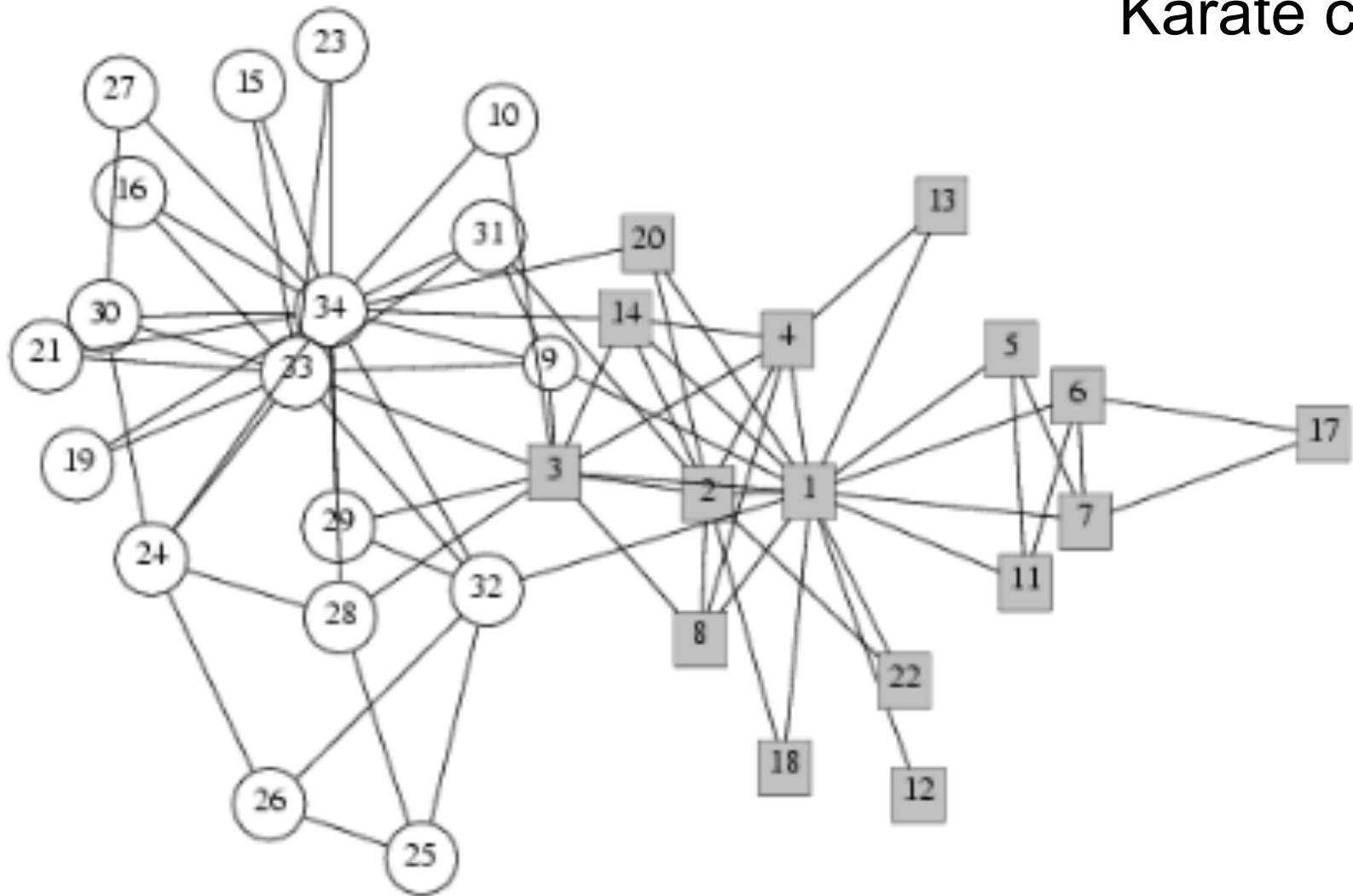
Docs and words are *exchangeable*.

A Matrix Can Also Represent a Graph

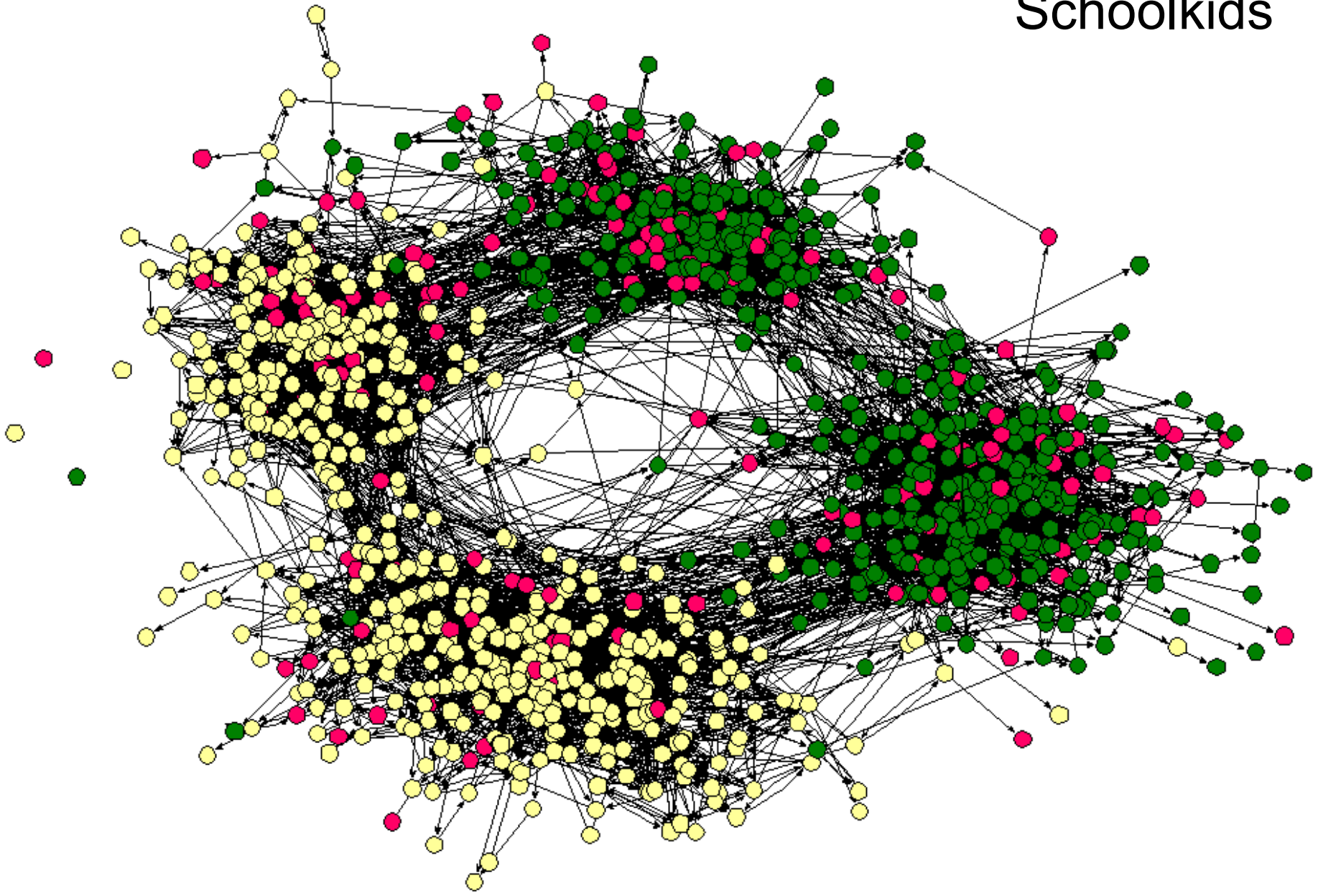
	A	B	C	D	E	F	G	H	I	J
A		1		1		1				
B	1		1							
C		1								
D	1					1				
E						1				
F	1			1	1					
G									1	
H							1		1	1
I							1	1		1
J								1	1	



