

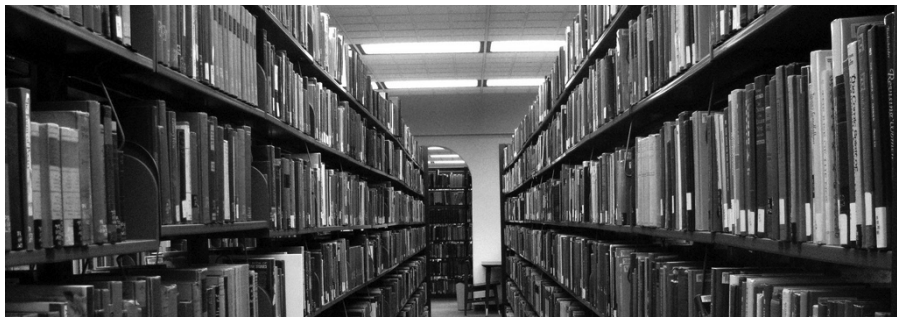
Topic Models for Understanding History

David M. Blei
Columbia University

joint work with Allison Chaney, Hanna Wallach, and Matt Connelly



- ▶ **ORGANIZE**
- ▶ **VISUALIZE**
- ▶ **SUMMARIZE**
- ▶ **SEARCH**
- ▶ **PREDICT**
- ▶ **UNDERSTAND**



TOPIC MODELING

1. **Discover** the thematic structure
2. **Annotate** the documents
3. **Use** the annotations to visualize, organize, summarize, ...

1

Game
Season
Team
Coach
Play
Points
Games
Giants
Second
Players

2

Life
Know
School
Street
Man
Family
Says
House
Children
Night

3

Film
Movie
Show
Life
Television
Films
Director
Man
Story
Says

4

Book
Life
Books
Novel
Story
Man
Author
House
War
Children

5

Wine
Street
Hotel
House
Room
Night
Place
Restaurant
Park
Garden

6

Bush
Campaign
Clinton
Republican
House
Party
Democratic
Political
Democrats
Senator

7

Building
Street
Square
Housing
House
Buildings
Development
Space
Percent
Real

8

Won
Team
Second
Race
Round
Cup
Open
Game
Play
Win

9

Yankees
Game
Mets
Season
Run
League
Baseball
Team
Games
Hit

10

Government
War
Military
Officials
Iraq
Forces
Iraqi
Army
Troops
Soldiers

11

Children
School
Women
Family
Parents
Child
Life
Says
Help
Mother

12

Stock
Percent
Companies
Fund
Market
Bank
Investors
Funds
Financial
Business

13

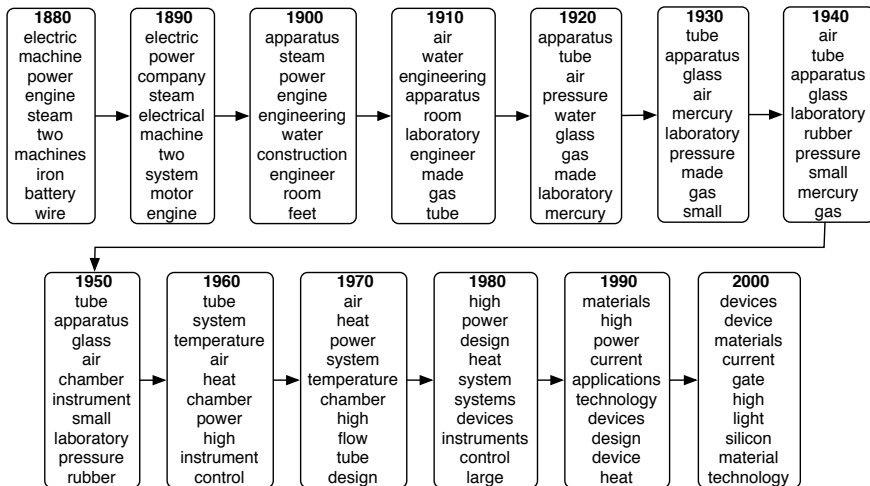
Church
War
Women
Life
Black
Political
Catholic
Government
Jewish
Pope

