Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

# Parametric Bayesian Models: Part II

## Mingyuan Zhou and Lizhen Lin

Department of Information, Risk, and Operations Management
Department of Statistics and Data Sciences
The University of Texas at Austin

Machine Learning Summer School, Austin, TX
January 08, 2015

Parametric
Bayesian
Models: Part
II

Mingyuan
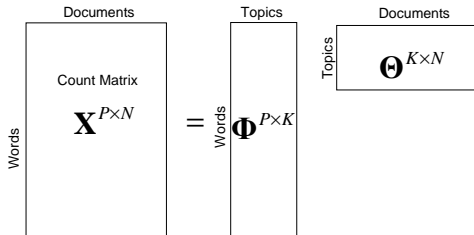Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Outline for Part II

- Bayesian modeling of count data
  - Poisson, gamma, and negative binomial distributions
  - Bayesian inference for the negative binomial distribution
  - Regression analysis for counts
- Latent variable models for discrete data
  - Latent Dirichlet allocation
  - Poisson factor analysis



- Relational network analysis

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Count data is common

- Nonnegative and discrete:
    - Number of auto insurance claims / highway accidents / crimes
    - Consumer behavior, labor mobility, marketing, voting
    - Photon counting
    - Species sampling
    - Text analysis
    - Infectious diseases, Google Flu Trends
    - Next generation sequencing (statistical genomics)
- Mixture modeling can be viewed as a count-modeling problem
    - Number of points in a cluster (mixture model, we are modeling a count vector)
    - Number of words assigned to topic $k$ in document $j$ (we are modeling a $K \times J$ latent count matrix in a topic model/mixed-membership model)

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Count data is common

- Nonnegative and discrete:
    - Number of auto insurance claims / highway accidents / crimes
    - Consumer behavior, labor mobility, marketing, voting
    - Photon counting
    - Species sampling
    - Text analysis
    - Infectious diseases, Google Flu Trends
    - Next generation sequencing (statistical genomics)
- Mixture modeling can be viewed as a count-modeling problem
    - Number of points in a cluster (mixture model, we are modeling a count vector)
    - Number of words assigned to topic $k$ in document $j$ (we are modeling a $K \times J$ latent count matrix in a topic model/mixed-membership model)

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Poisson distribution

**Siméon-Denis Poisson**
(21 June 1781 – 25 April 1840)

"Life is good for only two things:
doing mathematics and teaching it."



http://en.wikipedia.org

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline
Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

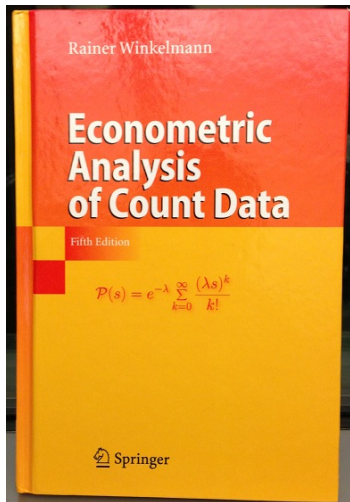Main
references

# Poisson distribution

**Siméon-Denis Poisson**
(21 June 1781 – 25 April 1840)

"Life is good for only two things:
doing mathematics and teaching it."

http://en.wikipedia.org

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

- Poisson distribution $x \sim \text{Pois}(\lambda)$
  - Probability mass function:

  $$P(x|\lambda) = \frac{\lambda^x e^{-\lambda}}{x!}, \quad x \in \{0, 1, \dots\}$$

  - The mean and variance is the same: $\mathbb{E}[x] = \text{Var}[x] = \lambda$.
  - Restrictive to model over-dispersed (variance greater than the mean) counts that are commonly observed in practice.
  - A basic building block to construct more flexible count distributions.
- Overdispersed count data are commonly observed due to
  - Heterogeneity: difference between individuals
  - Contagion: dependence between the occurrence of events

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Mixed Poisson distribution

$$x \sim \text{Pois}(\lambda), \ \lambda \sim f_\Lambda(\lambda)$$

- Mixing the Poisson rate parameter with a positive distribution leads to a mixed Poisson distribution.
- A mixed Poisson distribution is always over-dispersed.
  - Law of total expectation:

    $$\mathbb{E}[x] = \mathbb{E}[\mathbb{E}[x|\lambda]] = \mathbb{E}[\lambda].$$

  - Law of total variance:

    $$\text{Var}[x] = \text{Var}[\mathbb{E}[x|\lambda]] + \mathbb{E}[\text{Var}[x|\lambda]] = \text{Var}[\lambda] + \mathbb{E}[\lambda].$$

  - Thus $\text{Var}[x] > \mathbb{E}[x]$ unless $\lambda$ is a constant.
- The gamma distribution is a popular choice as it is conjugate to the Poisson distribution.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

- Mixing the gamma distribution with the Poisson distribution as

$$x \sim \text{Pois}(\lambda), \ \lambda \sim \text{Gamma}\left(r, \frac{p}{1-p}\right),$$

where $p/(1-p)$ is the gamma scale parameter, leads to the negative binomial distribution $x \sim \text{NB}(r, p)$ with probability mass function

$$P(x|r, p) = \frac{\Gamma(x+r)}{x!\Gamma(r)} p^x (1-p)^r, \quad x \in \{0, 1, \ldots\}$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
**Count
distributions**
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Compound Poisson distribution