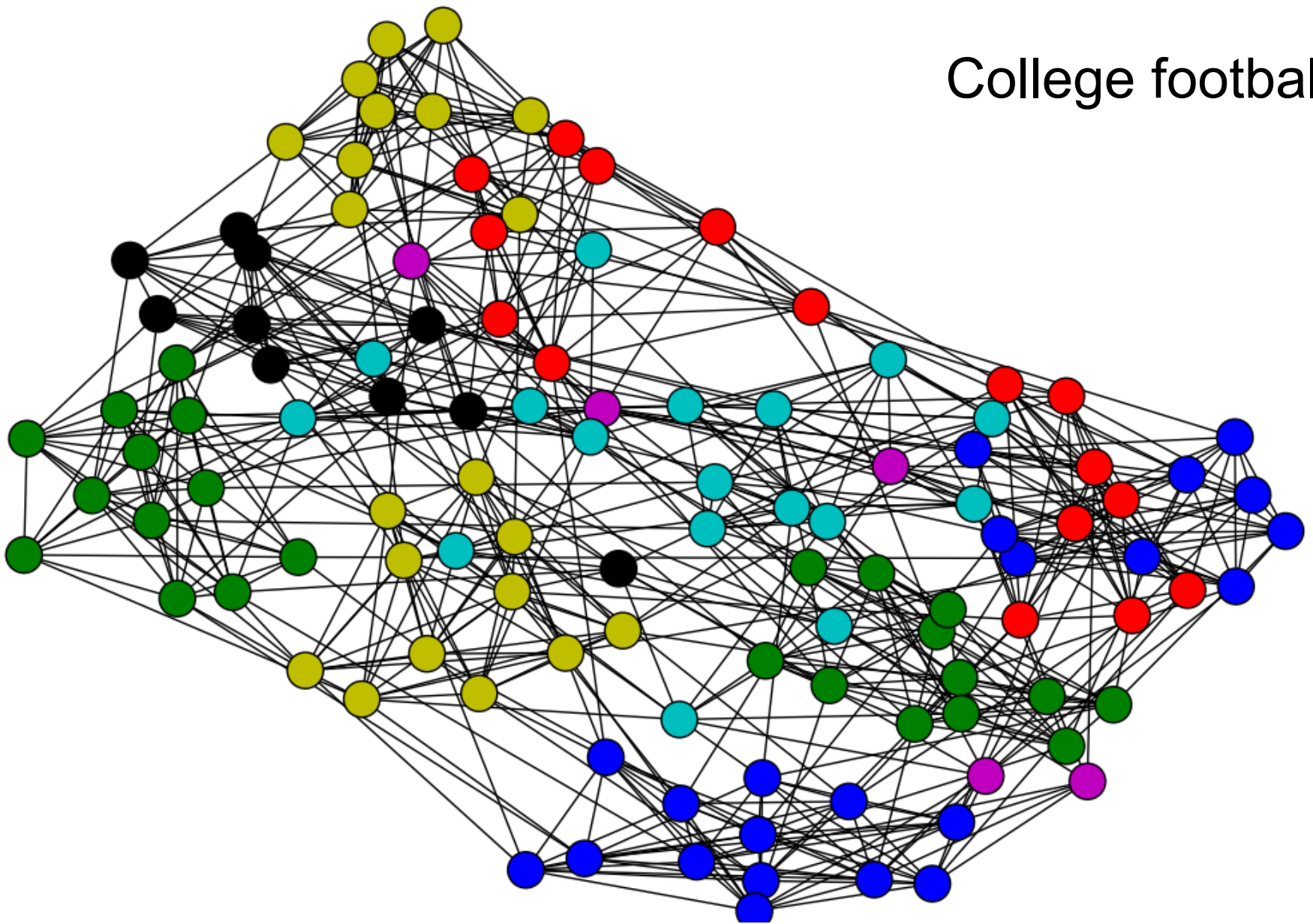
Karate club

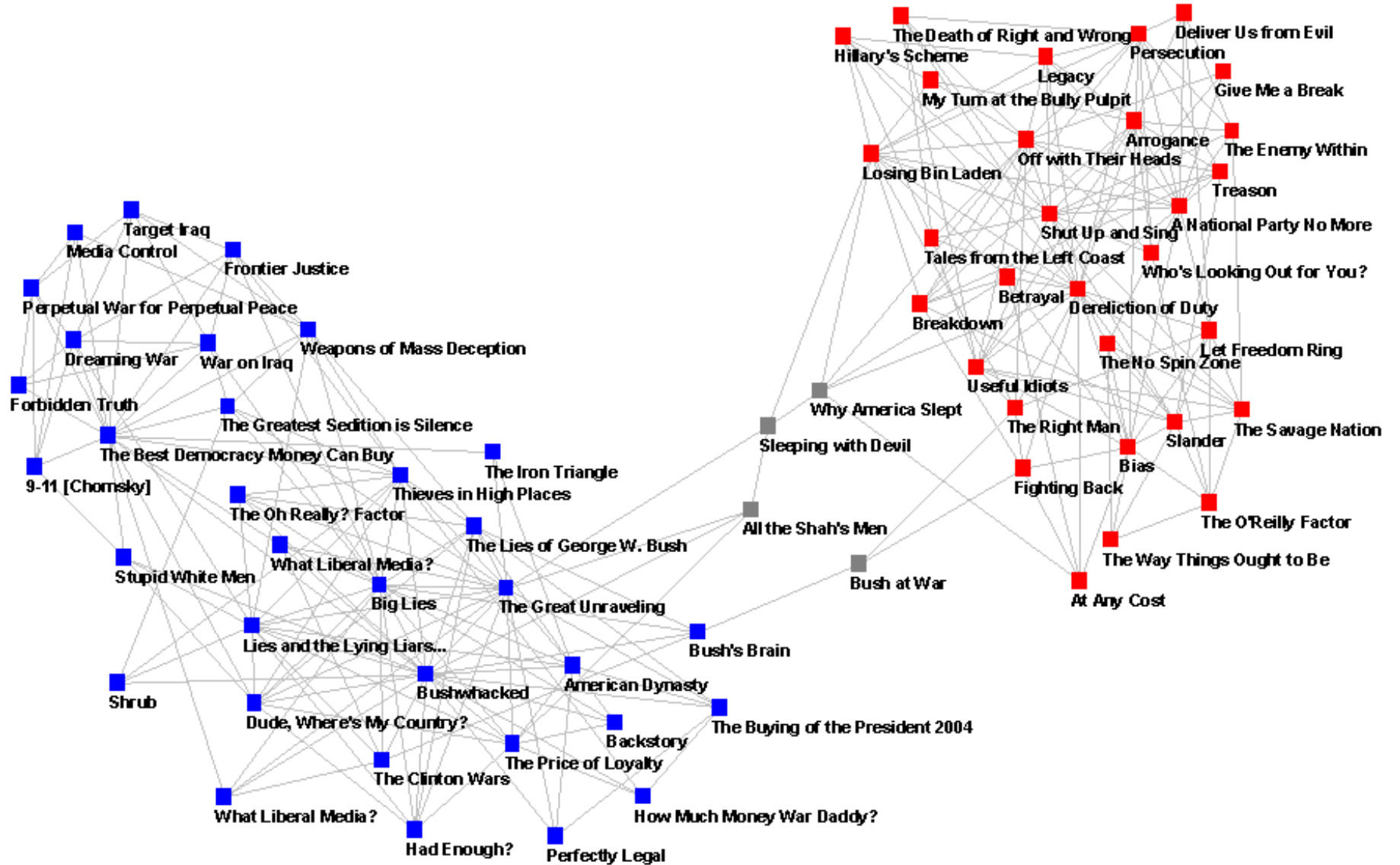


Schoolkids

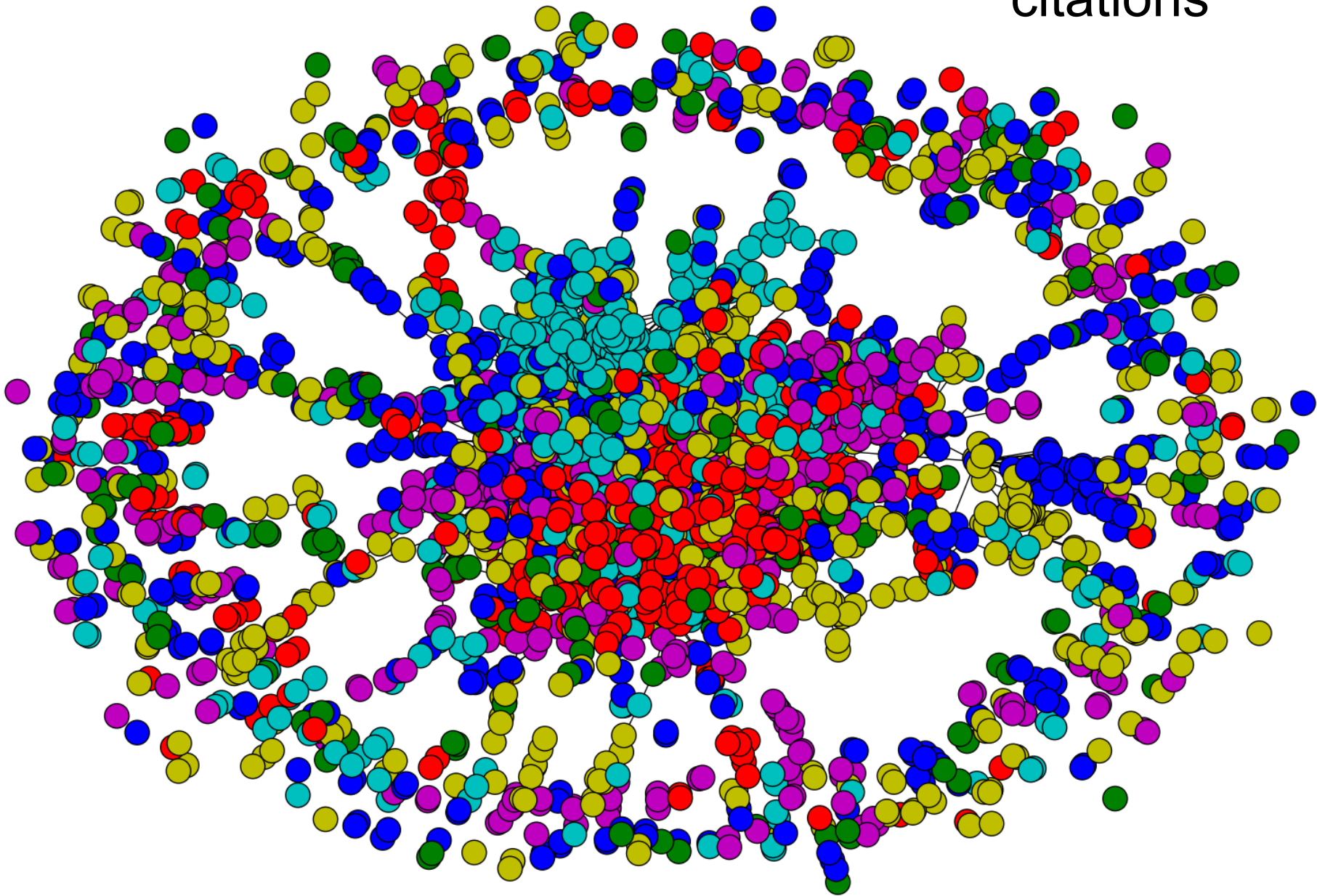


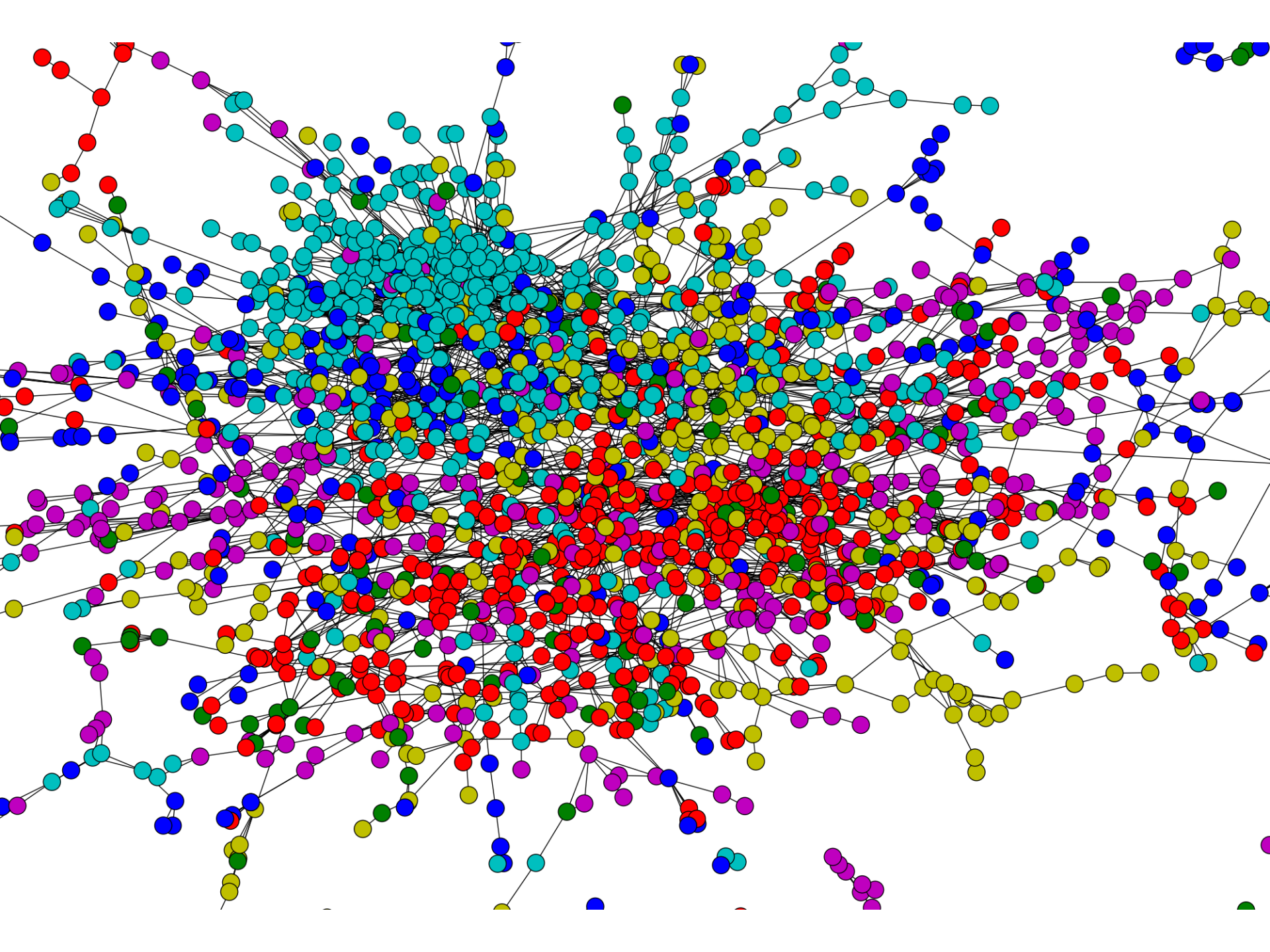
