# **DSC250: Advanced Data Mining**

**Topic Models** 

**Zhiting Hu** Lecture 6, Jan 24, 2025



HALICIOĞLU DATA SCIENCE INSTITUTE

**In-class paper presentation** 

# Adversarial Examples are not Bugs, they are Features



Andrew Ilyas\*



Logan Engstrom\*



Dimitris Tsipras\*



Brandon Tran



Shibani Santurkar\*



Aleksander Mądry

# **Topic Models**

## Outline

- Representations of Text and Topics
- Topic Model v1: Multinomial Mixture Model
- Topic Model v2: Probabilistic Latent Semantic Analysis (pLSA)
- Topic Model v3: Latent Dirichlet Allocation (LDA)

Slides adapted from:

- Y. Sun, CS 247: Advanced Data Mining
- M. Gormley, 10-701 Introduction to Machine Learning

## Motivation

Suppose'you're'given'a'massive'corpora'and'asked'to'carry'but'the" following'tasks"

- **Organize**'the'documents'Into'**thematic&ategories&**
- **Describe**'the'evolution'bf'those'tategories'**bver&ime**&
- Enable'a'domain'expert'to'analyze&nd&nderstand&he'content"
- Find'**relationships**'between'the'tategories"
- Understandenow'authorship'Influnces"the" ont ent "



## **Motivation**

Suppose you're given a massive corpora and asked to carry out the following tasks

- Organize the documents into thematic categories
- **Describe** the evolution of those categories over time
- Enable a domain expert to analyze and understand the content
- Find **relationships** between the categories
- Understand how **authorship** influences the content

#### **Topic Modeling:**

A method of (usually unsupervised) discovery of latent or hidden structure in a corpus

- Applied primarily to text corpora, but **techniques are more general**
- Provides a **modeling toolbox**
- Has prompted the exploration of a variety of new inference methods to accommodate large-scale datasets

#### **Topic Modeling: Examples**



Figure from (Blei, 2011), shows topics and top words learned automatically from reading 17,000 Science articles

#### **Topic Modeling: Examples**

Dirichlet;multinomial&egression&DMR)&opic&nodel&n&CML& (Mimno'&'McCallum,'2008)''

Topic 0 [0.152]



#### Topic 54 [0.051]





problem, optimization, problems, convex, convex optimization, linear, semidefinite programming, formulation, sets, constraints, proposed, margin, maximum margin, optimization problem, linear programming, programming, procedure, method, cutting plane, solutions

decision trees, trees, tree, decision tree, decision, tree ensemble, junction tree, decision tree learners, leaf nodes, arithmetic circuits, ensembles modts, skewing, ensembles, anytime induction decision trees, trees trees, random forests, objective decision trees, tree learners, trees grove, candidate split

inference, approximate inference, exact inference, markov chain, models, approximate, gibbs sampling, variational, bayesian, variational inference, variational bayesian, approximation, sampling, methods, exact, bayesian inference, dynamic bayesian, process, mcmc, efficient

http://'www.cs.umass.edu/~mimno/icml100.html"

8



### **Other Applications of Topic Models**

Spacial'LDA"

(Wang'&'Grimson,'2007)"



### **Other Applications of Topic Models**

• Word Sense Induction

(Brody & Lapata, 2009)

Senses of *drug* (WSJ) 1. U.S., administration, federal, against, war, dealer 2. patient, people, problem, doctor, company, abuse 3. company, million, sale, maker, stock, inc. 4. administration, food, company, approval, FDA

Senses of *drug* (BNC) 1. patient, treatment, effect, anti-inflammatory 2. alcohol, treatment, patient, therapy, addiction 3. patient, new, find, effect, choice, study 4. test, alcohol, patient, abuse, people, crime 5. trafficking, trafficker, charge, use, problem 6. abuse, against, problem, treatment, alcohol 7. people, wonder, find, prescription, drink, addict 8. company, dealer, police, enforcement, patient

11

• Selectional Preference<sup>7. people, wonder, find, prescription, drink, addie 8. company, dealer, police, enforcement, patient</sup>