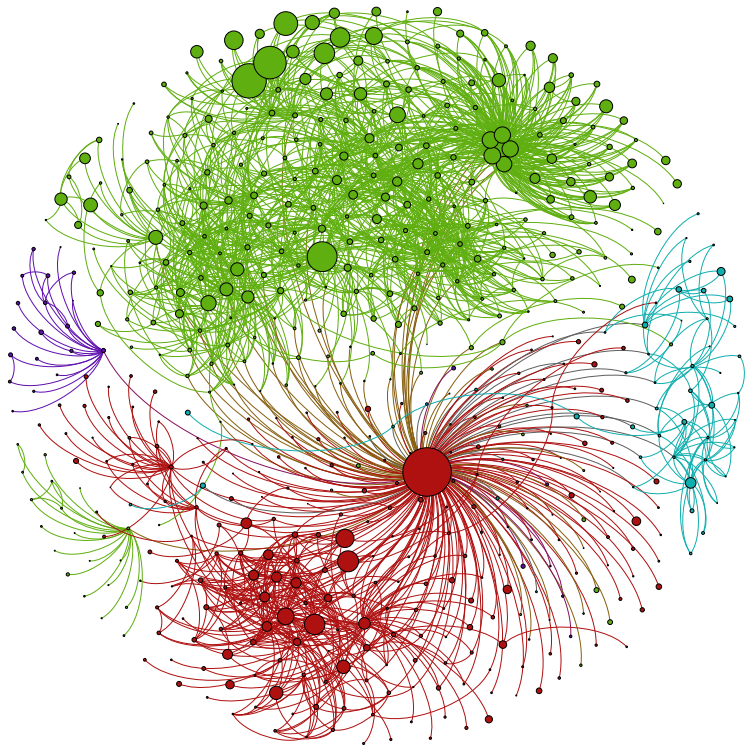
14

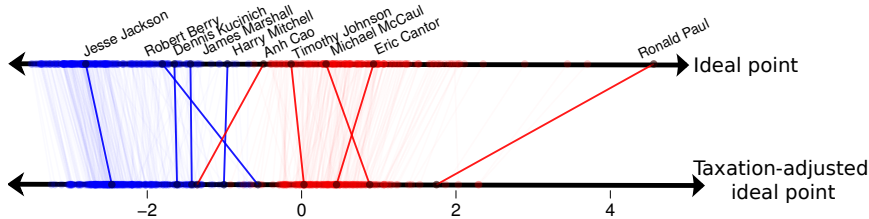
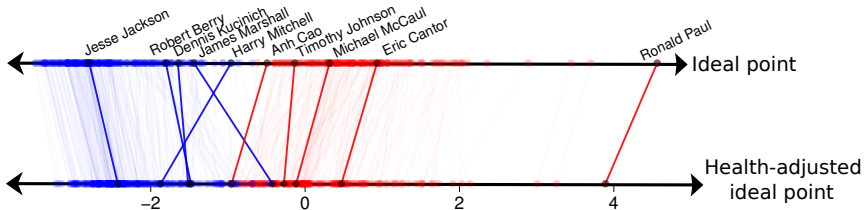
Art
Museum
Show
Gallery
Works
Artists
Street
Artist
Paintings
Exhibition

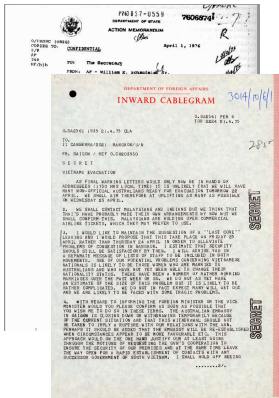
15

Police
Yesterday
Man
Officer
Officers
Case
Found
Charged
Street
Shot









- ▶ Historians want to identify important events from primary sources.
- ▶ Example: Embassies send cables to each other during the 1970s
- ▶ Goal: Use topic models to discover **events** in this data set

This talk

1. Introduction to topic modeling
2. Topic models for understanding history
3. The bigger picture: Using probability models to solve problems with data

Introduction to Topic Modeling

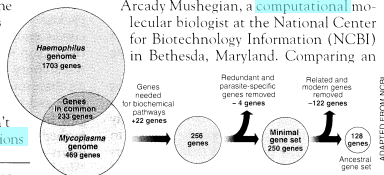
Seeking Life's Bare (Genetic) Necessities

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



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

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

SCIENCE • VOL. 272 • 24 MAY 1996

Documents exhibit multiple topics.

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
...

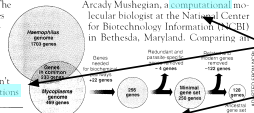
data 0.02
number 0.02
computer 0.01
...

Documents

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, a Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a little 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



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

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

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments

Latent Dirichlet Allocation

Topics



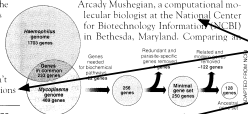
Documents

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,* 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, Uppsala University in Sweden, who lectured at the meeting. 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 the

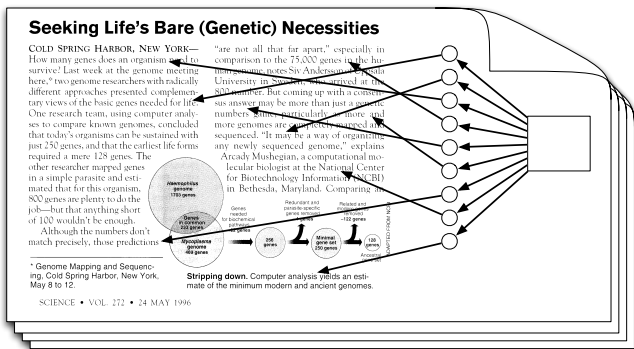


* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

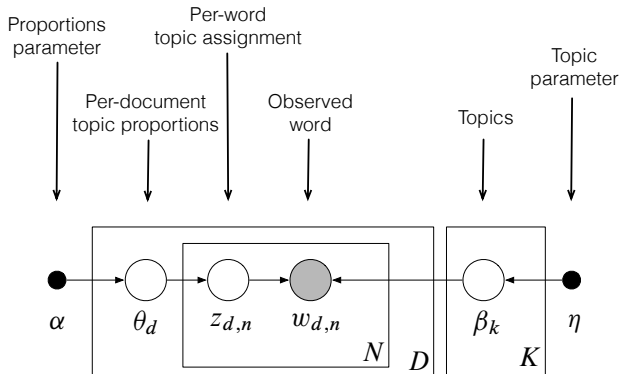
Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments

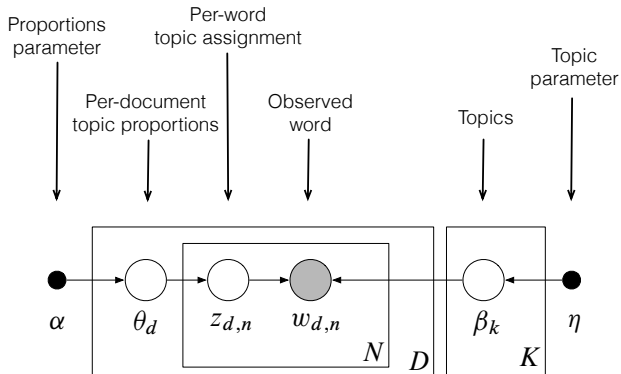


Latent Dirichlet Allocation



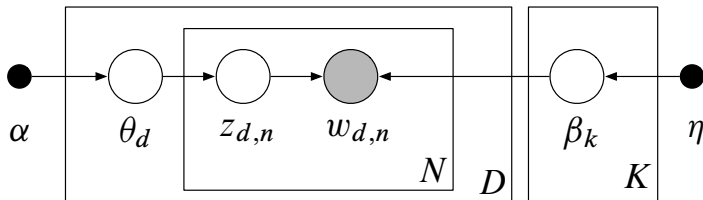
LDA as a graphical model

- Nodes are random variables; edges indicate dependence.
- Shaded nodes are observed; unshaded nodes are hidden.
- Plates indicate replicated variables.

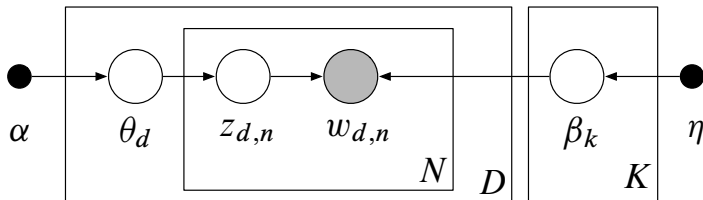


LDA as a graphical model

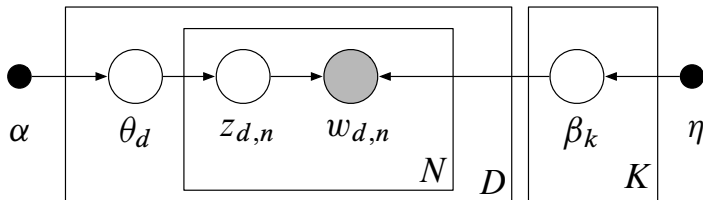
- Encodes independence assumptions about the variables
- Defines a factorization of the joint probability distribution
- Connects to algorithms for computing with data



- The joint defines a posterior, $p(\theta, z, \beta \mid w)$.
- From a collection of documents, infer
 - Per-word topic assignment $z_{d,n}$
 - Per-document topic proportions θ_d
 - Per-corpus topic distributions β_k
- Then use posterior expectations to perform the task at hand: information retrieval, document similarity, exploration, and others.



- ▶ Mean field variational methods (Blei et al., 2001, 2003)
- ▶ Expectation propagation (Minka and Lafferty, 2002)
- ▶ Collapsed Gibbs sampling (Griffiths and Steyvers, 2002)
- ▶ Distributed sampling (Newman et al., 2008; Ahmed et al., 2012)
- ▶ Collapsed variational inference (Teh et al., 2006)
- ▶ Stochastic inference (Hoffman et al., 2010, 2013; Mimno et al., 2012)
- ▶ Factorization inference (Arora et al., 2012; Anandkumar et al., 2012)



- ▶ LDA in R [<https://cran.r-project.org/web/packages/lda/>]
- ▶ GenSim [<https://radimrehurek.com/gensim>]
- ▶ Mallet [<http://mallet.cs.umass.edu>]
- ▶ Vowpal Wabbit [<http://hunch.net/~vw/>]
- ▶ Apache Spark [<http://spark.apache.org/>]
- ▶ SciKit Learn [<http://scikit-learn.org/>]



- ▶ **Data:** The OCR'ed collection of *Science* from 1990–2000
 - 17K documents
 - 11M words
 - 20K unique terms (stop words and rare words removed)
- ▶ **Model:** 100-topic LDA model using variational inference.

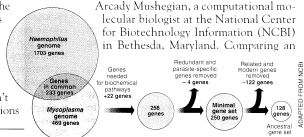
Seeking Life's Bare (Genetic) Necessities

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.

