

# Latent Dirichlet Allocation

---

Johan Alenlöv, Linköping University

2023-10-31

# Outline – Lecture 5

**Aim:** Introduce the Latent Dirichlet Allocation model (LDA) and show how to perform inference in this model through collapsed Gibbs sampling.

## Outline:

1. Topic modelling
2. The graphical model and plate notation.
3. Exploring a corpus using the LDA model.
4. Gibbs sampling in LDA.

# Topic modelling

SVENSKA DAGBLADET

Nyheter Näringsliv Kultur Ledare Debatt Tidningen

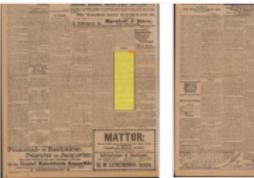
Sök



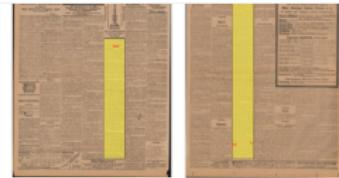
Meny



„generalmajor Gadd den 24 dennes kl. i e m Det litterära Nobelpriiset Med anledning af meddeleendet i en tysk...“



„...salunda skulle röjt dem Svensk-amerikansk aspirant till Nobelpriiset Hemlandets i Chicago redaktion...“



„...överläggningar om Nobelpriiset Skog smedels fonden och järnvägsbyggetet Km:t har medgivit att Orsa...“



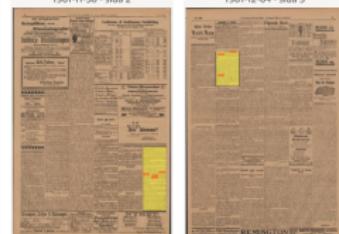
„...Samarows längesedan vissnade lagrar icke ha lämnat honom någon ro Eller är det månné Nobelpriiset...“



„...rykesundervisningen i Tyskland Österrike-Ungern Schweiz Belgien och Holland Det litterära Nobelpriiset för 1901...“



„...jordbruks O et litterara Nobelpriiset Biskop Billing samt professorerna Tegnér Söderwall och Rudin...“



„...E Åkerberg i Offerdilgrundane af vederbörande prövudelares beslut angående Nobelpriiset i fysik...“



„...af mikromästaren Georg Karsson Hur C D W förberedd Nobelpriiset ”jag föreställer mig att särskilt...“



„...hvard Posttidningen för i går har att säga Posttidningen yrtrar Det litterära Nobelpriiset Dtn...“



„...meddelade beslutet rörande Nobelpriiset i fysik hvilket skedde i följande tal något lågmäldt hålljet...“



„...Ryssland — Utsiktarna för dess genomförande främstellationen om Wreschenaffären Det litterära Nobelpriiset...“



„...Nobelpriiset på litteraturens område utan förutfattad mening mot tysk andebildning och tyskt...“



„...adress till Leo Tolstoj Postverkets stater 1903 Ny svensk Grönlandsexpedition Oet litterära Nobelpriiset...“

# Topic modelling

omskatt. Bland de  
hos Röntgenstråla-  
ren om s. den på  
ärkekraka avhandla-  
des. Närmede prak-  
tispar äro i olika  
ljuset, ärö de dä-  
men med den skil-  
luset fullt ogenom-  
skatt passeras av X-  
strå kroppar fullstän-  
digt ex. metaller  
ta, men tra, papp,  
papper och likn. och  
med ejvervaren  
en en framkän-  
ningsräckig kropp  
kommit i dessa vif-  
tämman därigenom  
troppeledes belyses  
samma på upptakten  
av spokdomarna äfven i sådana fall, då  
bakterierna redan lyckats vinna fläts och  
komma till utveckling inom organismen.  
Ett, och hittills det mest glimrande beviset  
på hvad i detta afseende kan vinnas erhövt  
der differen.

Så långt tillbaka i tiden, som kändes  
om mänskans sjukdomar kan ledas, harfa  
deteriorerat och förfallit varit ed  
af mänskighetskänsla. Tidigare  
då viselelsen satzat, så att den skenhant  
önskade, men stodas har den efter någon  
tids förlöpp åter blössat upp och ofta spridit  
sig som härliga epidemier af större  
eller mindre omfattning. Sedan fera Ar-  
titionen har den nu varit gängse inom den  
civiliserade världens annas i  
statens syfte.

Ira nägra ord på  
räntgen, gick därp-  
åde hr Röntgen vid  
i honom till kron-  
smade nu åt pris-  
hallen portfölj —  
spänningarna bröto  
sämen yträde nägra  
diken därpå under-  
vände tillbaka till  
Herr upp till talare-  
tal.

## vant Hoff.

Kemi har aka-  
ren vid Berlin Univer-  
sitetens vart  
bambrytande arbe-  
det trycket och öfver-  
sättningsen på atome-  
onriden har vant' s  
ens teori viktiga  
tan Dalton tid.  
Efter beträffar, har  
nu till en Pasteur  
hypotesen, som be-  
gägnat orientering  
elementarismerna,  
koförmeningarnas be-  
teende om kolatomen  
släggande af stereo-

gränande och omgestaltande inbytande på  
de medicinska vetenskapernas olika grenar.  
Huru mäktigt den inverkat på uppfattningen  
om den allmänna hälsovärden, och huru den  
tryckt sin prägel på nästan allt, som i detta  
härseende görs och läses, torde vara allbe-  
kannt. För den storstäende utvecklingen af  
kirurgien och därmed bestihedens gen-  
omgång, och ihjelycket stor del bakteriolo-  
gien att tacka.