College football

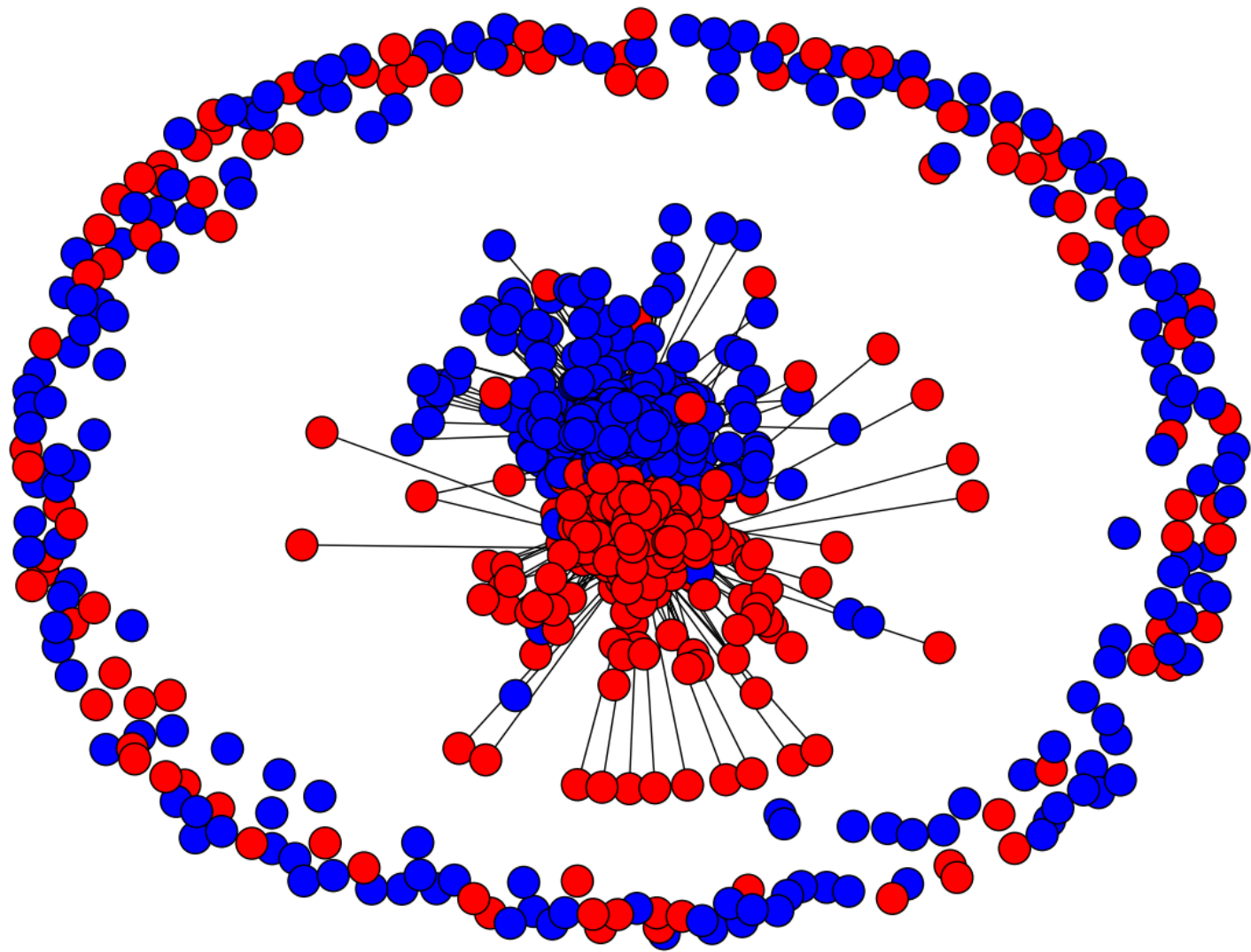


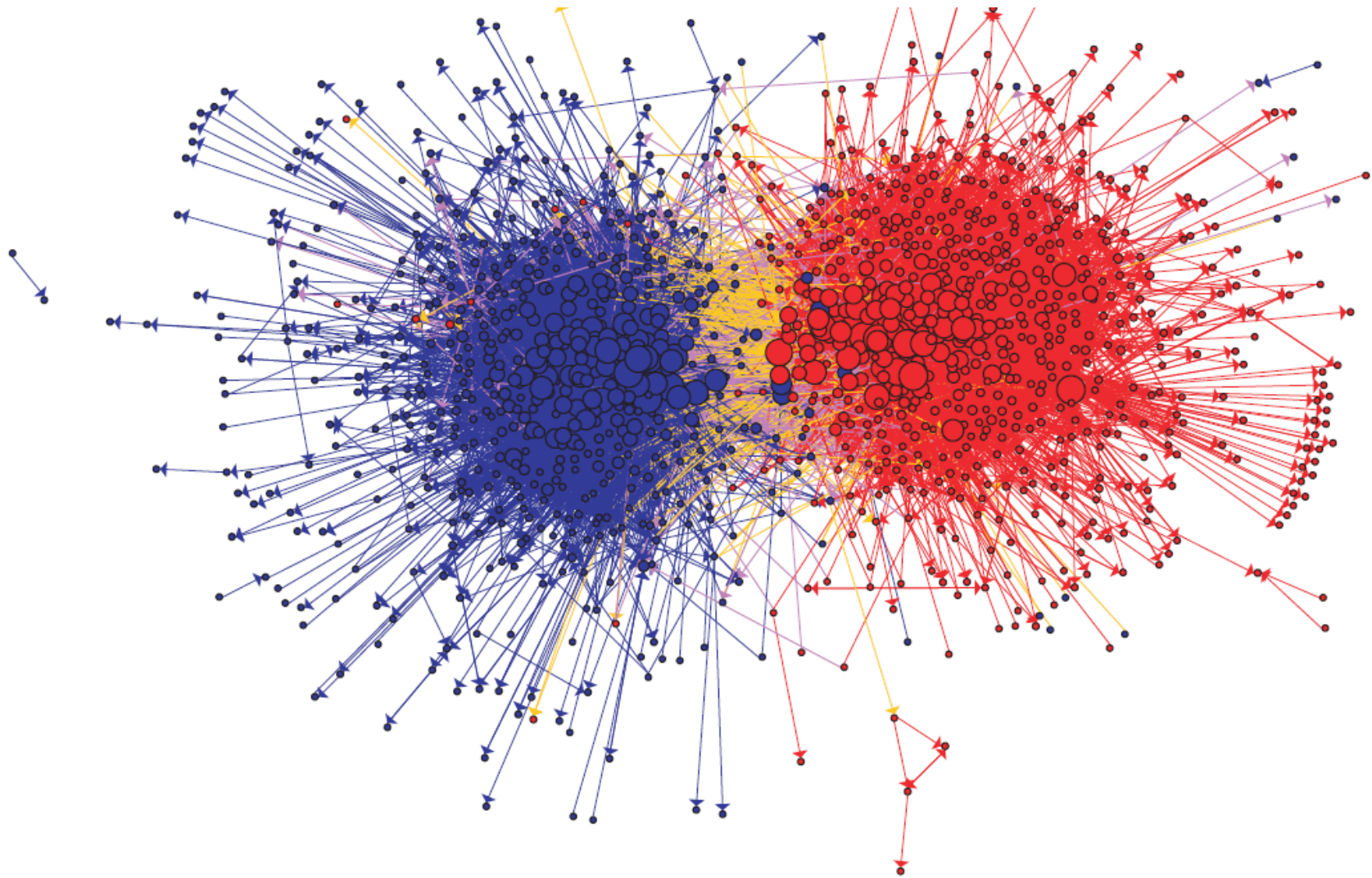


citations



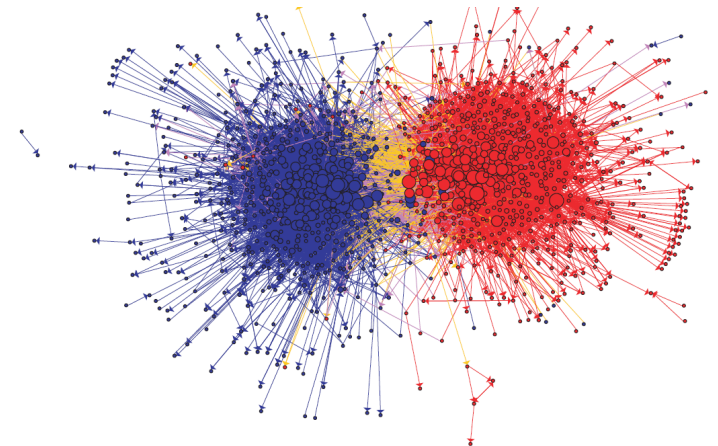