#### (Ritter et al., 2010)

Topic t	Arg1	Relations which assign highest probability to $t$	Arg2
18	The residue - The mixture - The reaction	was treated with, is	EtOAc - CH2Cl2 - H2O - CH.sub.2Cl.sub.2
	mixture - The solution - the mixture - the re-	treated with, was	- H.sub.2O - water - MeOH - NaHCO3 -
	action mixture - the residue - The reaction -	poured into, was	Et2O - NHC1 - CHCl.sub.3 - NHCl - drop-
	the solution - The filtrate - the reaction - The	extracted with, was	wise - CH2Cl.sub.2 - Celite - Et.sub.2O -
	product - The crude product - The pellet -	purified by, was di-	Cl.sub.2 - NaOH - AcOEt - CH2Cl2 - the
	The organic layer - Thereto - This solution	luted with, was filtered	mixture - saturated NaHCO3 - SiO2 - H2O
	- The resulting solution - Next - The organic	through, is disolved in,	- N hydrochloric acid - NHCl - preparative
	phase - The resulting mixture - C. )	is washed with	HPLC - to0 C

### **Text Data**

- Word/term
- Document
  - A sequence of words
- Corpus
  - A collection of

documents



#### **Represent a Document**

- Most common way: Bag-of-Words
  - Ignore the order of words
  - keep the count

acl Human machine nterface for Lab ABC computer applications

- c2: A survey of aser opinion of computer system response time
- c3: The EPS user interface management system
- c4: System and human system engineering testing of EPS
- c5: Relation of user-perceived response time to error measurement

m1: The generation of random, binary, unordered trees

- m2: The intersection graph of paths in trees
- m3: Graph minors IV: Widths of trees and well-quasi-ordering

m4: Graph minors: A survey



c1 c2 c3 c4 c5 m1 m2 m3 m4 0 0 0 0 human 0 0 0 0 0 0 interface 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 computer 0 0 0 user 0 0 0 0 0 system 0 3 0 0 response 0 0 0 0 0 0 time 0 EPS 0 0 0 0 0 0 0 survey 0 0 0 0 trees 0 0 0 0 0 graph 0 0 0 0 0 minors

Vector space model

#### **Represent a Document**

- Represent the doc as a vector where each entry corresponds to a different word and the number at that entry corresponds to how many times that word was present in the document (or some function of it)
  - Number of words is huge
  - Select and use a smaller set of words that are of interest
  - E.g. uninteresting words: 'and', 'the' 'at', 'is', etc. These are called <u>stop-</u> words
  - <u>Stemming</u>: remove endings. E.g. (learn', learning', 'learnable', 'learned', could be substituted by the single stem 'learn'
  - Other simplifications can also be invented and used
  - The set of different remaining words is called <u>dictionary</u> or <u>vocabulary</u>. Fix an ordering of the terms in the dictionary so that you can operate them by their index.
  - Can be extended to bi-gram, tri-gram, or so



### **Limitations of Bag-of-Words**

- Dimensionality
  - High dimensionality

V-dim

- Sparseness
  - Most of the entries are zero
- Shallow representation
  - The vector representation does not capture semantic relations between words
    - Ex: "Tom loves Kate."

#### **Represent a Topic**

# • A topic is represented by a word distribution

### • Relate to an issue

$\boldsymbol{\mathcal{C}}$	universe	0.0439	drug	0.0672		cells	0.0675		sequence	0.0818	years	0.156
S	galaxies 🥜	0.0375	patients	0.0493		stem	0.0478		sequences	0.0493	million	0.0556
	clusters	0.0279	drugs	0.0444		human	0.0421		genome	0.033	ago	0.045
(	matter	0.0233	clinical	0.0346		cell	0.0309		dna	0.0257	time	0.0317
$\mathbf{N}$	galaxy	0.0232	treatment	0.028		gene	0.025		sequencing	0.0172	age	0.0243
	cluster	0.0214	trials	0.0277		tissue	0.0185		map	0.0123	year	0.024
	cosmic	0.0137	therapy	0.0213		cloning	0.0169		genes	0.0122	record	0.0238
	dark	0.0131	trial	0.0164		transfer	0.0155		chromosome	0.0119	early	0.0233
	light	0.0109	disease	0.0157		blood	0.0113		regions	0.0119	billion	0.0177
1	density	0.01	medical	0.00997		embryos	0.0111		human	0.0111	history	0.0148
	bacteria	0.0983	male	0.0558	1	theory	0.0811	1[	immune	0.0909	stars	0.0524
	bacterial	0.0561	females	0.0541		physics	0.0782	Ш	response	0.0375	star	0.0458
	resistance	0.0431	female	0.0529		physicists	0.0146	Ш	system	0.0358	astrophys	0.0237
	coli	0.0381	males	0.0477		einstein	0.0142	Ш	responses	0.0322	mass	0.021
	strains	0.025	sex	0.0339		university	0.013	Ш	antigen	0.0263	disk.	0.0173
	microbiol	0.0214	reproductive	0.0172		gravity	0.013	Ш	antigens	0.0184	black	0.0161
	microbial	0.0196	offspring	0.0168		black	0.0127	Ш	immunity	0.0176	gas	0.0149
	strain	0.0165	sexual	0.0166		theories	0.01		immunology	0.0145	stellar	0.0127
	salmonella	0.0163	reproduction	0.0143		aps	0.00987		antibody	0.014	astron	0.0125
	resistant	0.0145	eggs	0.0138		matter	0.00954		autoimmune	0.0128	hole	0.00824
			-		-							