Xren på de rent medicinska vetenskap-  
grenarna har bakteriologi redan givit  
morgna frukter af det största värde, och an-  
talet af dem, som äro stodda åt utveckling,  
läter sig icke beräkna.

Gemensam kunskaps-  
om bakterierna, som  
sjukdomsalster och genom inblicken  
om mänskligas hälsa och förfallit, har  
sina spökdomarna äfven i sådana fall, då  
bakterierna redan lyckats vinna fläts och  
komma till utveckling inom organismen.  
Ett, och hittills det mest glimrande beviset  
på hvad i detta afseende kan vinnas erhövt  
der differen.

Jag torde ej behövta mäla de skräckbilder,  
som den framtidens och den fortvänande hvil-  
ken fördömer givet, och den framtidens  
fortbrottet har den med emmien efter den an-  
dra. Numera är förhållanden betydligt  
förändrade, och taflan skulle nu kunna mäs-  
sa med anvisning af betydligt ljusare  
färger.

Visserligen innebär differen alltjämt en  
hotande fara — och så torde väl förhållanden  
alltid förblifva. Närpeligen kan man väl  
hoppas att någonstans kommande dag  
den hell och förtrollande sätter en ny  
sjukdomsalster, som berodas på bak-  
terier, men & ändras sedan. Bestredes detta  
af framstående fackmän. Något positivt  
känd förelig emellertid icke, och den vet-  
enskapliga utredningen af frågan saknas.  
Å minstre kunde man säga sig äga  
någon bestämd kändedom om arten av  
sjukdomsalstende parasiten.

Nyskända, emellertid Loeflers  
sin utvärdering och tycke  
elementarismerna, koförmeningarnas  
beteende om kolatomen  
släggande af stereole-

## Sully-Prudhomme pristagare.

Ea passas om historieforskare och filosoffer.

Bokförfattaren af Alfred Nobels hufrud-  
skaliga Huvverksamhet hörde den naturigt,  
att den författare som berättade för den  
dottern, hvilken växte så mycket och rätt-  
vist uppehende, i främsta rummet ville till-  
godose naturforskingen samt helona upp-  
täcker inom atskiljorna af dess grenar. Hans  
kosmopolitiska strävan fördrevne gjorde honom  
och till van af fredstanken och folkhövdinge-  
saken. Men han unnade dock i sina  
testamentariska bestämmelser ett rum åf  
historie. Väntades också om att han gjort  
ordningsförflytteri efter den sista förfarandet  
till hvilken han känna störjande drag-  
ning. Historien har honom tacksam där-  
för att även dess mälemlen varit föremål  
för hans omtanke, och om hon kommit sät-  
inom de svenska prägrossperna, fram-  
skimlar här måhända den riktiga berättelse  
grund. Sånnan och sedan författaren  
måste ha varit en mänsklig upp-  
blommar. I alla händelser erahis huvudsak-  
pristagande vid dessan, nyan tider littera-  
reux floraux i lin, som i ytre värde  
överflörs förra tidera gulvial.

Emellertid modifier utdelningen af det lit-  
teraturpriset **Nobelpriset** sinn alldeles särskilda  
svårigheter. "Litteratur" är ett vidsträckt  
begrepp. Nobelpriset har föreskrifts-  
skriften, att dardär skola vid prisfriaren be-  
dömmande innefatta en blott litteratur-  
verk, men & även annan hvilke genre  
form, och framfällassättning. Såga litterari-  
tär värde. Men härigenom vidgat fältet och svår-  
igheterna växa. Det redan är van-  
ligt att afgöra, huruvida, förtärtat att de  
förfatningar är ungefärligen i lika hög grad  
förfärtiga, priset skall lämnas åt en lyrik  
eller epik eller dramatisk skild, se inveck-  
lade, men & annan författare, som det är svårt  
att sätta varsin motsättning. Såna huvudsak-  
litteratorna blifva då, som man si-  
ger, inkommensurabla. Man kan emellertid  
trösta sig därmed att, eftersom prisutdelningar  
är ora årliga, mängden förfatningar, som den  
ärata fått vika för en annan och lika  
stor, vid ett följande skall tillväxla sig ve-  
skrönigt erkännande.

År 1883 betecknades en vändpunkt i diffe-  
ren för att förfatningens underlag, som  
är af en och annan antagts, att diffe-  
ren vore en sjukdom, som berodas på bak-  
terier, men & ändras sedan. Bestredes detta  
af framstående fackmän. Något positivt  
känd förelig emellertid icke, och den vet-  
enskapliga utredningen af frågan saknas.  
Å minstre kunde man säga sig äga  
någon bestämd kändedom om arten af  
sjukdomsalstende parasiten.

Nyskända, emellertid Loeflers  
sin utvärdering och tycke  
elementarismerna, koförmeningarnas  
beteende om kolatomen  
släggande af stereole-

nägra ord med pristagarene, var festen slut,  
och man troppade längsamt af.