Stochastic Block Models

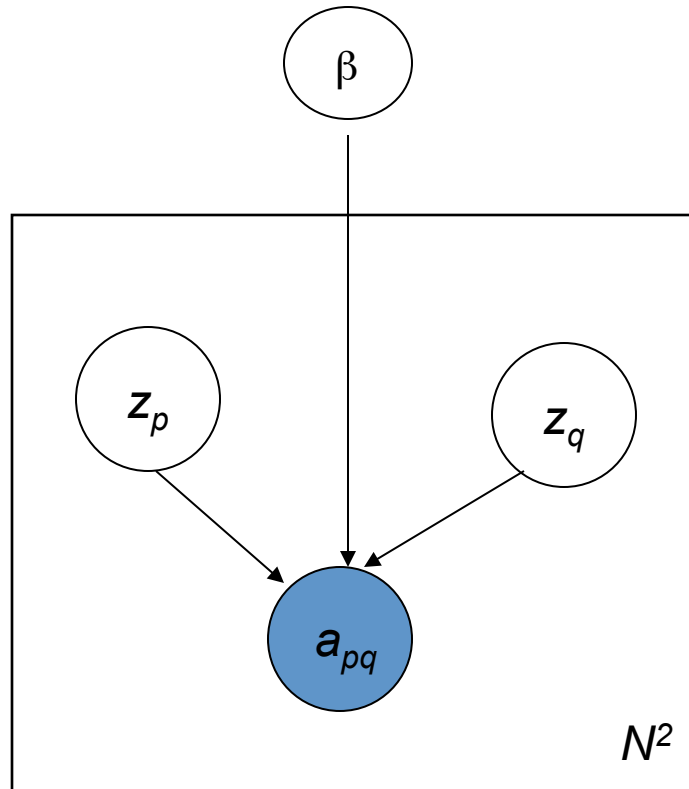
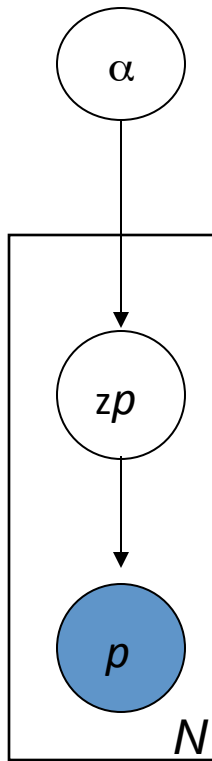
- “Stochastic block model”, aka “Block-stochastic matrix”:
 - Draw n_i nodes in block i
 - With probability p_{ij} , connect pairs (u,v) where u is in block i , v is in block j
 - Special, simple case: $p_{ii}=q_i$ and $p_{ij}=s$ for all $i \neq j$