TOPIC 42 LUCIC 41 scheed exptain main compared to the scheme data werkening compared to the scheme data provide scheme data compared to the scheme conversion balance opport the tit is cutants nodgall the Sch ta the the scheme data ratio scheme to ta the scheme data to the scheme ratio scheme to ta the scheme data to the scheme ratio scheme to ta the scheme data to the scheme ratio scheme to ta the scheme to the scheme ratio scheme to ta the scheme to the scheme ratio scheme to the scheme to the scheme ratio scheme to the scheme to the scheme to the scheme ratio scheme to the scheme to the scheme to the scheme scheme to the scheme to the scheme to the scheme to the scheme scheme to the sc continue ing centers ke ovement transportation and plication strengths ogionally

TOPIC 45

TOPIC 45 Consulty identified installed fibrary anneally ingle grant (Second Second S

opened updated numerous beginning

single

ovalues mechanisms fund ec extense modernizationgacilly. Proposed benefits amount anyte: increase fee miduel scal benefit statup private Project fund sco statup private Project investments nead et Capita Cost investments anyte increase

program: term costs mainten an ce

or year contena enformance projects program operating conditions ity funding tax fin ancial is that affacilities service years model

managed improvements

TOPIC 48

nductedies si bili tylisted addre a

#### TOPIC 43

TOPIC 43 methological lander of americal physical and the second land and the second physical and the second second second second physical properties of an accertification method second second second second second second physical properties of a second second second physical physical second second second second second physical second second second second second second physical second second second second second second second physical second second second second second second second physical second second second second second second second second physical second sec The second secon

13

CO com: dos

pathasio<sub>sa</sub> Pateisod<sup>4</sup> cust Pateiso pt rolatio sintro:

u ctiler

sumo

appr

#### TOPIC 46

supplement locate soloway <sup>12</sup> Splanse, heat, Norsa Bog Heathploader <sup>1</sup> Splanskie, heat OWN prises <sup>1</sup> Splanskie, shortOWN prises <sup>1</sup> Splanskie, shortOWN prises <sup>1</sup> Splanskie, shortOWN prises <sup>1</sup> Splanskie, <sup>1</sup> Splanskie, <sup>1</sup> Splanskie <sup>1</sup> Splanskie, <sup>1</sup> Splanskie heigh/located 🔒 adopted policies reas residents irpor an figure Special by highway ight provide pha drain age Arr Newoll& El additio\_\_\_

#### TOPIC 49

govenn requirer indivi prioriti es bringro ensuring en recen ad  $16^{k}$ 

### **Topic Models**

### Topic modeling

- Get topics automatically from a corpus
- Assign documents to topics automatically
- Most frequently used topic models
  - pLSA
  - LDA

"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

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.

### **Notations**

- Word, document, topic
  - $\circ w, d, z$
- Word count in document:
  - c(w,d): number of times word w occurs in document d
  - $\circ$  or  $x_{dn}$ : number of times the *n*th word in the vocabulary occurs in document d
- Word distribution for each topic (  $eta_z$  )
  - $\circ \ \beta_{zw} : p(w|z)$



### **Recap: Multinomial distribution**

- Multinomial distribution
  - Discrete random variable x that takes one of M values  $\{1, \dots, M\}$

$$\circ p(\mathbf{x}=i) = \pi_{i}, \qquad \sum_{i} \pi_{i} = 1$$

- Out of n independent trials, let  $k_i$  be the number of times x = i was observed
- The probability of observing a vector of occurrences  $\mathbf{k} = [k_1, ..., k_M]$  is given by the *multinomial distribution* parametrized by  $\boldsymbol{\pi}$

$$p(\mathbf{k}|\boldsymbol{\pi}, \mathbf{n}) = p(k_1, \dots, k_m | \pi_1, \dots, \pi_m, \mathbf{n}) = \frac{\mathbf{n}!}{k_1! k_2! \dots k_m!} \prod_{i=1}^{n} \pi_i^{k_i}$$