## Prisdiplomens lydelse.

Diplomen är affärtade på svenska med  
bilagd översättning till resp. pristagares  
språk. De för fysik och kemil utdelade di-  
plomen är konstnärligt textade af fröken  
Sofie Giesberg, och bokbindarierbetet är ut-  
fört af den Becka firmen. Diplomen för  
medicin och litteratur ha textats af hr Axel  
Lindgren, och bokbindarierbetet beträf-  
fande dessa är af hr G. Hedberg.

Diplomen lyda på följande vis. I fysik:

Kongliga Svenska Vetenskaps-Akademien  
har vid sitt sammansättande den 12 nov. 1901  
i enlighet med föreskrifterna i det af

ALFRED NOBEL

den 27 nov. 1895 upprättade testamente be-  
slutat att tilldela det pris som detta är  
bortgivet af den, som inom fysikens om-  
råden har gjort den viktigaste upptäckt eller  
upfinning, till

WILHELM CONRAD RÖNTGEN,

såsom ett erkännande af den utomordentliga  
förtjänst han inlägt genom upptäckten af  
de egendomliga strålarna som sedermera up-  
skalats efter honom.

Stockholm den 19 dec. 1901.

C. T. Odhner,  
K. Vet. Åks Præs.

Chr. Aurivillius,  
K. Vet. Åks sekreterare.

I kemi (början som föregående):

Kongliga Svenska Vetenskaps-Akademien  
har vid sitt sammansättande den 12 nov. 1901  
i enlighet med föreskrifterna i det af

ALFRED NOBEL

den 27 nov. 1895 upprättade testamente be-  
slutat att tilldela det pris som detta är  
bortgivet åt den, som har gjort den viktigaste  
kemiiska upptäckt eller förbättring, till

JACOBUS HERMANNUS VANT HOFF

såsom ett erkännande af den utomordentliga  
förtjänst han inlägt genom upptäckten af  
lagarna för den kemiiska dynamiken och för  
det osmotiska trycket i lösningsar

Stockholm den 19 dec. 1901.

G. T. Odhner,  
K. Vet. Åks Præs.

Chr. Aurivillius,  
K. Vet. Åks sekreterare.

I medicin:

K. Karolinska Mediko-Birurgiska Institu-  
tel, hvilket enligt testamente, som den 27  
nov. 1895 upprättades af

ALFRED NOBEL

sig att med Nobelprius belöna den viktigaste

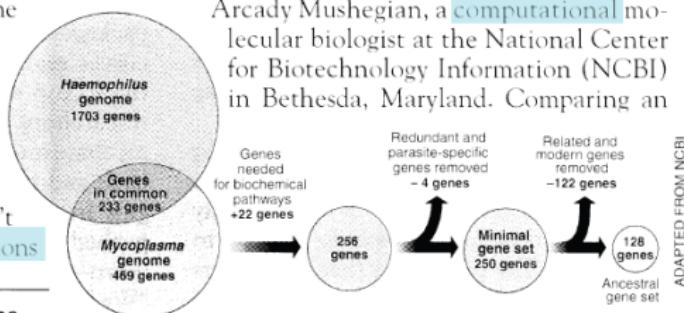
# Topic modelling

## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



ADAPTED FROM NCBI

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

# Topic modelling

## Topics

gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

data 0.02  
number 0.02  
computer 0.01  
...

## Documents

## Topic proportions and assignments

### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two dozen researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

“are not all that far apart,” especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson, a postdoctoral fellow at the University of Stockholm who arrived at the 800 number. But coming up with a consensus answer may be more than just a numbers game. As more and more genomes are mapped and sequenced, “It may be a way of organizing any newly sequenced genome,” explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing all

*Homo sapiens* genome: 27,000 genes  
Genes in common: 12,000 genes  
Stripped genome: 489 genes

*Paramecium* genome: 2,000 genes  
Genes in common: 1,000 genes  
Stripped genome: 128 genes

*Minimal gene set*: 250 genes  
Genes in common: 128 genes  
Stripped genome: 128 genes

*All life*: 128 genes  
Genes in common: 128 genes

*Stripping down*. Computer analysis yields an estimate of the minimum modern and ancient genomes.

\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

SCIENCE • VOL. 272 • 24 MAY 1996

<sup>1</sup>Credit to D. Blei for the picture

# Topic modelling

Topics



Documents

## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,<sup>1</sup> two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

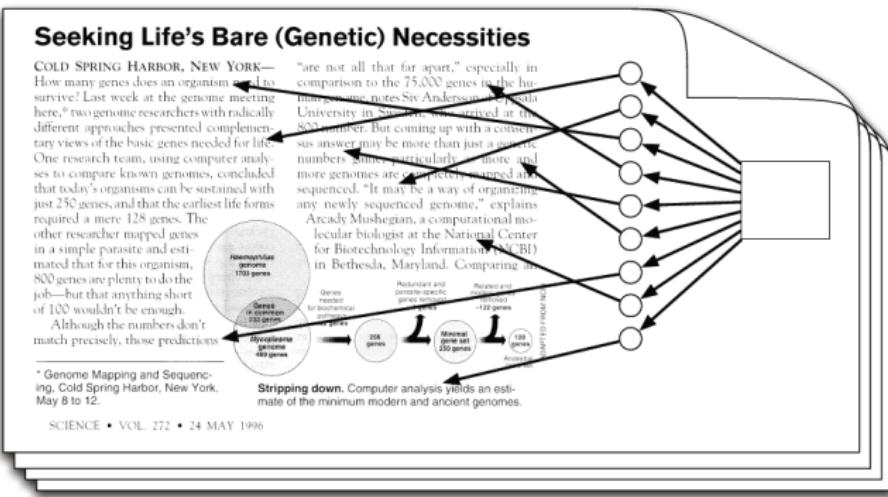
Although the numbers don't match precisely, those predictions

<sup>1</sup> Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Sri Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a academic numbers game, particularly if more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an-

arcade

Topic proportions and assignments



<sup>1</sup>Credit to D. Blei for the picture

# Probabilistic modeling

Directly trying to model topics from data is a hard problem!