Stochastic Block models:

assume 1) nodes w/in a block z and

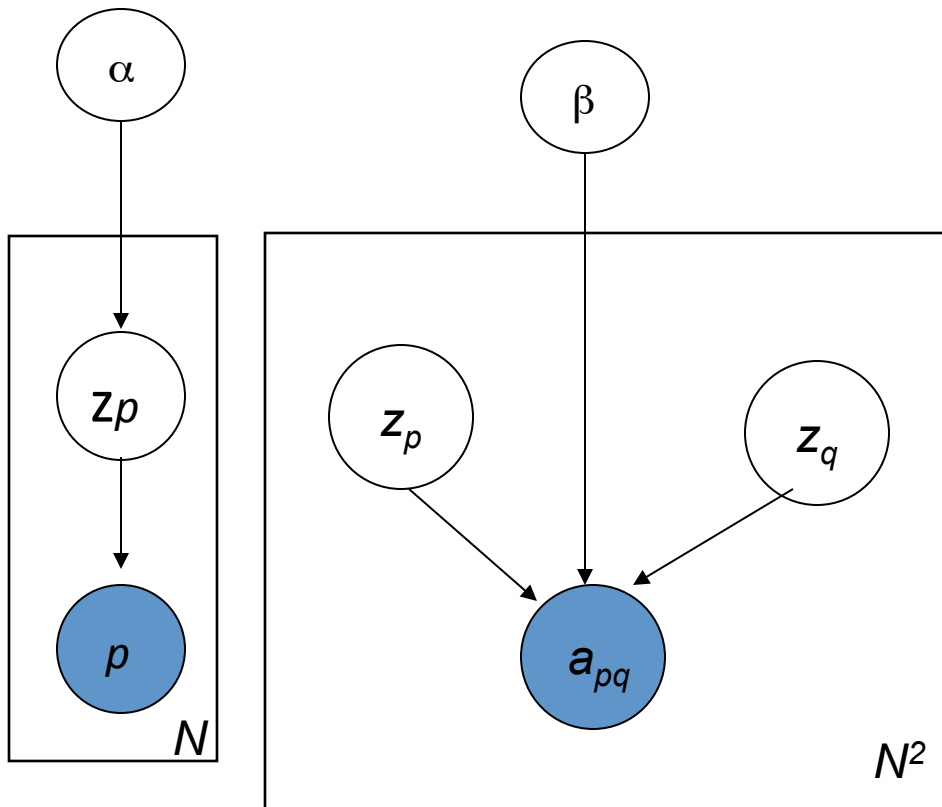
2) edges between blocks z_p, z_q are *exchangeable*



Stochastic Block models:

assume 1) nodes w/in a block z and

2) edges between blocks z_p, z_q are *exchangeable*

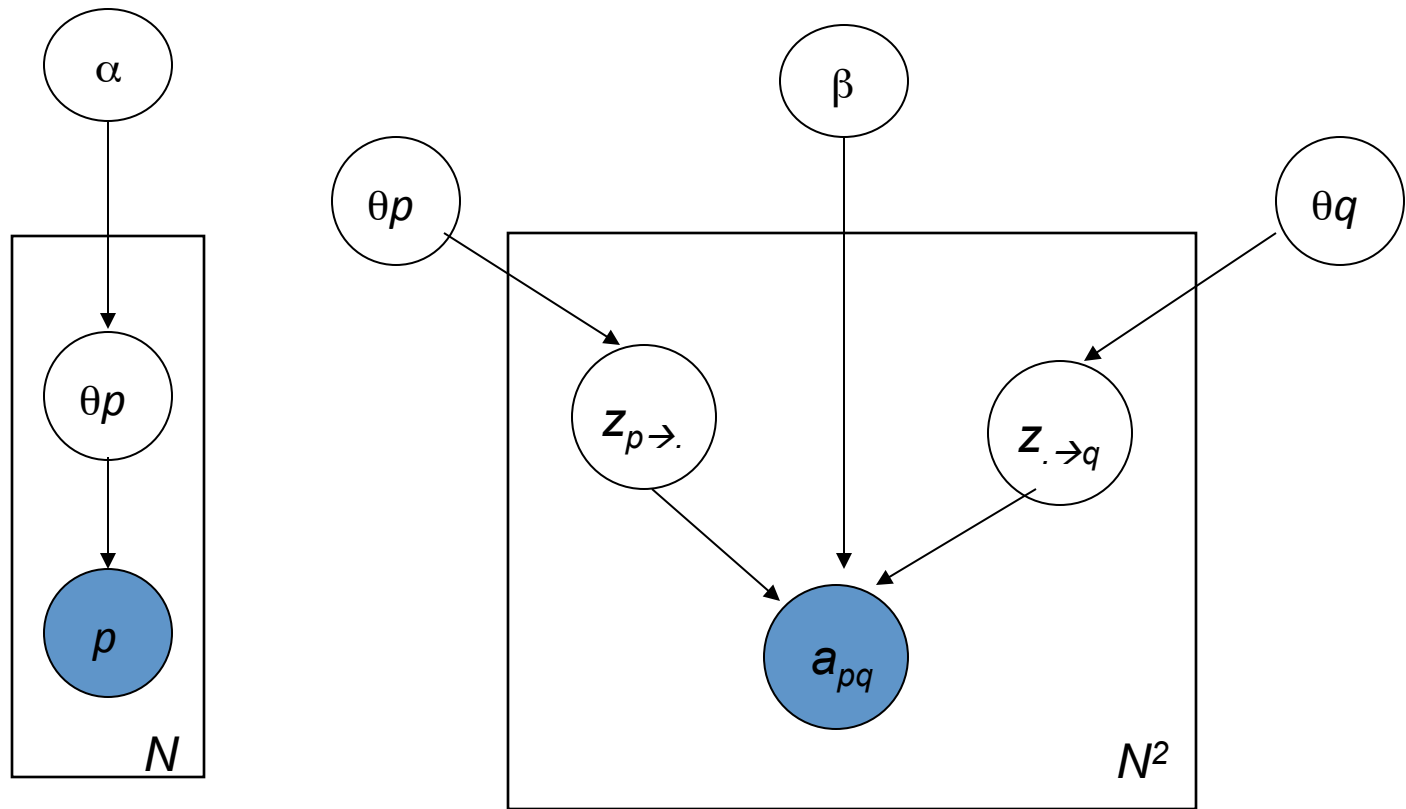


Gibbs sampling:

- Randomly initialize z_p for each node p .
- For $t = 1 \dots$
 - For each node p
 - Compute z_p given other z' s
 - Sample z_p