- A compound Poisson distribution is the summation of a Poisson random number of *i.i.d.* random variables.
- If $x = \sum_{i=1}^{n} y_i$, where $n \sim \text{Pois}(\lambda)$ and $y_i$ are *i.i.d.* random variable, then $x$ is a compound Poisson random variable.
- The negative binomial random variable $x \sim \text{NB}(r, p)$ can also be generated as a compound Poisson random variable as

$$x = \sum_{i=1}^{l} u_i, \ l \sim \text{Pois}[-r \ln(1 - p)], \ u_i \sim \text{Log}(p)$$

where $u \sim \text{Log}(p)$ is the logarithmic distribution with probability mass function

$$P(u|p) = \frac{-1}{\ln(1 - p)} \frac{p^u}{u}, \quad u \in \{1, 2, \cdots\}.$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
**Negative
binomial
distribution**
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Negative binomial distribution

$$m \sim \text{NB}(r, p)$$

- $r$ is the dispersion parameter
- $p$ is the probability parameter
- Probability mass function

$$f_M(m|r, p) = \frac{\Gamma(r + m)}{m!\Gamma(r)} p^m (1 - p)^r = (-1)^m \binom{-r}{m} p^m (1 - p)^r$$

- It is a gamma-Poisson mixture distribution
- It is a compound Poisson distribution
- Its variance $\frac{rp}{(1-p)^2}$ is greater that its mean $\frac{rp}{1-p}$
- $\text{Var}[m] = \mathbb{E}[m] + \frac{(\mathbb{E}[m])^2}{r}$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
**Negative
binomial
distribution**
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

- The conjugate prior for the negative binomial probability parameter $p$ is the beta distribution: if $m_i \sim \text{NB}(r, p)$, $p \sim \text{Beta}(a_0, b_0)$, then

$$(p|-) = \text{Beta}\left(a_0 + \sum_{i=1}^{n} m_i, b_0 + nr\right)$$

- The conjugate prior for the negative binomial dispersion parameter $r$ is unknown, but we have a simple data augmentation technique to derive closed-form Gibbs sampling update equations for $r$.

- If we assign $m$ customers to tables using a Chinese restaurant process with concentration parameter $r$, then the random number of occupied tables $l$ follows the Chinese Restaurant Table (CRT) distribution

$$f_L(l|m, r) = \frac{\Gamma(r)}{\Gamma(m+r)}|s(m, l)|r^l, \quad l = 0, 1, \cdots, m.$$

$|s(m, l)|$ are unsigned Stirling numbers of the first kind.

- The joint distribution of the customer count $m \sim \text{NB}(r, p)$ and table count is the Poisson-logarithmic bivariate count distribution

$$f_{M,L}(m, l|r, p) = \frac{|s(m, l)|r^l}{m!}(1 - p)^r p^m.$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
**Negative
binomial
distribution**
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Poisson-logarithmic bivariate count distribution

- Probability mass function:

$$f_{M,L}(m, l; r, p) = \frac{|s(m, l)| r^l}{m!}(1 - p)^r p^m.$$

- It is clear that the gamma distribution is a conjugate prior for $r$ to this bivariate count distribution.

The joint distribution of the customer count and table count are equivalent:

Draw NegBino($r$, $p$) customers

Draw Poisson($-r \ln (1 - p)$) tables

Assign customers to tables using a Chinese restaurant process with concentration parameter $r$

Draw Logarithmic($p$) customers on each table

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Bayesian inference for the negative binomial distribution

Negative binomial count modeling:

$$m_i \sim \text{NegBino}(r, p), \ p \sim \text{Beta}(a_0, b_0), \ r \sim \text{Gamma}(e_0, 1/f_0).$$

- Gibbs sampling via data augmetantion:

$$(p|-) \sim \text{Beta}\left(a_0 + \sum_{i=1}^{n} m_i, b_0 + nr\right);$$

$$(\ell_i|-) = \sum_{t=1}^{m_i} b_t, \ b_t \sim \text{Bernoulli}\left(\frac{r}{t+r-1}\right);$$

$$(r|-) \sim \text{Gamma}\left(e_0 + \sum_{i=1}^{n} \ell_i, \frac{1}{f_0 - n\ln(1-p)}\right).$$

- Expectation-Maximization

- Variational Bayes

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Bayesian inference for the negative binomial distribution

Negative binomial count modeling:

$$m_i \sim \text{NegBino}(r, p), \ p \sim \text{Beta}(a_0, b_0), \ r \sim \text{Gamma}(e_0, 1/f_0).$$

- Gibbs sampling via data augmetantion:

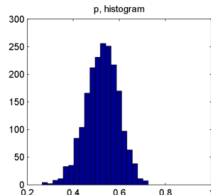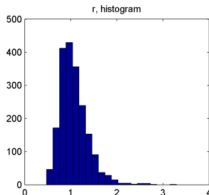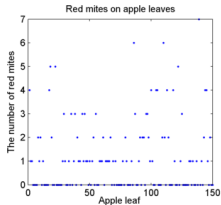$$(p|-) \sim \text{Beta}\left(a_0 + \sum_{i=1}^n m_i, b_0 + nr\right);$$

$$(\ell_i|-) = \sum_{t=1}^{m_i} b_t, \ b_t \sim \text{Bernoulli}\left(\frac{r}{t+r-1}\right);$$

$$(r|-) \sim \text{Gamma}\left(e_0 + \sum_{i=1}^n \ell_i, \frac{1}{f_0 - n\ln(1-p)}\right).$$

- Expectation-Maximization
- Variational Bayes

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
**Negative
binomial
distribution**
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

- Gibbs sampling: $\mathbb{E}[r] = 1.076$, $\mathbb{E}[p] = 0.525$.
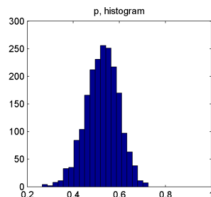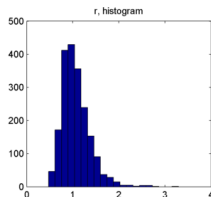


- Expectation-Maximization: $r : 1.025$, $p : 0.528$.

- Variational Bayes: $\mathbb{E}[r] = 0.999$, $\mathbb{E}[p] = 0.534$.