# Probabilistic modeling

Directly trying to model topics from data is a hard problem!

- In **probabilistic modelling** we construct a **generative probabilistic process** that could have generated the data.  
In our case we let the topic assignments be a **hidden variable**.

# Probabilistic modeling

Directly trying to model topics from data is a hard problem!

- In **probabilistic modelling** we construct a **generative probabilistic process** that could have generated the data.  
In our case we let the topic assignments be a **hidden variable**.
- We then infer the structure using **posterior inference**.

# Probabilistic modeling

Directly trying to model topics from data is a hard problem!

- In **probabilistic modelling** we construct a **generative probabilistic process** that could have generated the data.  
In our case we let the topic assignments be a **hidden variable**.
- We then infer the structure using **posterior inference**.

Instead of directly modelling the topics, we model how to generate a corpus and then using posterior inference to find the topics.

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,
  - a. Draw a **distribution over words** for that topic  $\beta_k \sim \text{Dir}(\eta)$

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,
  - a. Draw a **distribution over words** for that topic  $\beta_k \sim \text{Dir}(\eta)$
2. For each document,  $d = 1, \dots, D$

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,
  - a. Draw a **distribution over words** for that topic  $\beta_k \sim \text{Dir}(\eta)$
2. For each document,  $d = 1, \dots, D$ 
  - a. Draw a vector of **topic proportions**  $\theta_d \sim \text{Dir}(\alpha)$

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,
  - a. Draw a **distribution over words** for that topic  $\beta_k \sim \text{Dir}(\eta)$
2. For each document,  $d = 1, \dots, D$ 
  - a. Draw a vector of **topic proportions**  $\theta_d \sim \text{Dir}(\alpha)$
  - b. For each **word position**,  $n = 1, \dots, N_d$

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,
  - a. Draw a **distribution over words** for that topic  $\beta_k \sim \text{Dir}(\eta)$
2. For each document,  $d = 1, \dots, D$ 
  - a. Draw a vector of **topic proportions**  $\theta_d \sim \text{Dir}(\alpha)$
  - b. For each **word position**,  $n = 1, \dots, N_d$ 
    - i. Draw a **topic assignment**  $z_{d,n} \sim \text{Cat}(\theta_d)$

# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,
  - a. Draw a **distribution over words** for that topic  $\beta_k \sim \text{Dir}(\eta)$
2. For each document,  $d = 1, \dots, D$ 
  - a. Draw a vector of **topic proportions**  $\theta_d \sim \text{Dir}(\alpha)$
  - b. For each **word position**,  $n = 1, \dots, N_d$ 
    - i. Draw a **topic assignment**  $z_{d,n} \sim \text{Cat}(\theta_d)$
    - ii. Draw a **word**  $w_{d,n} \sim \text{Cat}(\beta_{z_{d,n}})$

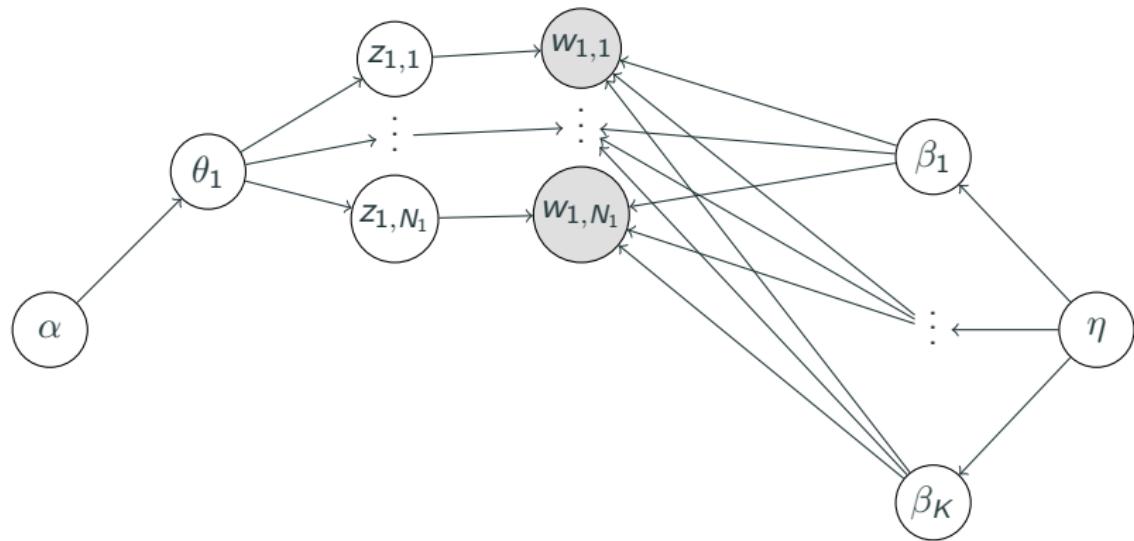
# Probabilistic document modelling

Some notation: We will use bold symbols  $\theta$  for the collection of all the random variables, the first sub index  $\theta_i$  will result in a vector and the second subindex  $\theta_{i,j}$  the position in that vector.