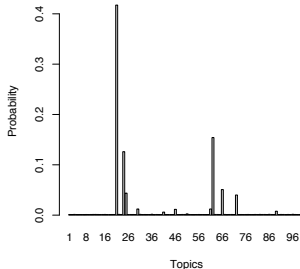
Although the numbers don't match precisely, those predictions

"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



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.

human	evolution	disease	computer
genome	evolutionary	host	models
dna	species	bacteria	information
genetic	organisms	diseases	data
genes	life	resistance	computers
sequence	origin	bacterial	system
gene	biology	new	network
molecular	groups	strains	systems
sequencing	phylogenetic	control	model
map	living	infectious	parallel
information	diversity	malaria	methods
genetics	group	parasite	networks
mapping	new	parasites	software
project	two	united	new
sequences	common	tuberculosis	simulations

1

Game
Season
Team
Coach
Play
Points
Games
Giants
Second
Players

2

Life
Know
School
Street
Man
Family
Says
House
Children
Night

3

Film
Movie
Show
Life
Television
Films
Director
Man
Story
Says

4

Book
Life
Books
Novel
Story
Man
Author
House
War
Children

5

Wine
Street
Hotel
House
Room
Night
Place
Restaurant
Park
Garden

6

Bush
Campaign
Clinton
Republican
House
Party
Democratic
Political
Democrats
Senator

7

Building
Street
Square
Housing
House
Buildings
Development
Space
Percent
Real

8

Won
Team
Second
Race
Round
Cup
Open
Game
Play
Win

9

Yankees
Game
Mets
Season
Run
League
Baseball
Team
Games
Hit

10

Government
War
Military
Officials
Iraq
Forces
Iraqi
Army
Troops
Soldiers

11

Children
School
Women
Family
Parents
Child
Life
Says
Help
Mother

12

Stock
Percent
Companies
Fund
Market
Bank
Investors
Funds
Financial
Business

13

Church
War
Women
Life
Black
Political
Catholic
Government
Jewish
Pope

14

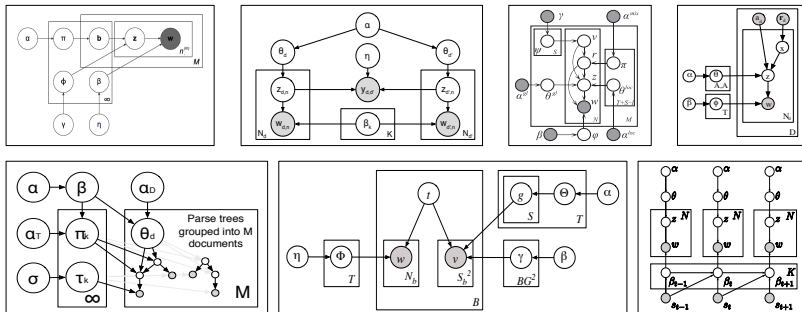
Art
Museum
Show
Gallery
Works
Artists
Street
Artist
Paintings
Exhibition

15

Police
Yesterday
Man
Officer
Officers
Case
Found
Charged
Street
Shot

How does LDA “work”?

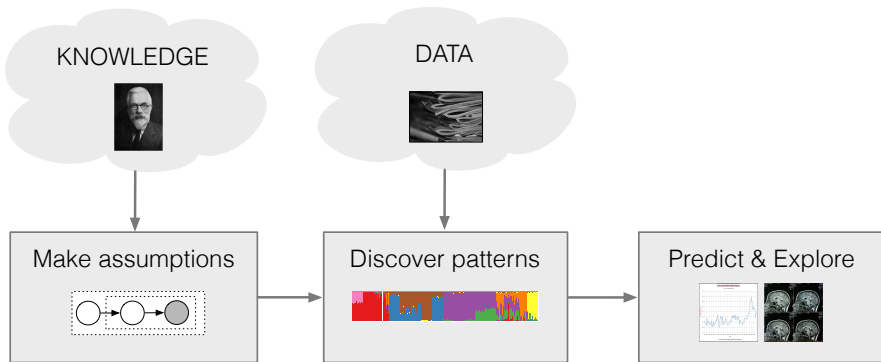
- ▶ LDA trades off two goals.
 1. In each **document**, allocate its words to **few topics**.
 2. In each **topic**, assign high probability to **few terms**.
- ▶ These goals are at odds.
 - Putting a document in a single topic makes #2 hard:
All of its words must have probability under that topic.
 - Putting very few words in each topic makes #1 hard:
To cover a document’s words, it must assign many topics to it.
- ▶ Trading off these goals finds groups of tightly co-occurring words.



- ▶ Organizing and finding patterns in text is important in the sciences, humanities, industry, and culture.
- ▶ LDA is a simple building block that enables many applications. Topic modeling is an active field of research.
- ▶ Algorithmic improvements let us fit models to massive data.

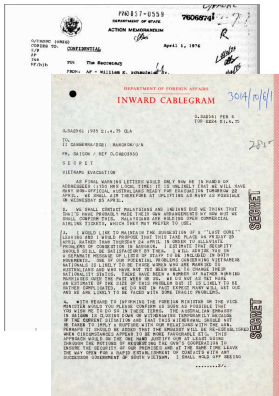
- LDA is a simple building block that enables many applications. Topic modeling is an active field of research.

- Algorithmic improvements let us fit models to massive data.



- ▶ Case study in **text analysis with probability models**
- ▶ Topic modeling research
 - develops new models.
 - develops new inference algorithms.
 - develops new applications, visualizations, tools.