- For this example, variational Bayes inference correctly identifies the modes but underestimates the posterior variances of model parameters.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
**Negative
binomial
distribution**
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

- Gibbs sampling: $\mathbb{E}[r] = 1.076$, $\mathbb{E}[p] = 0.525$.
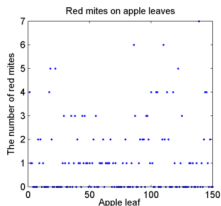


- Expectation-Maximization: $r : 1.025$, $p : 0.528$.
- Variational Bayes: $\mathbb{E}[r] = 0.999$, $\mathbb{E}[p] = 0.534$.

- For this example, variational Bayes inference correctly identifies the modes but underestimates the posterior variances of model parameters.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
**Negative
binomial
distribution**
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

- Gibbs sampling: $\mathbb{E}[r] = 1.076$, $\mathbb{E}[p] = 0.525$.
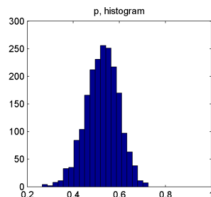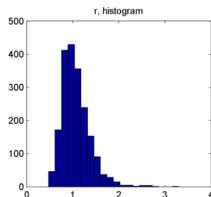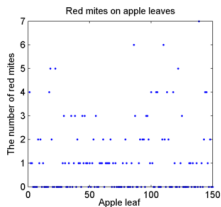
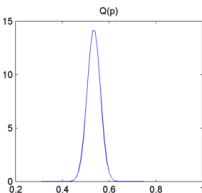

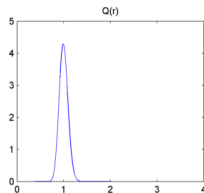- Expectation-Maximization: $r : 1.025$, $p : 0.528$.
- Variational Bayes: $\mathbb{E}[r] = 0.999$, $\mathbb{E}[p] = 0.534$.



- For this example, variational Bayes inference correctly
  identifies the modes but underestimates the posterior
  variances of model parameters.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

# Negative binomial gamma chain

NegBino-Gamma-Gamma-...

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

# Negative binomial gamma chain

NegBino-Gamma-Gamma-...

*Augmentation*

(CRT, NegBino)-Gamma-Gamma-...

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
**Negative
binomial
distribution**
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Negative binomial gamma chain

NegBino-Gamma-Gamma-...

*Augmentation*

(CRT, NegBino)-Gamma-Gamma-...

*Equivalence*

(Log,  Poisson)-Gamma-Gamma-...

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

# Negative binomial gamma chain

NegBino-Gamma-Gamma-...

*Augmentation*

(CRT, NegBino)-Gamma-Gamma-...

*Equivalence*

(Log, Poisson)-Gamma-Gamma-...

*Marginalization*

NegBino-Gamma-...

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

# Negative binomial gamma chain

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline
Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Poisson and multinomial distributions

- Suppose that $x_1, \ldots, x_K$ are independent Poisson random variables with

$$x_k \sim \text{Pois}(\lambda_k), \quad x = \sum_{k=1}^{K} x_k.$$

Set $\lambda = \sum_{k=1}^{K} \lambda_k$; let $(y, y_1, \ldots, y_K)$ be random variables such that

$$y \sim \text{Pois}(\lambda), \ (y_1, \ldots, y_k)|y \sim \text{Mult}\left(y; \frac{\lambda_1}{\lambda}, \ldots, \frac{\lambda_K}{\lambda}\right).$$

Then the distribution of $\boldsymbol{x} = (x, x_1, \ldots, x_K)$ is the same as the distribution of $\boldsymbol{y} = (y, y_1, \ldots, y_K)$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Multinomial and Dirichlet distributions

- Model:

$$(x_{i1}, \ldots, x_{ik}) \sim \text{Multinomial}(n_i, p_1, \ldots, p_k),$$

$$(p_1, \ldots, p_k) \sim \text{Dirichlet}(\alpha_1, \ldots, \alpha_k) = \frac{\Gamma(\sum_{j=1}^k \alpha_j)}{\prod_{j=1}^k \Gamma(\alpha_j)} \prod_{j=1}^k p_j^{\alpha_j - 1}$$

- The conditional posterior of $(p_1, \ldots, p_k)$ is Dirichlet distributed as

$$(p_1, \ldots, p_k | -) \sim \text{Dirichlet}\left(\alpha_1 + \sum_i x_{i1}, \ldots, \alpha_k + \sum_i x_{ik}\right)$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression
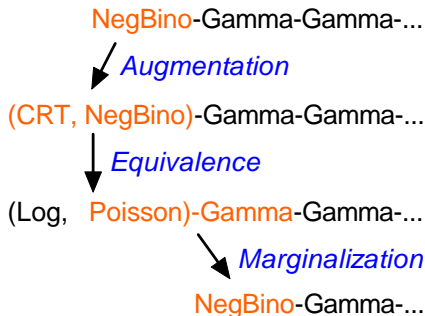
Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Gamma and Dirichlet distributions

- Suppose that random variables $y$ and $(y_1, \ldots, y_K)$ are independent with

$$y \sim \text{Gamma}(\gamma, 1/c), \quad (y_1, \ldots, y_K) \sim \text{Dir}(\gamma p_1, \cdots, \gamma p_K)$$

where $\sum_{k=1}^{K} p_k = 1$; Let

$$x_k = y y_k$$

then $\{x_k\}_{1,K}$ are independent gamma random variables with

$$x_k \sim \text{Gamma}(\gamma p_k, 1/c).$$

- The proof can be found in arXiv:1209.3442v1

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
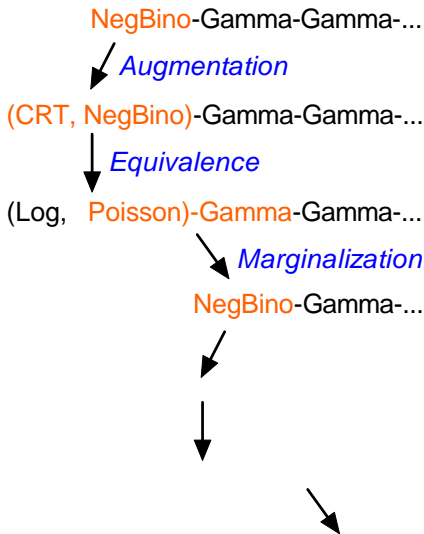binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Relationships between various distributions