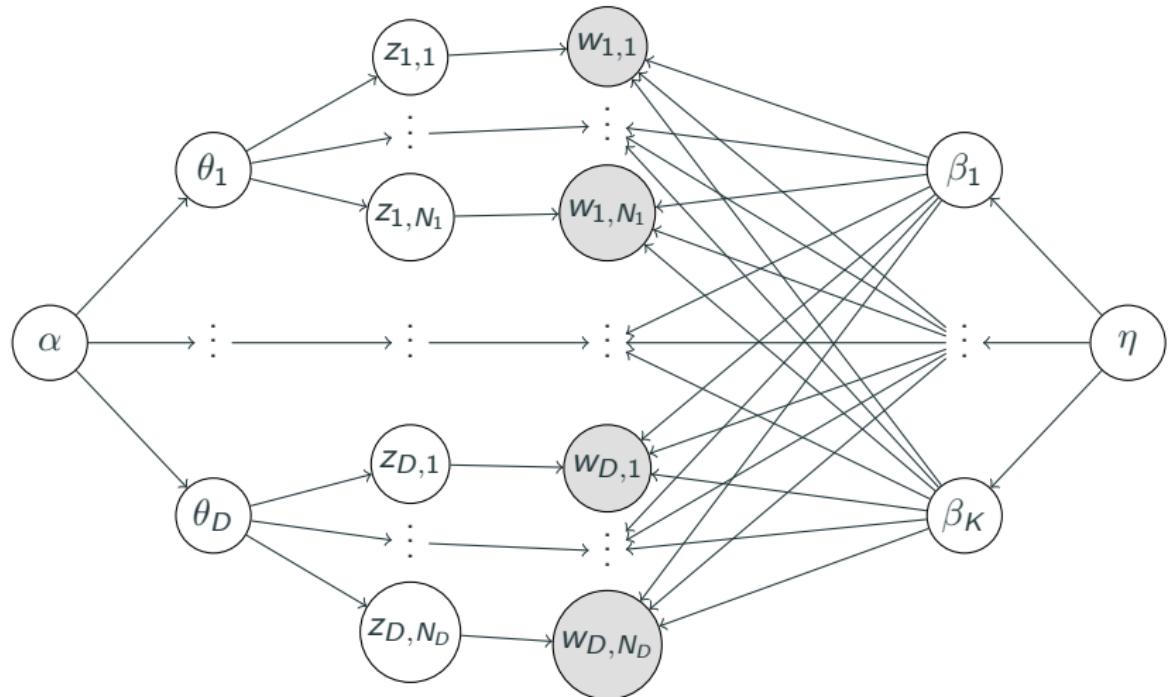
We begin by describing the **generative probabilistic process**.

1. For each topic,  $k = 1, \dots, K$ ,
  - a. Draw a **distribution over words** for that topic  $\beta_k \sim \text{Dir}(\eta)$   
—  $1 \times V$  **probability vector**
2. For each document,  $d = 1, \dots, D$ 
  - a. Draw a vector of **topic proportions**  $\theta_d \sim \text{Dir}(\alpha)$   
—  $1 \times K$  **probability vector**
  - b. For each **word position**,  $n = 1, \dots, N_d$ 
    - i. Draw a **topic assignment**  $z_{d,n} \sim \text{Cat}(\theta_d)$  — **integer** between 1 and  $K$
    - ii. Draw a **word**  $w_{d,n} \sim \text{Cat}(\beta_{z_{d,n}})$  — **integer** between 1 and  $V$