- E.g., describing a text document by the frequency of occurrence of every distinct word
- For n = 1, a.k.a. categorical distribution
  - $p(x = i \mid \pi) = \pi_i$  In k = [k<sub>1</sub>,...,k<sub>M</sub>]: k<sub>i</sub> = 1, and k<sub>j</sub> = 0 for all j ≠ i → a.k.a., one-hot representation of i

- For documents with bag-of-words representation
  - $x_d = (x_{d1}, x_{d2}, ..., x_{dN}), x_{dn}$  is the number of words for nth word in the vocabulary
- Generative model

For documents with bag-of-words representation

•  $x_d = (x_{d1}, x_{d2}, ..., x_{dN}), x_{dn}$  is the number of words for nth word in the vocabulary

Generative model

Formulating the statistical relationship between words, documents and latent topics as a generative process describing how documents are created

- For documents with bag-of-words representation
  - $x_d = (x_{d1}, x_{d2}, ..., x_{dN}), x_{dn}$  is the number of words for nth word in the vocabulary
- Generative model
  - For each document
    - Sample its cluster label *z*~*Categorical*(π)
      π = (π<sub>1</sub>, π<sub>2</sub>, ..., π<sub>K</sub>), π<sub>k</sub> is the proportion of jth cluster
      p(z = k) = π<sub>k</sub>
    - Sample its word vector  $x_d \sim multinomial(\beta_z)$  $\beta_z = (\beta_{z1}, \beta_{z2}, ..., \beta_{zN}), \beta_{zn}$  is the parameter associate with nth word in the vocabulary

• 
$$p(\mathbf{x}_d|z=k) = \frac{(\sum_n x_{dn})!}{\prod_n x_{dn}!} \prod_n \beta_{kn}^{x_{dn}} \propto \prod_n \beta_{kn}^{x_{dn}}$$

N: Vocabushy Size N= length of doc

Graphical Model

Generative model

- Plates indicate replicated variables.
- Shaded nodes are observed; unshaded nodes are hidden.

- For each document
  - Sample its cluster labe  $(z^2 Categorical(\pi))$ 
    - $\boldsymbol{\pi} = (\pi_1, \pi_2, ..., \pi_K), \pi_k$  is the proportion of jth cluster

• 
$$p(z=k) = \pi_k$$

- Sample its word vector  $x_d$  multinomial  $(\beta_z)$ 
  - $\beta_z = (\beta_{z1}, \beta_{z2}, ..., \beta_{zN}), \beta_{zn}$  is the parameter associate with nth word in the vocabulary

• 
$$p(\mathbf{x}_d|z=k) = \frac{(\sum_n x_{dn})!}{\prod_n x_{dn}!} \prod_n \beta_{kn}^{x_{dn}} \propto \prod_n \beta_{kn}^{x_{dn}}$$





- Plates indicate replicated variables.
- Shaded nodes are observed; unshaded nodes are hidden.

- Generative model
  - For each document



 $\pi = (\pi_1, \pi_2, \dots, \pi_K), \pi_k$  is the proportion  $p(z = k) = \pi_k$ 

• Sample its word vector  $x_d \sim multinomia (\beta_z)$ 

•  $\beta_z = (\beta_{z1}, \beta_{z2}, ..., \beta_{zN}), \beta_{zn}$  is the parameter associate with nth word in the vocabulary

$$p(\mathbf{x}_d|z=k) = \frac{(\sum_n x_{dn})!}{\prod_n x_{dn}!} \prod_n \beta_{kn}^{x_{dn}} \propto \prod_n \beta_{kn}^{x_{dn}}$$



#### **Likelihood Function**

$$L = \prod_{d} p(\mathbf{x}_{d}) = \prod_{d} \sum_{k} p(\mathbf{x}_{d}, z = k)$$
$$= \prod_{d} \sum_{k} p(\mathbf{x}_{d} | z = k) p(z = k)$$
$$= \prod_{d} \frac{(\sum_{n} x_{dn})!}{\prod_{n} x_{dn}!} \sum_{k} p(z = k) \prod_{n} \beta_{kn}^{x_{dn}}$$