Count Modeling   Mixture Modeling   Latent Gaussian

Gaussian — Logit   Logarithmic — Poisson — Multinomial

Polya-Gamma — Negative Binomial

Chinese Restaurant   Beta   Gamma — Dirichlet

Bernoulli

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
**Count regression**

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Poisson regression

- Model:

$$y_i \sim \text{Pois}(\lambda_i), \quad \lambda_i = \exp(\boldsymbol{x}_i^T \boldsymbol{\beta})$$

- Model assumption:

$$\text{Var}[y_i|\boldsymbol{x}_i] = \mathbb{E}[y_i|\boldsymbol{x}_i] = \exp(\boldsymbol{x}_i^T \boldsymbol{\beta}).$$

- Poisson regression does not model over-dispersion.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline
Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Poisson regression with multiplicative random effects

- Model:

$$y_i \sim \text{Pois}(\lambda_i), \quad \lambda_i \sim \exp(\mathbf{x}_i^T \boldsymbol{\beta}) \epsilon_i$$

- Model property:

$$\text{Var}[y_i|\mathbf{x}_i] = \mathbb{E}[y_i|\mathbf{x}_i] + \frac{\text{Var}[\epsilon_i]}{\mathbb{E}^2[\epsilon_i]} \mathbb{E}^2[y_i|\mathbf{x}_i].$$

- Negative binomial regression (gamma random effect):

$$\epsilon_i \sim \text{Gamma}(r, 1/r) = \frac{r^r}{\Gamma(r)} \epsilon_i^{r-1} e^{-r\epsilon_i}.$$

- Lognormal-Poisson regression (lognormal random effect):

$$\epsilon_i \sim \ln \mathcal{N}(0, \sigma^2).$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
**Count regression**

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Poisson regression with multiplicative random effects

- Model:

$$y_i \sim \mathsf{Pois}(\lambda_i), \quad \lambda_i \sim \exp(\boldsymbol{x}_i^T \boldsymbol{\beta}) \epsilon_i$$

- Model property:

$$\mathsf{Var}[y_i|\boldsymbol{x}_i] = \mathbb{E}[y_i|\boldsymbol{x}_i] + \frac{\mathsf{Var}[\epsilon_i]}{\mathbb{E}^2[\epsilon_i]} \mathbb{E}^2[y_i|\boldsymbol{x}_i].$$

- Negative binomial regression (gamma random effect):

$$\epsilon_i \sim \mathsf{Gamma}(r, 1/r) = \frac{r^r}{\Gamma(r)} \epsilon_i^{r-1} e^{-r\epsilon_i}.$$

- Lognormal-Poisson regression (lognormal random effect):

$$\epsilon_i \sim \ln \mathcal{N}(0, \sigma^2).$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Poisson regression with multiplicative random effects

- Model:

$$y_i \sim \text{Pois}(\lambda_i), \quad \lambda_i \sim \exp(\boldsymbol{x}_i^T \boldsymbol{\beta}) \epsilon_i$$

- Model property:

$$\text{Var}[y_i|\boldsymbol{x}_i] = \mathbb{E}[y_i|\boldsymbol{x}_i] + \frac{\text{Var}[\epsilon_i]}{\mathbb{E}^2[\epsilon_i]} \mathbb{E}^2[y_i|\boldsymbol{x}_i].$$

- Negative binomial regression (gamma random effect):

$$\epsilon_i \sim \text{Gamma}(r, 1/r) = \frac{r^r}{\Gamma(r)} \epsilon_i^{r-1} e^{-r\epsilon_i}.$$

- Lognormal-Poisson regression (lognormal random effect):

$$\epsilon_i \sim \ln \mathcal{N}(0, \sigma^2).$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Lognormal and gamma mixed negative binomial regression

- Model (Zhou et al., ICML2012):

$$y_i \sim \text{NegBino}(r, p_i), \quad r \sim \text{Gamma}(a_0, 1/h)$$

- Bayesian inference with the Polya-Gamma distribution.
- Model properties:

$$\text{Var}[y_i|x_i] = \mathbb{E}[y_i|x_i] + \left(e^{\sigma^2}(1 + r^{-1}) - 1\right)\mathbb{E}^2[y_i|x_i].$$

- Special cases:
  - Negative binomial regression: $\sigma^2 = 0$;
  - Lognormal-Poisson regression: $r \to \infty$;
  - Poisson regression: $\sigma^2 = 0$ and $r \to \infty$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Lognormal and gamma mixed negative binomial regression

- Model (Zhou et al., ICML2012):

$$y_i \sim \text{NegBino}\left(r, p_i\right), \quad r \sim \text{Gamma}(a_0, 1/h)$$

$$\psi_i = \log\left(\frac{p_i}{1 - p_i}\right) = \boldsymbol{x}_i^T \boldsymbol{\beta} + \ln \epsilon_i, \quad \epsilon_i \sim \ln \mathcal{N}(0, \sigma^2)$$

- Bayesian inference with the Polya-Gamma distribution.