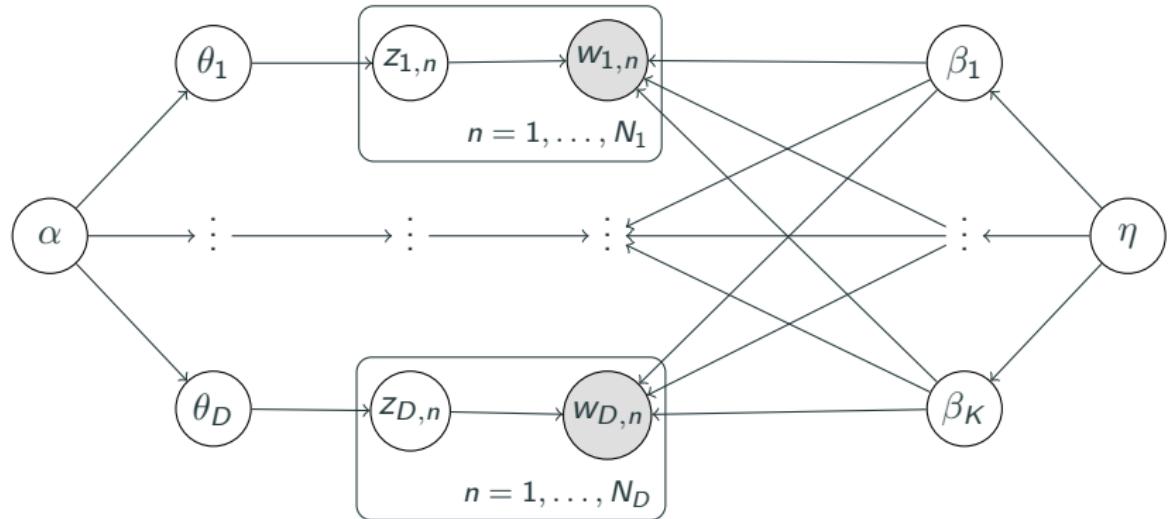
# Graphical model



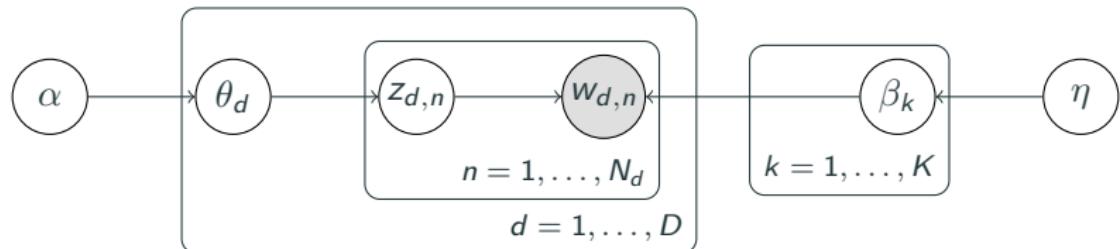
# Graphical model



# Graphical model



# Graphical model



This is known as **plate notation** where each box indicates repetition of those nodes the number of times stated.

## The problem

Given a corpus of documents. We wish to infer the distribution of topics on each document and the words within each topic.

## The problem

Given a corpus of documents. We wish to infer the distribution of topics on each document and the words within each topic.

Specifically we wish to calculate the **posterior** distribution

$$p(\theta, z, \beta | w, \eta, \alpha) = \frac{p(\theta, z, \beta, w | \eta, \alpha)}{p(w | \eta, \alpha)}$$

## The problem

Given a corpus of documents. We wish to infer the distribution of topics on each document and the words within each topic.

Specifically we wish to calculate the **posterior** distribution

$$p(\theta, z, \beta | w, \eta, \alpha) = \frac{p(\theta, z, \beta, w | \eta, \alpha)}{p(w | \eta, \alpha)}$$

Before going in to how to sample from that distribution, let's see what we can do when we have it. Introduce:

$$\hat{\beta}_k = \mathbb{E}[\beta_k | w]$$

$$\hat{\theta}_d = \mathbb{E}[\theta_d | w]$$

$$\hat{z}_{d,n}^k = \mathbb{E}[z_{d,n} = k | w]$$

With these quantities we can explore the corpus in the following way.

## Exploring a corpus

**Visualizing a topic:** The simplest way to visualize a topic is to look at the estimated probabilities in  $\hat{\beta}_k$  and order them.

## Exploring a corpus

**Visualizing a topic:** The simplest way to visualize a topic is to look at the estimated probabilities in  $\hat{\beta}_k$  and order them. A better method, inspired by the **term frequency–inverse document frequency** (TFIDF) score, is the following

$$\text{term-score}_{k,v} = \hat{\beta}_{k,v} \log \left( \frac{\hat{\beta}_{k,v}}{\left( \prod_{j=1}^K \hat{\beta}_{j,v} \right)^{1/K}} \right)$$

The first term gives the **term frequency** and is then scaled by weighting down terms that are popular in many topics.

**Finding similar documents:** Using the **posterior topic proportions**  $\hat{\theta}_d$

we have a representation of the document.

To find similar documents we would like to compare these proportions. The Hellinger distance can be used to calculate such a similarity

$$\text{document-similarity}_{d,f} = \sum_{k=1}^K \left( \sqrt{\hat{\theta}_{d,k}} - \sqrt{\hat{\theta}_{f,k}} \right)^2.$$

**Finding similar documents:** Using the **posterior topic proportions**  $\hat{\theta}_d$  we have a representation of the document.

To find similar documents we would like to compare these proportions. The Hellinger distance can be used to calculate such a similarity

$$\text{document-similarity}_{d,f} = \sum_{k=1}^K \left( \sqrt{\hat{\theta}_{d,k}} - \sqrt{\hat{\theta}_{f,k}} \right)^2.$$

**Visualizing a document:** Finally we would like to use the estimates  $\hat{\theta}_d$  and  $\hat{z}_{d,n}^k$  to visualize the document.

# Exploring a corpus

## Chance and Statistical Significance in Protein and DNA Sequence Analysis

Samuel Karlin and Volker Brendel

Top words from the top topics (by term score)

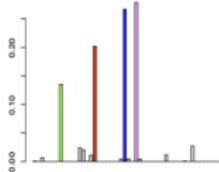
sequence  
region  
pcr  
identified  
fragments  
two  
genes  
three  
cdna  
analysis

measured  
average  
range  
values  
different  
size  
three  
calculated  
two  
low

residues  
binding  
domains  
helix  
cys  
regions  
structure  
terminus  
terminal  
site

computer  
methods  
number  
two  
principle  
design  
access  
processing  
advantage  
important

Expected topic proportions



Abstract with the most likely topic assignments

Statistical approaches help in the determination of significant configurations in protein and nucleic acid sequence data. Three recent statistical methods are discussed: (i) score-based sequence analysis that provides a means for characterizing anomalies in local sequence text and for evaluating sequence comparisons; (ii) quantile distributions of amino acid usage that reveal general compositional biases in proteins and evolutionary relations; and (iii) r-scan statistics that can be applied to the analysis of spacings of sequence markers.