Limitations of Multinomial Mixture Model

 All the words in the same documents are sampled from the same topic



"Sene

### **Limitations of Multinomial Mixture Model**

Mixture'vs.'Admixture'





Diagrams'from'Wallach,'JHU'2011,'Slides"

#### Topic Model v2: Probabilistic Latent Semantic Analysis (pLSA)

	" <u>Arts</u> "	"Budgets"	"Children"	"E <u>duca</u> tion"
	NEW	MILLION	CHILDREN	SCHOOL
	FILM	TAX	WOMEN	STUDENTS
	SHOW	PROGRAM	PEOPLE	SCHOOLS
	MUSIC	BUDGET	CHILD	EDUCATION
1	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

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.

#### **Generative Model for pLSA**

- For each position in d,  $n = 1, ..., N_d$ 
  - Generate the topic for the position as  $z_n | d \sim Categorical(\theta_d), i.e., p(z_n = k | d) = \theta_{dk}$ (Note, 1 trial multinomial)
    - Generate the word for the position as

$$w_n | z_n \sim Categorical(\beta_{z_n}), i.e., p(w_n = w | z_n) = \beta_{z_n w}$$

#### **Generative Model for pLSA**

- For each position in d,  $n = 1, ..., N_d$ 
  - Generate the topic for the position as  $z_n | d \sim Categorical(\theta_d), i.e., p(z_n = k | d) = \theta_{dk}$ (Note, 1 trial multinomial)
    - Generate the word for the position as

 $w_n | z_n \sim Categorical(\boldsymbol{\beta}_{z_n}), i.e., p(w_n = w | z_n) = \beta_{z_n w}$ 



### **Generative Model for pLSA**

• For each position in d,  $n = 1, ..., N_d$ 

- Generate the topic for the position as  $z_n | d \sim Categorical(\theta_d), i.e., p(z_n = k | d) = \theta_{dk}$  (Note, 1 trial multinomial)
  - Generate the word for the position as

 $w_n | z_n \sim Categorical(\boldsymbol{\beta}_{z_n}), i.e., p(w_n = w | z_n) = \beta_{z_n w}$ 





![](_page_32_Picture_0.jpeg)

### **Likelihood Function**

![](_page_33_Figure_1.jpeg)

# Probability of a word w

$$p(w|d,\theta,\beta) = \sum_{k} p(w,z=k|d,\theta,\beta)$$
$$= \sum_{k} p(w|z=k,d,\theta,\beta)p(z=k|d,\theta,\beta) = \sum_{k} \beta_{kw}\theta_{dk}$$

### **Likelihood Function**

![](_page_34_Figure_1.jpeg)

# Probability of a word w

$$p(w|d,\theta,\beta) = \sum_{k} p(w,z=k|d,\theta,\beta)$$
$$= \sum_{k} p(w|z=k,d,\theta,\beta)p(z=k|d,\theta,\beta) = \sum_{k} \beta_{kw}\theta_{dk}$$

# Likelihood of a corpus

#### **Re-arrange the Likelihood Function**

Group the same word from different positions together

$$\max \log L = \sum_{dw} c(w, d) \log \sum_{z} \theta_{dz} \beta_{zw}$$
  
s.t.  $\sum_{z} \theta_{dz} = 1$  and  $\sum_{w} \beta_{zw} = 1$ 

### Limitations of pLSA

- Not a proper generative model
  - $\boldsymbol{\theta}_d$  is treated as a parameter
  - Cannot model new documents

# • Solution:

• Make it a proper generative model by adding priors to  $\theta$  and  $\beta$ 

## Limitations of pLSA

- Not a proper generative model
  - $\boldsymbol{\theta}_d$  is treated as a parameter
  - Cannot model new documents

- Solution:
  - Make it a proper generative model by adding priors to  $\theta$  and  $\beta$

Topic Model v3: Latent Dirichlet Allocation (LDA)

#### **Review: Dirichlet Distribution**

• Dirichlet distribution:  $\theta \sim Dirichlet(\alpha)$ 

• *i.e.*, 
$$p(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_{k} \alpha_{k})}{\prod_{k} \Gamma(\alpha_{k})} \prod_{k} \theta_{k}^{\alpha_{k}-1}$$
, where  $\alpha_{k} > 0$   
•  $\Gamma(\cdot)$  is gamma function:  $\Gamma(z) = \int_{0}^{\infty} e^{-t} t^{z-1} dt$   
•  $\Gamma(z+1) = z\Gamma(z)$ 