- Model properties:

$$\text{Var}[y_i|\boldsymbol{x}_i] = \mathbb{E}[y_i|\boldsymbol{x}_i] + \left(e^{\sigma^2}(1 + r^{-1}) - 1\right)\mathbb{E}^2[y_i|\boldsymbol{x}_i].$$

- Special cases:

  - Negative binomial regression: $\sigma^2 = 0$;
  - Lognormal-Poisson regression: $r \to \infty$;
  - Poisson regression: $\sigma^2 = 0$ and $r \to \infty$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Lognormal and gamma mixed negative binomial regression

- Model (Zhou et al., ICML2012):

$$y_i \sim \text{NegBino}\,(r, p_i)\,, \quad r \sim \text{Gamma}(a_0, 1/h)$$

$$\psi_i = \log\left(\frac{p_i}{1 - p_i}\right) = \boldsymbol{x}_i^T \boldsymbol{\beta} + \ln \epsilon_i, \quad \epsilon_i \sim \ln \mathcal{N}(0, \sigma^2)$$

- Bayesian inference with the Polya-Gamma distribution.
- Model properties:

$$\text{Var}[y_i | \boldsymbol{x}_i] = \mathbb{E}[y_i | \boldsymbol{x}_i] + \left(e^{\sigma^2}(1 + r^{-1}) - 1\right) \mathbb{E}^2[y_i | \boldsymbol{x}_i].$$

- Special cases:
  - Negative binomial regression: $\sigma^2 = 0$;
  - Lognormal-Poisson regression: $r \to \infty$;
  - Poisson regression: $\sigma^2 = 0$ and $r \to \infty$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Lognormal and gamma mixed negative binomial regression

- Model (Zhou et al., ICML2012):

$$y_i \sim \text{NegBino}\left(r, p_i\right), \quad r \sim \text{Gamma}(a_0, 1/h)$$

$$\psi_i = \log\left(\frac{p_i}{1 - p_i}\right) = \boldsymbol{x}_i^T \boldsymbol{\beta} + \ln \epsilon_i, \quad \epsilon_i \sim \ln \mathcal{N}(0, \sigma^2)$$

- Bayesian inference with the Polya-Gamma distribution.
- Model properties:

$$\text{Var}[y_i|\boldsymbol{x}_i] = \mathbb{E}[y_i|\boldsymbol{x}_i] + \left(e^{\sigma^2}(1 + r^{-1}) - 1\right)\mathbb{E}^2[y_i|\boldsymbol{x}_i].$$

- Special cases:
  - Negative binomial regression: $\sigma^2 = 0$;
  - Lognormal-Poisson regression: $r \to \infty$;
  - Poisson regression: $\sigma^2 = 0$ and $r \to \infty$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Lognormal and gamma mixed negative binomial regression

- Model (Zhou et al., ICML2012):

$$y_i \sim \text{NegBino}\,(r, p_i), \quad r \sim \text{Gamma}(a_0, 1/h)$$

$$\psi_i = \log\left(\frac{p_i}{1 - p_i}\right) = \boldsymbol{x}_i^T \boldsymbol{\beta} + \ln \epsilon_i, \quad \epsilon_i \sim \ln \mathcal{N}(0, \sigma^2)$$

- Bayesian inference with the Polya-Gamma distribution.
- Model properties:

$$\text{Var}[y_i | \boldsymbol{x}_i] = \mathbb{E}[y_i | \boldsymbol{x}_i] + \left( e^{\sigma^2}(1 + r^{-1}) - 1 \right) \mathbb{E}^2[y_i | \boldsymbol{x}_i].$$

- Special cases:
  - Negative binomial regression: $\sigma^2 = 0$;
  - Lognormal-Poisson regression: $r \to \infty$;
  - Poisson regression: $\sigma^2 = 0$ and $r \to \infty$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data
Motivations
Count
distributions
Negative
binomial
distribution
Relationships
between
distributions
Count regression

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

# Count regression example

- Count regression on the NASCAR dataset:

| Model | Poisson | NB | LGNB | LGNB |
| Parameters | (MLE) | (MLE) | (VB) | (Gibbs) |
|---|---|---|---|---|
| $\sigma^2$ | N/A | N/A | 0.1396 | 0.0289 |
| $r$ | N/A | 5.2484 | 18.5825 | 6.0420 |
| $\beta_0$ | -0.4903 | -0.5038 | -3.5271 | -2.1680 |
| $\beta_1$ (Laps) | 0.0021 | 0.0017 | 0.0015 | 0.0013 |
| $\beta_2$ (Drivers) | 0.0516 | 0.0597 | 0.0674 | 0.0643 |
| $\beta_3$ (TrkLen) | 0.6104 | 0.5153 | 0.4192 | 0.4200 |

- Using Variational Bayes inference, we can calculate the correlation matrix for $(\beta_1, \beta_2, \beta_3)^T$ as

$$\begin{pmatrix} 1.0000 & -0.4824 & 0.8933 \\ -0.4824 & 1.0000 & -0.7171 \\ 0.8933 & -0.7171 & 1.0000 \end{pmatrix}$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Latent Dirichlet allocation (Blei et al., 2003)

- Hierarchical model:

$$x_{ji} \sim \text{Mult}(\phi_{z_{ji}})$$
$$z_{ji} \sim \text{Mult}(\boldsymbol{\theta}_j)$$
$$\phi_k \sim \text{Dir}(\eta, \dots, \eta)$$
$$\boldsymbol{\theta}_j \sim \text{Dir}\left(\frac{\alpha}{K}, \dots, \frac{\alpha}{K}\right)$$

- There are $K$ topics $\{\phi_k\}_{1,K}$, each of which is a distribution over the $V$ words in the vocabulary.

- There are $N$ documents in the corpus and $\boldsymbol{\theta}_j$ represents the proportion of the $K$ topics in the $j$th document.

- $x_{ji}$ is the $i$th word in the $j$th document.

- $z_{ji}$ is the index of the topic selected by $x_{ji}$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

- Denote $n_{vjk} = \sum_i \delta(x_{ji} = v)\delta(z_{ji} = k)$, $n_{v \cdot k} = \sum_j n_{vjk}$, $n_{jk} = \sum_v n_{vjk}$, and $n_{\cdot k} = \sum_j n_{jk}$.