Top Ten Similar Documents

- Exhaustive Matching of the Entire Protein Sequence Database
- How Big Is the Universe of Exons?
- Counting and Discounting the Universe of Exons
- Detecting Subtle Sequence Signals: A Gibbs Sampling Strategy for Multiple Alignment
- Ancient Conserved Regions in New Gene Sequences and the Protein Databases
- A Method to Identify Protein Sequences that Fold into a Known Three-Dimensional Structure
- Testing the Exon Theory of Genes: The Evidence from Protein Structure
- Predicting Coiled Coils from Protein Sequences
- Genome Sequence of the Nematode *C. elegans*: A Platform for Investigating Biology

2

<sup>2</sup>D. Blei and J. Lafferty. Topic models. In A. Srivastava and M. Sahami, editors, *Text Mining: Theory and Applications*. Taylor and Francis, 2009.

9/16

## Posterior inference in LDA

We wish to calculate expected values under the posterior distribution

$$p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta} | \mathbf{w}, \eta, \alpha) = \frac{p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta}, \mathbf{w} | \eta, \alpha)}{p(\mathbf{w} | \eta, \alpha)}$$

## Posterior inference in LDA

We wish to calculate expected values under the posterior distribution

$$p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta} | \mathbf{w}, \eta, \alpha) = \frac{p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta}, \mathbf{w} | \eta, \alpha)}{p(\mathbf{w} | \eta, \alpha)}$$

Looking at the numerator we have

$$p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta}, \mathbf{w} | \eta, \alpha) = p(\boldsymbol{\beta} | \eta)p(\boldsymbol{\theta} | \alpha)p(\mathbf{z} | \boldsymbol{\theta})p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta})$$

## Posterior inference in LDA

We wish to calculate expected values under the posterior distribution

$$p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta} | \mathbf{w}, \eta, \alpha) = \frac{p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta}, \mathbf{w} | \eta, \alpha)}{p(\mathbf{w} | \eta, \alpha)}$$

Looking at the numerator we have

$$\begin{aligned} p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta}, \mathbf{w} | \eta, \alpha) &= p(\boldsymbol{\beta} | \eta)p(\boldsymbol{\theta} | \alpha)p(\mathbf{z} | \boldsymbol{\theta})p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta}) \\ &= \prod_{k=1}^K p(\beta_k | \eta) \prod_{d=1}^D \left[ p(\theta_d | \alpha) \prod_{n=1}^{N_d} [p(z_{d,n} | \theta_d)p(w_{d,n} | z_{d,n}, \beta_{z_{d,n}})] \right] \end{aligned}$$

## Posterior inference in LDA

We wish to calculate expected values under the posterior distribution

$$p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta} | \mathbf{w}, \eta, \alpha) = \frac{p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta}, \mathbf{w} | \eta, \alpha)}{p(\mathbf{w} | \eta, \alpha)}$$

Looking at the numerator we have

$$\begin{aligned} p(\boldsymbol{\theta}, \mathbf{z}, \boldsymbol{\beta}, \mathbf{w} | \eta, \alpha) &= p(\boldsymbol{\beta} | \eta)p(\boldsymbol{\theta} | \alpha)p(\mathbf{z} | \boldsymbol{\theta})p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta}) \\ &= \prod_{k=1}^K p(\beta_k | \eta) \prod_{d=1}^D \left[ p(\theta_d | \alpha) \prod_{n=1}^{N_d} [p(z_{d,n} | \theta_d)p(w_{d,n} | z_{d,n}, \beta_{z_{d,n}})] \right] \end{aligned}$$

For the denominator we have

$$p(\mathbf{w} | \alpha, \eta) = \int \int \sum_{\mathbf{z}} p(\boldsymbol{\beta} | \eta)p(\boldsymbol{\theta} | \alpha)p(\mathbf{z} | \boldsymbol{\theta})p(\mathbf{w} | \mathbf{z}, \boldsymbol{\beta}) d\boldsymbol{\theta} d\boldsymbol{\beta},$$

which is intractable.

## Gibbs sampling in LDA

We will now look at the different conditionals needed for the Gibbs sampler.

- $\theta$  We notice that the Dirichlet is conjugate prior to the categorical and have that

$$\theta_d | \beta, \mathbf{z}, \mathbf{w}, \alpha, \eta \sim \text{Dir}(\alpha + c_d)$$

- $\beta$  We notice that the Dirichlet is conjugate prior to the categorical and have that

$$\beta_k | \theta, \mathbf{z}, \mathbf{w}, \alpha, \eta \sim \text{Dir}(\eta + \tilde{c}_k)$$

- $\mathbf{z}$  This is a discrete random variable, so we calculate the probabilities directly.

$$p(z_{d,n} = k | \theta, \beta, \mathbf{w}, \alpha, \eta) \propto \theta_{d,k} \beta_{k,w_{d,n}}$$

Here we have that the elements of  $c_d$  and  $\tilde{c}_k$  is given by

$$c_{d,k} = \sum_n \mathbb{1}\{z_{d,n} = k\}, \quad \tilde{c}_{k,v} = \sum_d \sum_n \mathbb{1}\{w_{d,n} = v\} \mathbb{1}\{z_{d,n} = k\}$$

## Collapsed Gibbs sampling

---

- Performing Gibbs in this way requires us (at each iteration) to sample from  $D \times K + K \times V + \sum_d N_d$  variables. Sampling from such a high-dimensional distribution using Gibbs sampling usually results in poor mixing.