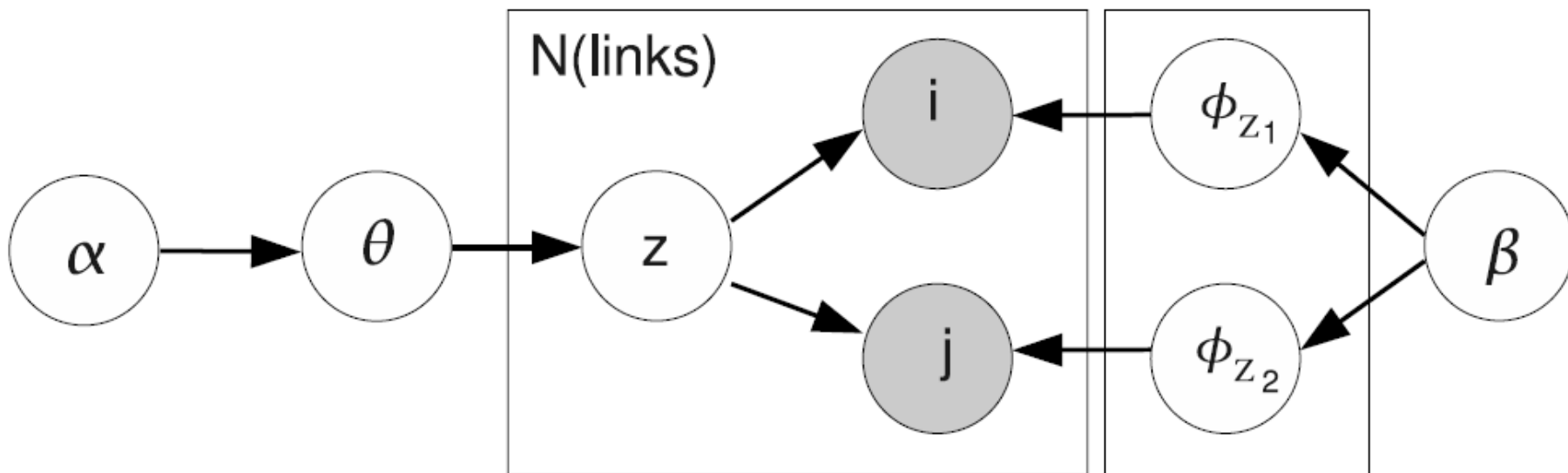
See: Snijders & Nowicki, 1997, Estimation and Prediction for Stochastic Blockmodels for Groups with Latent Graph Structure

Mixed Membership Stochastic Block models



Airoldi et al, JMLR 2008

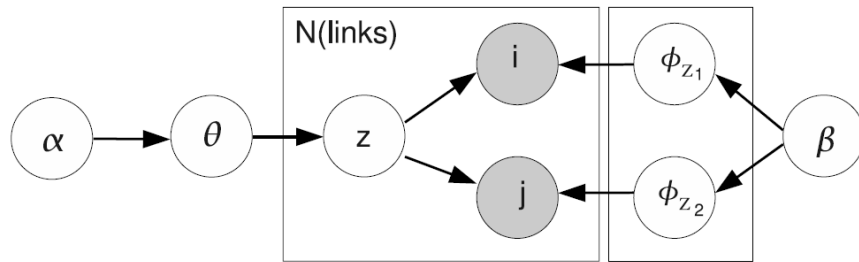
Another mixed membership block model



$$p(z_l | \{z\}^{-l}, \{(i, j)\}^{-l}, \alpha, \beta) \propto$$

$$(n_z^{-l} + \alpha) \cdot \frac{(q_{z_1 i}^{-l} + \beta)(q_{z_2 j}^{-l} + \beta)}{(q_{z_1 \cdot}^{-l} + M\beta)(q_{z_2 \cdot}^{-l} + M\beta + \delta_z)},$$

Another mixed membership block model



$z=(z_i, z_j)$ is a pair of block ids

$n_z = \#pairs\ z$

$q_{z_1, i} = \#links\ to\ i\ from\ block\ z_1$

$q_{z_1, \cdot} = \#outlinks\ in\ block\ z_1$

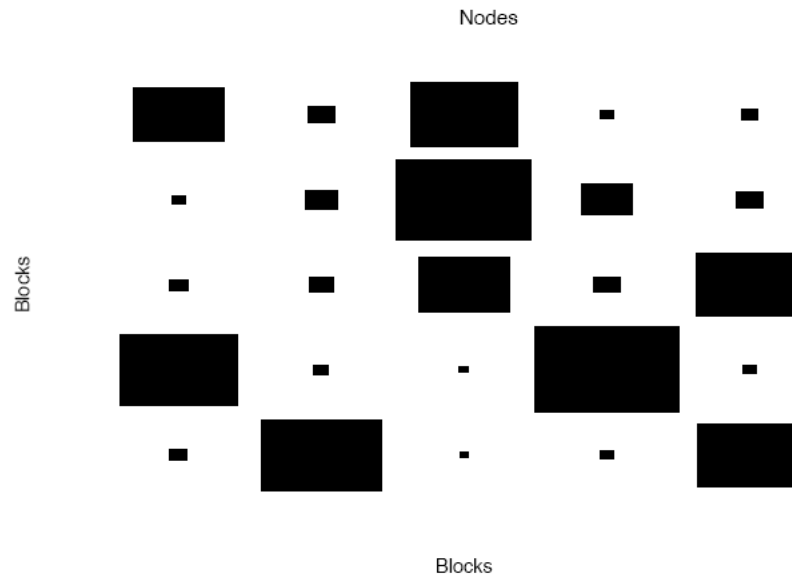
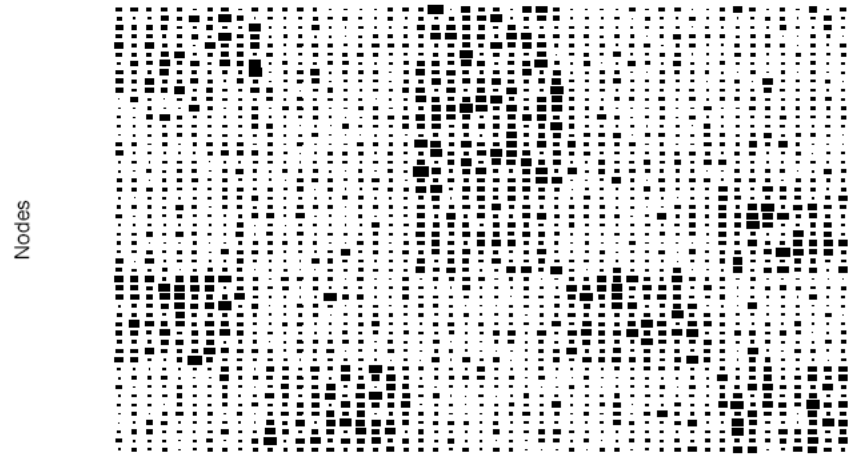
$\delta = indicator\ for\ diagonal$