- Blocked Gibbs sampling:

$$P(z_{ji} = k | -) \propto \phi_{x_{ji}k}\theta_{jk}, \quad k \in \{1, \ldots, K\}$$

$$(\phi_k | -) \sim \text{Dir}(\eta + n_{1 \cdot k}, \ldots, \eta + n_{V \cdot k})$$

$$(\boldsymbol{\theta}_j | -) \sim \text{Dir}\left(\frac{\alpha}{K} + n_{j1}, \ldots, \frac{\alpha}{K} + n_{jK}\right)$$

- Variational Bayes inference (Blei et al., 2003).

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

- Collapsed Gibbs sampling (Griffiths and Steyvers, 2004):
  - Marginalizing out both the topics $\{\phi_k\}_{1,K}$ and the topic proportions $\{\theta_j\}_{1,N}$.
  - Sample $z_{ji}$ conditioning on all the other topic assignment indices $\boldsymbol{z}^{-ji}$:

$$P(z_{ji} = k|\boldsymbol{z}^{-ji}) \propto \frac{\eta + n_{x_{ji}\cdot k}^{-ji}}{V\eta + n_{\cdot k}^{-ji}} \left(n_{jk}^{-ji} + \frac{\alpha}{K}\right), \quad k \in \{1, \dots, K\}$$

- This is easy to understand as

$$P(z_{ji} = k|\phi_k, \boldsymbol{\theta}_j) \propto \phi_{x_{ji}k}\theta_{jk}$$

$$P(z_{ji} = k|\boldsymbol{z}^{-ji}) = \iint P(z_{ji} = k|\phi_k, \boldsymbol{\theta}_j)P(\phi_k, \boldsymbol{\theta}_j|\boldsymbol{z}^{-ji})d\phi_k d\boldsymbol{\theta}_j$$

$$P(\phi_k|\boldsymbol{z}^{-ji}) = \text{Dir}(\eta + n_{1\cdot k}^{-ji}, \dots, \eta + n_{V\cdot k}^{-ji})$$

$$P(\boldsymbol{\theta}_j|\boldsymbol{z}^{-ji}) = \text{Dir}\left(\frac{\alpha}{K} + n_{j1}^{-ji}, \dots, \frac{\alpha}{K} + n_{jK}^{-ji}\right)$$

$$P(\phi_k, \boldsymbol{\theta}_j|\boldsymbol{z}^{-ji}) = P(\phi_k|\boldsymbol{z}^{-ji})P(\boldsymbol{\theta}_j|\boldsymbol{z}^{-ji})$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

- In latent Dirichlet allocation, the words in a document are assumed to be exchangeable (bag-of-words assumption).
- Below we will relate latent Dirichlet allocation to Poisson factor analysis and show it essentially tries to factorize the term-document word count matrix under the Poisson likelihood:

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Poisson factor alaysis

- Factorize the term-document word count matrix $\mathbf{M} \in \mathbb{Z}_+^{V \times N}$ under the Poisson likelihood as

$$\mathbf{M} \sim \text{Pois}(\mathbf{\Phi\Theta})$$

  where $\mathbb{Z}_+ = \{0, 1, \dots\}$ and $\mathbb{R}_+ = \{x : x > 0\}$.

- $m_{vj}$ is the number of times that term $v$ appears in document $j$.

- Factor loading matrix: $\mathbf{\Phi} = (\phi_1, \dots, \phi_K) \in \mathbb{R}_+^{V \times K}$.

- Factor score matrix: $\mathbf{\Theta} = (\theta_1, \dots, \theta_N) \in \mathbb{R}_+^{K \times N}$.

- A large number of discrete latent variable models can be united under the Poisson factor analysis framework, with the main differences on how the priors for $\phi_k$ and $\theta_j$ are constructed.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

## Two equivalent augmentations

- Poisson factor analysis

$$m_{vj} \sim \text{Pois}\left(\sum_{k=1}^{K} \phi_{vk}\theta_{jk}\right)$$

- Augmentation 1:

$$m_{vj} = \sum_{k=1}^{K} n_{vjk}, \ n_{vjk} \sim \text{Pois}(\phi_{vk}\theta_{jk})$$

- Augmentation 2:

$$m_{vj} \sim \text{Pois}\left(\sum_{k=1}^{K} \phi_{vk}\theta_{jk}\right), \ \zeta_{vjk} = \frac{\phi_{vk}\theta_{jk}}{\sum_{k=1}^{K} \phi_{vk}\theta_{jk}}$$

$$[n_{vj1}, \cdots, n_{vjK}] \sim \text{Mult}\left(m_{vj}; \zeta_{vj1}, \cdots, \zeta_{vjK}\right)$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Nonnegative matrix factorization and gamma-Poisson factor analysis

- Gamma priors on $\mathbf{\Phi}$ and $\mathbf{\Theta}$:
$$m_{vj} = \text{Pois}\left(\sum_{k=1}^{K} \phi_{vk}\theta_{jk}\right)$$

$$\phi_{vk} \sim \text{Gamma}(a_\phi, 1/b_\phi), \quad \theta_{jk} \sim \text{Gamma}(a_\theta, g_k/a_\theta).$$

- Expectation-Maximization (EM) algorithm:
$$\phi_{vk} = \phi_{vk} \frac{\frac{a_\phi-1}{\phi_{vk}} + \sum_{i=1}^{N} \frac{m_{vj}\theta_{jk}}{\sum_{k=1}^{K} \phi_{vk}\theta_{jk}}}{b_\phi + \theta_{.k}}.$$
$$\theta_{jk} = \theta_{jk} \frac{\frac{a_\theta-1}{\theta_{jk}} + \sum_{p=1}^{P} \frac{m_{vj}\phi_{vk}}{\sum_{k=1}^{K} \phi_{vk}\theta_{jk}}}{a_\theta/g_k + \phi_{.k}}.$$