## Collapsed Gibbs sampling

---

- Performing Gibbs in this way requires us (at each iteration) to sample from  $D \times K + K \times V + \sum_d N_d$  variables. Sampling from such a high-dimensional distribution using Gibbs sampling usually results in poor mixing.
- To improve the mixing of the algorithm we can **marginalize** some of the variables and sample from the remaining ones.
  - That is instead of sampling from  $p(\theta, z, \beta | w, \eta, \alpha)$  we only sample from  $p(z | w, \eta, \alpha)$
  - This will help with the mixing.
  - **Introduces dependencies** between the topic variables.

## Collapsed Gibbs sampling

---

- Performing Gibbs in this way requires us (at each iteration) to sample from  $D \times K + K \times V + \sum_d N_d$  variables. Sampling from such a high-dimensional distribution using Gibbs sampling usually results in poor mixing.
- To improve the mixing of the algorithm we can **marginalize** some of the variables and sample from the remaining ones.
  - That is instead of sampling from  $p(\theta, z, \beta | w, \eta, \alpha)$  we only sample from  $p(z | w, \eta, \alpha)$
  - This will help with the mixing.
  - **Introduces dependencies** between the topic variables.
- This can be compared to **Rao-Blackwellisation** where marginalization over variables help with estimation.

## Collapsed Gibbs sampling in LDA

In the LDA model we can integrate  $\theta$  and  $\beta$  and just keep the topic choices  $z$ . We only need to sample the  $z$  from the resulting marginalized distribution. This is known as **collapsed Gibbs sampling**.

For our collapsed Gibbs sampling we need

$$\begin{aligned} p(z_{d,n} = k | \mathbf{z}^{-(d,n)}, \mathbf{w}, \alpha, \eta) \\ \propto p(z_{d,n} = k | \mathbf{z}^{-(d,n)}, \alpha) p(w_{d,n} | z_{d,n} = k, \mathbf{w}^{-(d,n)}, \mathbf{z}^{-(d,n)}, \eta) \\ \propto \frac{\alpha + c_{d,k}^{-(d,n)}}{\sum_{k'=1}^K \alpha + c_{d,k'}^{-(d,n)}} \times \frac{\eta + \tilde{c}_{k,w_{d,n}}^{-(d,n)}}{\sum_{v=1}^V \eta + \tilde{c}_{k,v}^{-(d,n)}} \end{aligned}$$

Given a set of  $\mathbf{z}$  we can then estimate  $\hat{\theta}_d$  and  $\hat{\beta}_k$  by

$$\hat{\theta}_{d,k} = \frac{\alpha + c_{d,k}}{\sum_{k'=1}^K \alpha + c_{d,k'}} \quad \hat{\beta}_{k,v} = \frac{\eta + \tilde{c}_{k,v}}{\sum_{v'=1}^V \eta + \tilde{c}_{k,v'}}$$

# The algorithm

---

- Initialize  $\mathbf{z}$  randomly.
- Create all the counting matrices.
- At every iteration:
  - For every document  $d = 1, \dots, D$  and every word  $n = 1, \dots, N_d$ :
    - Set  $v = w_{d,n}$
    - Set  $k = z_{d,n}$
    - Set  $c_{d,k} -= 1$
    - Set  $\tilde{c}_{k,v} -= 1$
    - Sample  $k \sim \text{Cat}(\{p(z_{d,n} = k' | \mathbf{z}^{-(d,n)}, \mathbf{w}, \alpha, \eta)\}_{k'=1}^K)$
    - Set  $c_{d,k} += 1$
    - Set  $\tilde{c}_{k,v} += 1$

## Pre-processing of the documents

Before we can perform this algorithm, some pre-processing of the documents are needed.

**Tokenization:** Take each document, split it into words. Change all letters to lowercase and remove punctuation.

**Small words:** Remove words with fewer than 3 characters.

**Stopwords:** Remove any stopwords. These are some of the most common words for the language.

**Lemmatize:** Words in third person are changed to first person and verbs are changed into present.

**Stem:** Words are reduced to their root form.

After this is done you construct your dictionary.

There are many packages that will handle all of these steps for you, no need to reinvent the wheel.

# Extensions

- We don't have to work with text documents. There is nothing in the model that stops it from being used in images.

## Extensions

- We don't have to work with text documents. There is nothing in the model that stops it from being used in images.
- We can study how topics change over time, by splitting the dataset and constructing one LDA model for each time. The distributions over words  $\beta_k$  are then connected sequentially.

## Extensions

- We don't have to work with text documents. There is nothing in the model that stops it from being used in images.
- We can study how topics change over time, by splitting the dataset and constructing one LDA model for each time. The distributions over words  $\beta_k$  are then connected sequentially.
- You can probably think about something else that can be done with the model.