$M = \#nodes$

$$p(z_l | \{z\}^{-l}, \{(i, j)\}^{-l}, \alpha, \beta) \propto$$

$$(n_z^{-l} + \alpha) \cdot \frac{(q_{z_1 i}^{-l} + \beta)(q_{z_2 j}^{-l} + \beta)}{(q_{z_1 \cdot}^{-l} + M\beta)(q_{z_2 \cdot}^{-l} + M\beta + \delta_z)},$$

Another mixed membership block model



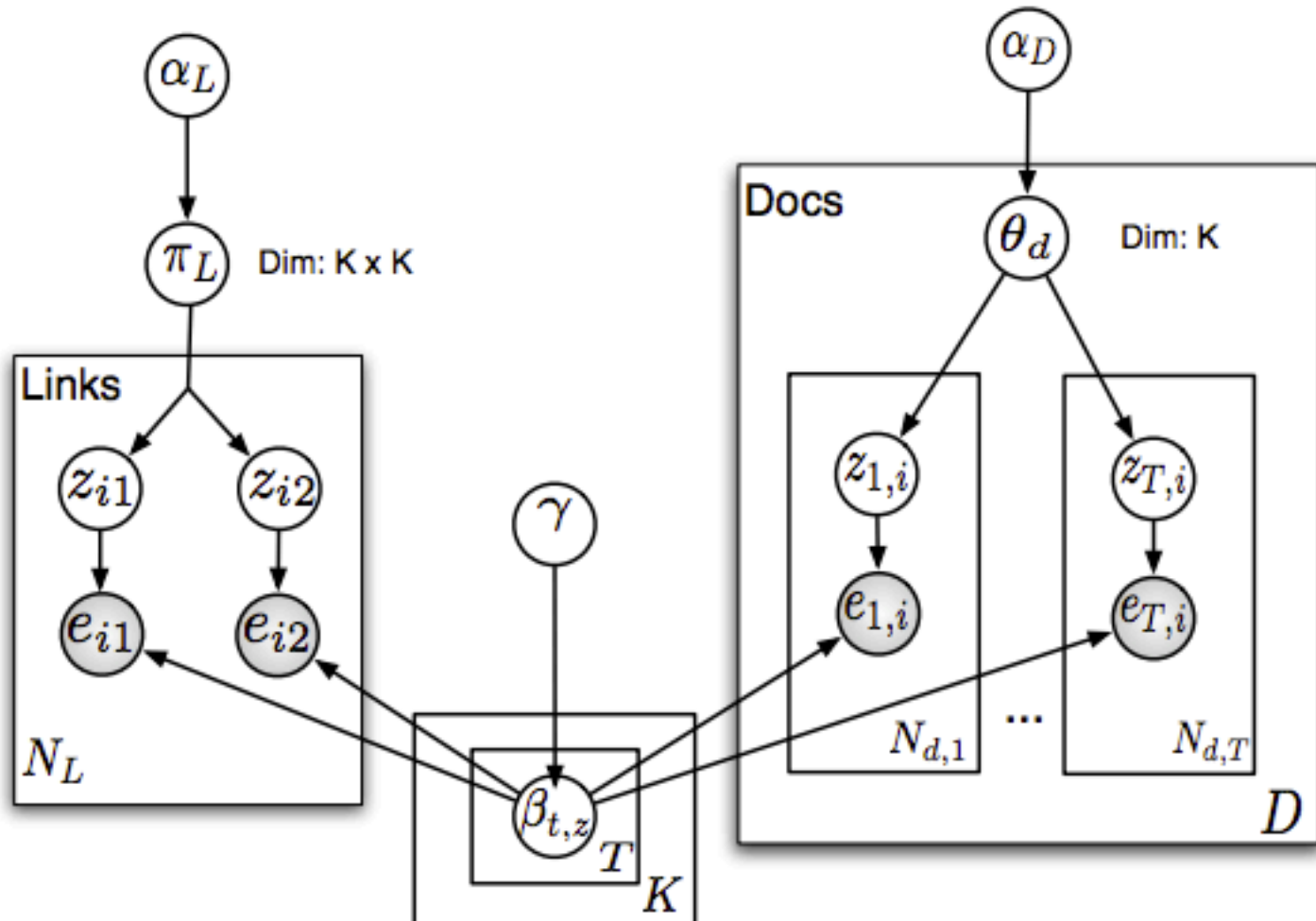
Sparse network model vs spectral clustering

(c) Best alignment: Social networks

Dataset	PSK	NCut	NJW
Karate	1.00	0.95	0.95
Dolphin	0.90	0.98	0.98
UMBC	0.95	0.95	0.96
AG	0.95	0.52	0.51
MSP	0.88	0.63	0.64
Senate	0.98	0.99	0.99
PolBook	0.78	0.82	0.80
Football	0.76	0.72	0.67
MGE-mail	0.28	0.59	0.56
CiteSeer	0.33	0.48	0.52
Cora	0.47	0.29	0.42
Average	0.75	0.72	0.73

Hybrid models

Balasubramanyan & Cohen, BioNLP 2012



Topic	Top Words & Proteins
<p>Protein Structure & Interactions</p> <p>(Publications Only)</p>	<p>Words: protein structure binding residues domain structural beta complex atp proteins alpha interactions folding structures form terminal peptide helix model interaction bound domains molecular changes conformational</p> <p>Proteins: CYC1 SSA1 HSP82 SUP35 HSP104 HSC82 SSA2 YDJ1 URE2 KAR2 SSB1 SSA4 GCN4 SSA3 SSB2 PGK1 PDI1 SSC1 HSP60 STI1 SIS1 RNQ1 SEC61 SSE1 CCP1</p>
<p>DNA Repair</p> <p>(Using MIPS PPI)</p>	<p>Words:dna recombination repair replication strand single double cells mutations stranded induced base uv mutants mutation homologous virus telomere human type yeast activity telomerase mutant dna_polymerase</p> <p>Proteins: RAD52 RAD51 RAD50 MRE11 RAD1 RAD54 SGS1 MSH2 RAD6 YKU70 REV3 POL30 RAD3 XRS2 RAD18 RAD2 POL3 RAD27 YKU80 RAD9 RFA1 TLC1 TEL1 EST2 HO</p>
<p>Vesicular Transport</p> <p>(Using Wetlab PPI)</p>	<p>Words:membrane protein transport proteins atp golgi er atpase membranes plasma_membrane vesicles cells endoplasmic_reticulum complex fusion ca2 dependent translocation vacuolar intracellular yeast lipid channel hsp90 vesicle</p> <p>Proteins: SSA1 HSP82 KAR2 PMA1 HSC82 SEC18 SSA2 YDJ1 SEC61 PEP4 HSP104 SEC23 VAM3 IRE1 SEC4 SSA4 SEC1 PMR1 PEP12 VMA3 VPH1 SSB1 VMA1 SAR1 HAC1</p>