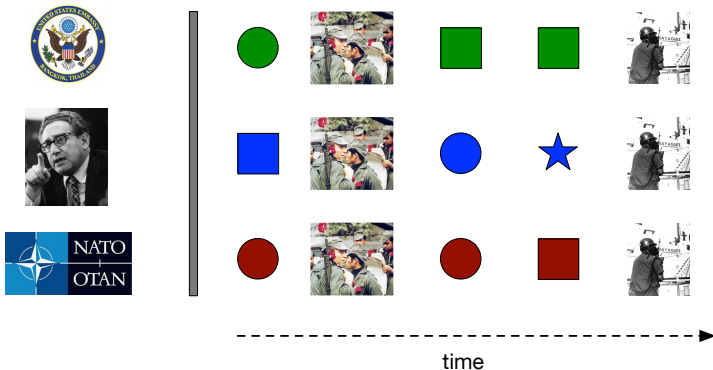
Topic Models for Understanding History



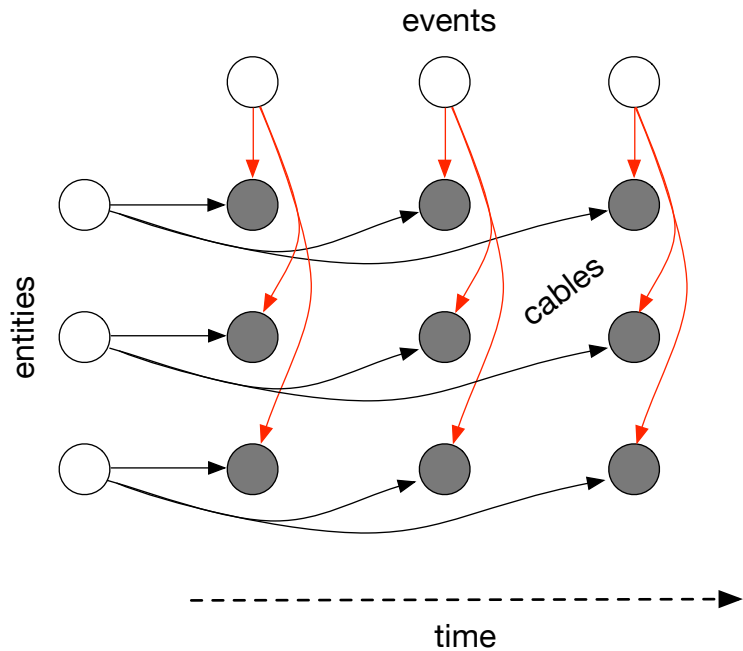
- Historians want to identify important events from primary sources.
- Example: Embassies send cables to each other during the 1970s
- Goal: Use topic models to find events in this data set

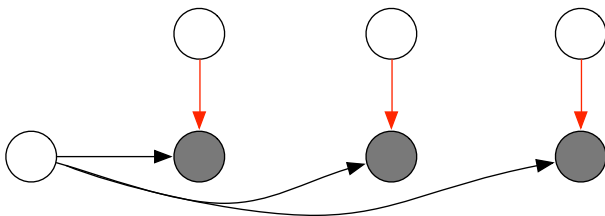
right human church freedom person inquiry violation religious discrimination cooper	minister depart arrive schedule party hour departure president reception airport	bank exchange room reserve central governor rate finance financial single	negotiation proposal position agreement question agree point negotiate issue western
fish bill vessel fishery zone airport mile fare dote water	university student health medical education professor school child american care	build construction facility plant unit extension supply area work cost	refugee status personnel resettlement name family parole swiss check grantee
soviet moscow ussr union brezhnev detente russian side pravda gromyko	control narcotic drug traffic rangoon indonesian extradition enforcement opium attorney	press article story news interview medium statement coverage carry american	arab israeli east talk egyptian middle peace minister palestinian settlement

Topic models (by themselves) are a start. But they don't identify events.



- ▶ Embassies typically discuss their *usual business*
- ▶ When a cable is about an **event**:
 - It diverges from the usual business of the sender
 - Multiple embassies discuss it
- ▶ *Usual business* is framed in terms of topics; **events** are framed in terms of words.



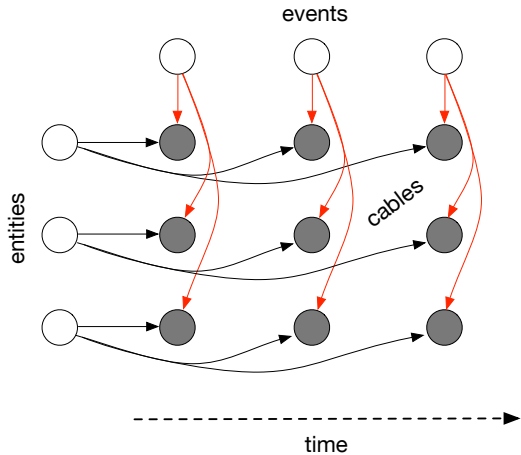


Hidden variables

- ▶ Topics
- ▶ Event description (per week)
- ▶ Topic description (per entity)
- ▶ Topic strength, event strength (per cable)