#### **Review: Dirichlet Distribution**

• Dirichlet distribution:  $\theta \sim Dirichlet(\alpha)$ 

• *i.e.*, 
$$p(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_{k} \alpha_{k})}{\prod_{k} \Gamma(\alpha_{k})} \prod_{k} \theta_{k}^{\alpha_{k}-1}$$
, where  $\alpha_{k} > 0$   
•  $\Gamma(\cdot)$  is gamma function:  $\Gamma(z) = \int_{0}^{\infty} e^{-t} t^{z-1} dt$   
•  $\Gamma(z+1) = z\Gamma(z)$ 

Simplex view:

• 
$$x = x_1(1,0,0) + x_2(0,1,0) + x_3(0,0,1)$$
  
• Where  $0 \le x_1, x_2, x_3 \le 1$  and  $x_1 + x_2 + x_3 = 1$ 

![](_page_39_Figure_5.jpeg)

#### **More Examples in the Simplex View**

![](_page_40_Figure_1.jpeg)

### **Topic Model v3: Latent Dirichlet Allocation (LDA)**

![](_page_41_Figure_1.jpeg)

 $\theta_d \sim Dirichlet(\alpha)$ : address topic distribution for unseen documents  $\beta_k \sim Dirichlet(\eta)$ : smoothing over words

### **Topic Model v3: Latent Dirichlet Allocation (LDA)**

![](_page_42_Figure_1.jpeg)

 $\theta_d \sim Dirichlet(\alpha)$ : address topic distribution for unseen documents  $\beta_k \sim Dirichlet(\eta)$ : smoothing over words

### **Generative Model for LDA**

![](_page_43_Figure_1.jpeg)

![](_page_43_Figure_2.jpeg)

![](_page_44_Picture_2.jpeg)

- The'generative&tory&egins'with'bnly'a'Dirichlet& prior&ver'the'topics."
- Each'topic&s'defind 'as'a'Multinomial&distribution'' over'the'vocabulary,'parameterized'by' β<sub>k</sub> '

![](_page_45_Figure_2.jpeg)

- The'generative&tory&egins'with'bnly'a'Dirichlet& prior&ver'the'topics."
- Each'topic&'defind 'as'a'Multinomial&istribution'' over'the'vocabulary,'parameterized'by' β<sub>k</sub> '

# LDA'for'Topic'Modeling"

![](_page_46_Figure_2.jpeg)

team, season, + hockey, player, + penguins, ice, ++ canadiens, + puck, +montreal, + stanley, +cup+

 A'topic'ls'visualized'as'lts'high&robability& words.'''

![](_page_47_Figure_2.jpeg)

- A'topic'ls'visualized'as'lts'high&robability& words.'''
- A'þedagogical'**label&**s'used'to'ldentify'the'topic."

![](_page_48_Figure_2.jpeg)

- A'topic'ls'visualized'as'lts'high'probability" words.'''
- A'þedagogical'**label**&s'used'to'ldentify'the'topic."

![](_page_49_Figure_2.jpeg)

![](_page_50_Figure_2.jpeg)

![](_page_51_Figure_2.jpeg)

![](_page_52_Figure_2.jpeg)

![](_page_53_Figure_2.jpeg)

![](_page_54_Figure_2.jpeg)

![](_page_55_Figure_2.jpeg)

![](_page_56_Figure_2.jpeg)

#### Joint Distribution for LDA

![](_page_57_Figure_1.jpeg)

 Joint distribution of latent variables and documents is:

$$p(\boldsymbol{\beta}_{1:K}, \boldsymbol{z}_{1:D}, \boldsymbol{\theta}_{1:D}, \boldsymbol{w}_{1:D} | \boldsymbol{\alpha}, \boldsymbol{\eta}) = \prod_{i=1}^{K} p(\boldsymbol{\beta}_{i} | \boldsymbol{\eta}) \prod_{d=1}^{D} p(\boldsymbol{\theta}_{d} | \boldsymbol{\alpha}) \left( \prod_{n=1}^{N} p(\boldsymbol{z}_{d,n} | \boldsymbol{\theta}_{d}) p(\boldsymbol{w}_{d,n} | \boldsymbol{\beta}_{1:K}, \boldsymbol{z}_{d,n}) \right)$$

# **Questions?**