- If we set $b_\phi = 0$, $a_\phi = a_\theta = 1$ and $g_k = \infty$, then the EM algorithm is the same as those of non-negative matrix factorization (Lee and Seung, 2000) with an objective function of minimizing the KL divergence $D_{KL}(\mathbf{M}||\mathbf{\Phi\Theta})$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Latent Dirichlet
allocation

Poisson factor
analysis

Relational
network
analysis

Main
references

# Latent Dirichlet allocation and Dirichlet-Poisson factor analysis

- Dirichlet priors on $\boldsymbol{\Phi}$ and $\boldsymbol{\Theta}$:
$$m_{vj} = \text{Pois}\left(\sum_{k=1}^{K} \phi_{vk}\theta_{jk}\right)$$

$$\boldsymbol{\phi}_k \sim \text{Dir}(\eta, \ldots, \eta), \quad \boldsymbol{\theta}_j \sim \text{Dir}(\alpha/K, \ldots, \alpha/K).$$

- One may show that both the block Gibbs sampling inference and variational Bayes inference of the Dirichlet-Poisson factor analysis model are the same as that of the Latent Dirichlet allocation.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Beta-gamma-Poisson factor analysis

- Hierachical model (Zhou et al., 2012, Zhou and Carin, 2014):

$$m_{vj} = \sum_{k=1}^{K} n_{vjk}, \ n_{vjk} \sim \text{Pois}(\phi_{vk}\theta_{jk})$$

$$\phi_k \sim \text{Dir}(\eta, \cdots, \eta),$$

$$\theta_{jk} \sim \text{Gamma}[r_j, p_k/(1 - p_k)],$$

$$r_j \sim \text{Gamma}(e_0, 1/f_0),$$

$$p_k \sim \text{Beta}[c/K, c(1 - 1/K)].$$

- $n_{jk} = \sum_{v=1}^{V} n_{vjk} \sim \text{NB}(r_j, p_k)$

- This parametric model becomes a nonparametric Bayesian model governed by the beta-negative binomial process as $K \to \infty$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Gamma-gamma-Poisson factor analysis

- Hierachical model (Zhou and Carin, 2014):

$$m_{vj} = \sum_{k=1}^{K} n_{vjk}, \ n_{vjk} \sim \text{Pois}(\phi_{vk}\theta_{jk})$$
$$\phi_k \sim \text{Dir}(\eta, \cdots, \eta),$$
$$\theta_{jk} \sim \text{Gamma}[r_k, p_j/(1-p_j)],$$
$$p_j \sim \text{Beta}(a_0, b_0),$$
$$r_k \sim \text{Gamma}(\gamma_0/K, 1/c).$$

- $n_{jk} \sim \text{NB}(r_k, p_j)$

- This parametric model becomes a nonparametric Bayesian model governed by the gamma-negative binomial process as $K \to \infty$.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Poisson factor analysis and mixed-membership modeling

- We may represent the Poisson factor analysis

$$m_{vj} = \sum_{k=1}^{K} n_{vjk}, \ n_{vjk} \sim \text{Pois}(\phi_{vk}\theta_{jk})$$

in terms of a mixed-membership model, whose group sizes are randomized, as

$$x_{ji} \sim \text{Mult}(\phi_{z_{ji}}), \ z_{ji} \sim \sum_{k=1}^{K} \frac{\theta_{jk}}{\sum_k \theta_{jk}} \delta_k, \ m_j \sim \text{Pois}\left(\sum_k \theta_{jk}\right),$$

where $i = 1, \ldots, m_j$ in the $j$th document, and
$n_{vjk} = \sum_{i=1}^{m_j} \delta(x_{ji} = v)\delta(z_{ji} = k)$.

- The likelihoods of the two representations are different update to a multinomial coefficient (Zhou, 2014).

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Connections to previous approaches

- Nonnegative matrix factorization (K-L divergence) (NMF)
- Latent Dirichlet allocation (LDA)
- GaP: gamma-Poisson factor model (GaP) (Canny, 2004)
- Hierarchical Dirichlet process LDA (HDP-LDA) (Teh et al., 2006)

| Poisson factor analysis priors on $\theta_{jk}$ | Infer $(p_k, r_j)$ | Infer $(p_j, r_k)$ | Support $K \to \infty$ | Related algorithms |
|---|---|---|---|---|
| gamma | $\times$ | $\times$ | $\times$ | NMF |
| Dirichlet | $\times$ | $\times$ | $\times$ | LDA |
| beta-gamma | $\checkmark$ | $\times$ | $\checkmark$ | GaP |
| gamma-gamma | $\times$ | $\checkmark$ | $\checkmark$ | HDP-LDA |

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Blocked Gibbs sampling

- Sample $z_{ji}$ from multinomial;
  $n_{vjk} = \sum_{i=1}^{m_j} \delta(x_{ji} = v)\delta(z_{ji} = k)$.
- Sample $\phi_k$ from Dirichlet
- For the beta-negative binomial model
  (beta-gamma-Poisson factor analysis)
    - Sample $l_{jk}$ from CRT$(n_{jk}, r_j)$
    - Sample $r_j$ from gamma
    - Sample $p_k$ from beta
    - Sample $\theta_{jk}$ from Gamma$(r_j + n_{jk}, p_k)$
- For the gamma-negative binomial model
  (gamma-gamma-Poisson factor analysis)
    - Sample $l_{jk}$ from CRT$(n_{jk}, r_k)$
    - Sample $r_k$ from gamma
    - Sample $p_j$ from beta
    - Sample $\theta_{jk}$ from Gamma$(r_k + n_{jk}, p_j)$
- Collapsed Gibbs sampling for the beta-negative binomial
  model can be found in (Zhou, 2014).