Observed variables

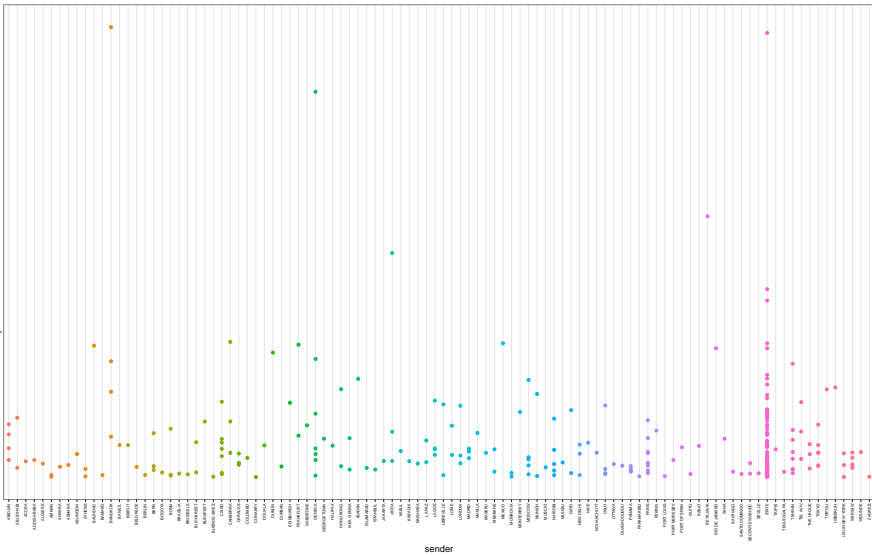
- ▶ Cables (per-week, per-entity)



To find events:

- ▶ Calculate the posterior of the hidden variables given the observed variables
- ▶ Examine cables where the event strength is high

top 5% of content about week 1976-07-01



1. THE FOLLOWING IS A COMPOSITE REPORT ON SOME OF THE ASPECTS OF THE RECENT HIJACKING OF THE AIR FRANCE AIRBUS BASED ON CONVERSATIONS WITH A NUMBER OF THE AMERICANS WHO WERE RELEASED. THEIR ACCOUNTS WERE GENERALLY CONGRUENT AND COINCIDE WITH THE REPORTS FROM THE GERMAN EMBASSY IN KAMPALA (REF. A). SINCE MANY ASPECTS ARE A MATTER OF IMPRESSION RATHER THAN FIRM OBSERVATION, THERE ARE SOME DIFFERENCES IN THE ACCOUNTS, WHICH WE NOTE. WE GO OVER SOME OF THE SAME GROUND COVERED IN THE ACCOUNTS IN REFS B, C AND D. THERE IS SOME CONFUSION OVER NOMENCLATURE FOR THE THREE ELEMENTS INVOLVED IN THE HIJACKING. ALL THE AMERICANS REFERRED TO MEMBERS OF ALL THREE AS PLO PRESUMABLY SINCE ENGLISH SPEAKING GERMANS WITH WHOM THEY HAD MOST CONTACT SAID THEY WORKED FOR PLO. BUT THE ARABS ON THE GROUND WHO WERE APPARENTLY IN CHARGE OF THE WHOLE CHE GUEVARA COMMANDO GROUP MADE THE DISTINCTION AND SPECIFIED THEY WORKED FOR PFLP AND NOT PLO. GERMANS PERHAPS DID NOT UNDERSTAND THE DISTINCTION. 2. RETAINED AMERICANS. SOME OF THE AMERICANS HAD IDENTIFIED GEORGE AND RENE KARFUNKEL AS AMERICANS AND WERE AWARE THAT THEY HAD ONLY AMERICAN PASSPORTS AND HAD BEEN IN ISRAEL ONLY BRIEFLY. NO ONE HAD A SATISFACTORY EXPLANATION OF WHY THE KARFUNKELS HAD NOT BEEN RELEASED WITH THE OTHER AMERICANS. THE KARFUNKELS APPEARED VERY ORTHODOX, ATE ONLY KOSHER FOOD...

1. FRG INFORMED EMBASSY AT 2130 LOCAL THAT FRENCH GOVERNMENT HAS PASSED IT THE FOLLOWING MESSAGE: A. THE HIJACKERS REJECT ANY EXCHANGE OUTSIDE OF ENTEBBE AIRPORT. B. THE EXCHANGE MUST TAKE PLACE UNDER THE SUPERVISION OF AMIN (OR OTHER HIGH UGANDAN OFFICIAL), TWO FRENCH REPS, AND SOMALI AMBASSADOR. C. HIJACKERS NOT PREPARED TO DISTINGUISH BETWEEN THE PRISONERS HELD IN ISRAEL AND THOSE HELD IN OTHER COUNTRIES. D. ALL HOSTAGES MUST BE EXCHANGED AGAINST ALL THE PRISONERS. E. THE HIJACKERS EXPECT AN ANSWER FROM ALL FOUR COUNTRIES HOLDING PRISONERS, NOT ONLY ISRAEL. F. THE HIJACKERS REFER TO THE COMPLETE LIST OF 53 "COMRADES." G. THE HIJACKERS INSIST AGAIN ON A PACKAGE DEAL: 53 "COMRADES" AGAINST ALL THE HOSTAGES AT ENTEBBE AIRPORT.

2. THE FRENCH PASSED A SECOND MESSAGE, THIS ONE FROM AMIN: AMIN TOLD THE FRENCH AMBASSADOR THAT HE EXPECTS ALL FOUR COUNTRIES TO COMMUNICATE TO HIM THE FLIGHT NUMBERS AND ETA OF ALL AIRCRAFT BRINGING PRISONERS TO UGANDA BEFORE THE END OF THE ULTIMATUM. (HE DID NOT SPECIFY AN HOUR.)

3. THE FRG CRISIS CENTER TOLD US THE FRG IS CONSULTING WITH THE OTHER GOVERNMENTS INVOLVED AT THE HIGHEST LEVEL. IT HAS NOT RPT NOT REACHED A DECISION ON RELEASE OF PRISONERS. HILLENBRAND

1. ACCORDING TO AS YET PROVISIONAL INFORMATION OBTAINED FROM IDF SPOKESMAN AND PASSENGERS, THERE WERE FOUR AMERICAN CITIZENS AMONG PASSENGERS FROM HIJACKED AIR FRANCE FLIGHT BROUGHT TO ISRAEL ON JULY 4. THEY ARE GEORGE AND RENE GARFUNKEL OF NEW YORK CITY; MRS. JANETTE ALMOG (HUSBAND ESRA ALMOG IS ISRAELI CITIZEN) OF MADISON WISCONSIN; AND MOSHE PERES OF NEW HAVEN CONNECTICUT. ALL RETURNED PASSENGERS ARE REPORTED WELL IN TEL AVIV. THOSE IN TRANSIT ARE BEING HOUSED AT PLAZA HOTEL IN TEL AVIV.
2. PASSENGERS KILLED DURING LIBERATION WERE JEAN-JACQUES MIMOUNI REPORTEDLY OF FRENCH NATIONALITY AND MRS IDA BOROCHOWITZ, AN ISRAELI OF RUSSIAN ORIGIN. ONE ISRAELI SOLDIER DIES IN FIGHT. NINE PERSONS REQUIRING MEDICAL CARE WERE LEFT IN NAIROBI DURING BRIEF STOP-OVER.
3. ACCORDING TO THE GARFUNKELS, FRENCH PILOT SAID THAT HIJACKERS ENTERED EMBARKATION AREA WITH SIX PACKAGES CLAIMING THAT THEY CONTAINED CANDY. AT TIME OF THEIR ENTRY ELECTRICITY ALLEGEDLY WENT OUT AND STOPPED SCREENING DEVICES FROM WORKING. RATHER THAN DELAY THE PLANE FOR EXAMINATION OF PACKAGES, THE SIX WERE HURRIED ON BOARD.
4. GARFUNKELS ARE LEAVING JULY 5 ON EL AL...

The way it really happened~The tense, action~packed story of the raid that startled the world.

"OPERATION THUNDERBOLT"



An INTER-OCEAN film "OPERATION THUNDERBOLT" A Golan-Globus Film of a G.S. Films Production

starring YEHOAM GAON ASSAF DAYAN-KLAUS KINSKY SYDIL DANNING-ORI LEVY-ARIK LAVI and MARK HEATH as Idi Amin

Produced by MENAHEM GOLAN and YORAM GLOBUS Music by DOV SELTZER Directed by MENAHEM GOLAN Eastmancolour " Distributed by EMI Films Limited

EMI

FROM THURS.
OCT. 20th

ABC1

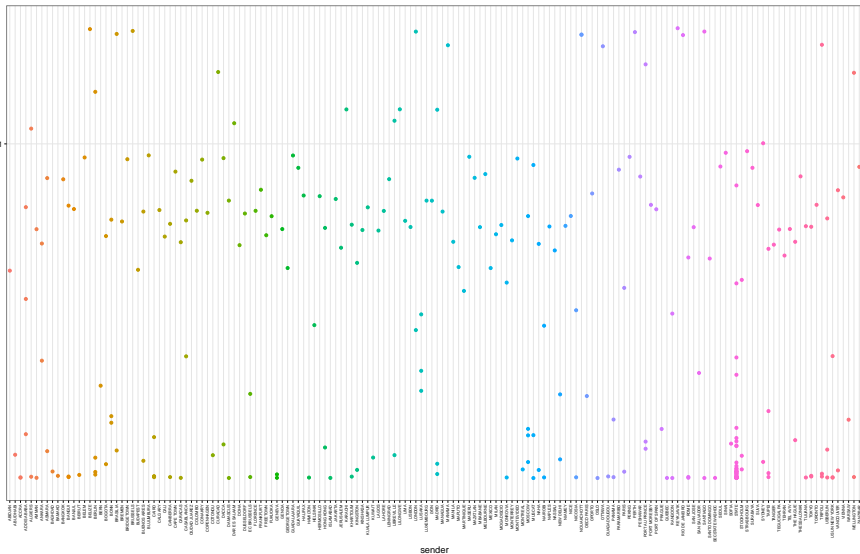
Shaftesbury Ave

Tel: 836 8861

Licensed Bar

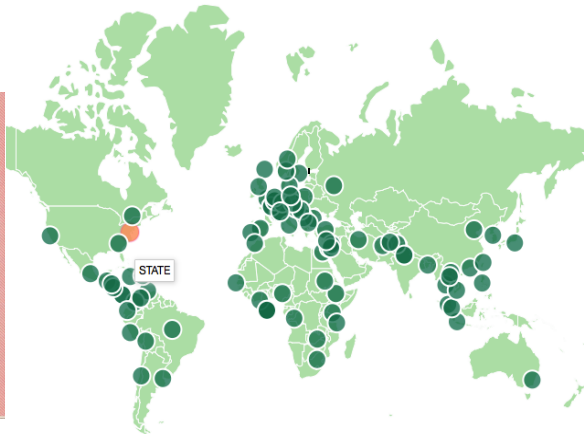
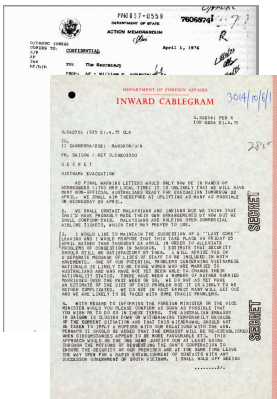
AND AT SELECTED **ABC** AND OTHER LEADING
CINEMAS FROM OCT. 27th. SEE LOCAL PRESS FOR DETAILS