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling
Latent Dirichlet
allocation
Poisson factor
analysis

Relational
network
analysis

Main
references

# Example application

- Example Topics of United Nation General Assembly Resolutions inferred by the gamma-gamma-Poisson factor analysis:

| Topic 1 | Topic 2 | Topic 3 | Topic 4 | Topic 5 |
|---------|---------|---------|---------|---------|
| trade | rights | environment | women | economic |
| world | human | management | gender | summits |
| conference | united | protection | equality | outcomes |
| organization | nations | affairs | including | conferences |
| negotiations | commission | appropriate | system | major |

- The gamma-negative binomial and beta-negative binomial models have distinct mechanisms on controlling the number of inferred factors.

- They produce state-of-the-art perplexity results when used for topic modeling of a document corpus (Zhou et al, 2012, Zhou and Carin 2014, Zhou 2014).

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Relational
network
analysis
Stochastic
blockmodel

Main
references

# Relational network

- A relational network (graph) is commonly used to describe the relationship between nodes, where a node could represent a person, a movie, a protein, etc.

- Two nodes are connected if there is an edge (link) between them.

- An undirected unweighted relational network with $N$ nodes can be equivalently represented with a sysmetric binary affinity matrix $B \in \{0,1\}^{N \times N}$, where $b_{ij} = b_{ji} = 1$ if an edge exists between nodes $i$ and $j$ and $b_{ij} = b_{ji} = 0$ otherwise.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Stochastic
blockmodel

Main
references

# Stochastic blockmodel

- Each node is assigned to a cluster.
- The probability for an edge to exist between two nodes is solely decided by the clusters that they are assigned to.
- Hierachical model:

$$
\begin{aligned}
b_{ij} &\sim \text{Bernoulli}(p_{z_i z_j}), \quad \text{for } j > i \\
p_{k_1 k_2} &\sim \text{Beta}(a_0, b_0), \\
z_i &\sim \text{Mult}(\pi_1, \ldots, \pi_K), \\
(\pi_1, \ldots, \pi_K) &\sim \text{Dir}(\alpha/K, \ldots, \alpha/K)
\end{aligned}
$$

- Blocked Gibbs sampling:

$$
P(z_i = k | -) = \pi_k \left\{ \prod_{j \neq i} p_{k z_j}^{b_{ij}} (1 - p_{k z_j})^{1 - b_{ij}} \right\}
$$

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Stochastic
blockmodel

Main
references

# Infinite relational model (Kemp et al., 2006)

- As $K \to \infty$, the stochastic block model becomes a nonparametric Bayesian model governed by the Chinese restaurant process (CRP) with concentration parameter $\alpha$:
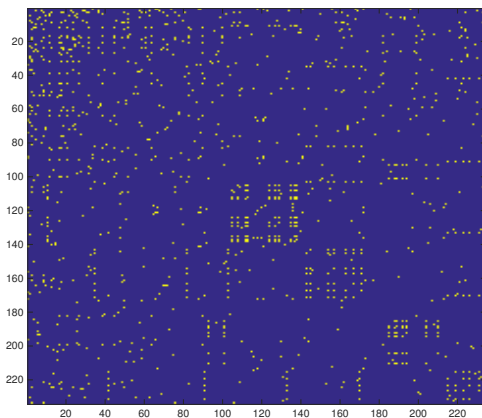
$$b_{ij} \sim \text{Bernoulli}(p_{z_i z_j}), \quad \text{for } i > j$$
$$p_{k_1 k_2} \sim \text{Beta}(a_0, b_0),$$
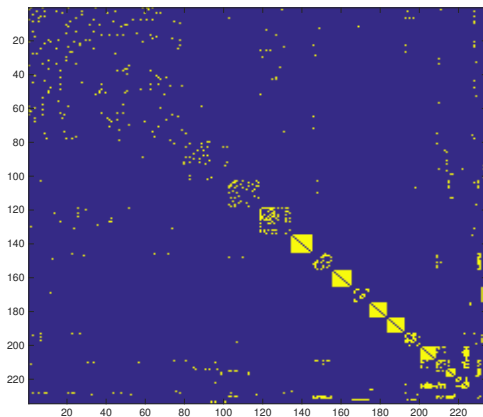$$(z_1, \ldots, z_N) \sim \text{CRP}(\alpha)$$

- Collapsed Gibbs sampling can be derived by marginalizing out $p_{k_1 k_2}$ and using the prediction rule of the Chinese restaurant process.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

The coauthor network of the top 234 NIPS authors.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Stochastic
blockmodel

Main
references

The reordered network using the stochastic blockmodel.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

The estimated link probabilities within and between blocks.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

Outline

Analysis of
count data

Count matrix
factorization
and topic
modeling

Relational
network
analysis

Main
references

D. Blei, A. Ng, and M. Jordan.
Latent Dirichlet allocation.
*J. Mach. Learn. Res.*, 2003.

T. L. Griffiths and M. Steyvers.
Finding scientific topics.
*PNAS*, 2004.

C. Kemp, J. B. Tenenbaum, T. L. Griffiths, T. Yamada, and N. Ueda.
Learning systems of concepts with an infinite relational model.
In *AAAI*, 2006.

D. D. Lee and H. S. Seung.
Algorithms for non-negative matrix factorization.
In *NIPS*, 2000.

Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei.
Hierarchical Dirichlet processes.
*JASA*, 2006.

M. Zhou, L. Hannah, D. Dunson, and L. Carin.
Beta-negative binomial process and Poisson factor analysis.
In *AISTATS*, 2012.

M. Zhou, L. Li, D. Dunson, and L. Carin.
Lognormal and gamma mixed negative binomial regression.
In *ICML*, 2012.

Parametric
Bayesian
Models: Part
II

Mingyuan
Zhou and
Lizhen Lin

M. Zhou and L. Carin.

Augment-and-conquer negative binomial processes.
In *NIPS*, 2012.

M. Zhou and L. Carin.

Negative binomial process count and mixture modeling.
*IEEE TPAMI*, 2014.

M. Zhou.

Beta-negative binomial process and exchangeable random partitions for mixed-membership modeling.
In *NIPS*, 2014.