top 5% of content about week 1976-12-16



1. THE GENERAL ACCOUNTING OFFICE IS ACTING UPON A REQUEST, BY THE HOUSE GOVERNMENT OPERATIONS COMMITTEE, TO GATHER INFORMATION ON THE OPERATION OF GAMING DEVICES IN ALL EMBASSY AND CONSULATE FACILITIES. GAMING DEVICES INCLUDE SLOT MACHINES, WHEELS OF CHANCE, DICE GAMES, ETC. BUT SPECIFICALLY EXCLUDE BINGO.
2. WE NEED TO HAVE FOLLOWING INFORMATION BY PRIORITY CABLE NO LATER THAN DECEMBER 30. REPLIES SHOULD BE DIRECTED TO LEAMON R. HUNT, DEPUTY ASSISTANT SECRETARY FOR OPERATIONS. NEGATIVE RESPONSES ARE REQUIRED.
 - A. NUMBER OF GAMING DEVICES, WHERE LOCATED, AND WHO MANAGES THEM.
 - B. OPERATING POLICIES.
 - C. WHETHER LEASING OR CONCESSIONS ARE PERMITTED, TO INCLUDE IDENTIFYING MACHINES LEASED OR BELONGING TO CONCESSIONAIRES.
 - D. NUMBER OF DEVICES PURCHASED AND PROCUREMENT POLICIES.
 - E. AMOUNT OF PROFITS DERIVED.
3. WE APPRECIATE YOUR PROMPT ATTENTION TO THIS REQUEST. KISSINGER

- ▶ 1976-12-21 | BANJUL | STATE | GAMING DEVICES
THERE ARE NO GAMING DEVICES IN EMBASSY FACILITIES. WYGANT
- ▶ 1976-12-22 | BREMEN | STATE | GAMING DEVICES
NO RPT NO GAMING DEVICES OF ANY TYPE AT AMCONSUL BREMEN. LONGMYER
- ▶ 1976-12-21 | BANGUI | STATE | GAMING DEVICES
THERE ARE NO RPT NO GAMING DEVICES IN EMBASSY FACILITIES. QUAINTON
- ▶ 1976-12-22 | BELIZE | STATE | GAMING DEVICES
NEGATIVE RESPONSE. WALSH
- ▶ 1976-12-21 | ABIDJAN | STATE | GAMING DEVICES
FOR: LEAMON R. HUNT, DEPUTY ASSISTANT SECRETARY FOR OPERATIONS NO
GAMING DEVICES EXIST AT THIS POST. STEARNS
- ▶ 1976-12-21 | ASMARA | STATE | GAMING DEVICES
CONSULATE GENERAL ASMARA NEITHER OPERATES OR OWNS ANY GAMING
DEVICES. WAUCHOPE
- ▶ 1976-12-21 | AMMAN | STATE | GAMING DEVICES
EMBASSY AMMAN SUBMITS NEGATIVE REPORT ON AVAILABILITY OF GAMING
DEVICES WITHIN MISSION. PICKERING



- ▶ Topic models can help us find events
- ▶ Extensions:
 - Network characteristics
 - Better characterize an event for fewer “false positives”
 - Word embeddings
 - Autocorrelated time series

Discussion: Modern Probabilistic Modeling

TOPIC
MODELING

The diagram consists of two nested ellipses. The outer ellipse is white and contains the text 'STATISTICS', 'MACHINE LEARNING', and 'DATA SCIENCE' at the bottom. The inner ellipse is shaded light gray and contains the text 'PROBABILISTIC MODELING' and 'TOPIC MODELING'. An arrow points from 'TOPIC MODELING' to a small black dot located within the gray ellipse.

PROBABILISTIC
MODELING

STATISTICS
MACHINE LEARNING
DATA SCIENCE

I. Assume our data come from a model with hidden patterns at work

Topics



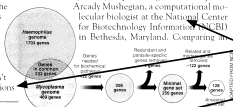
Documents

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,* 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

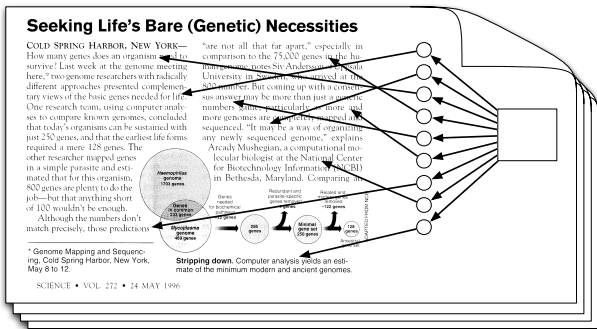
* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 6 to 12.

"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersen, a University of Stockholm researcher at the 800-gene level. 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 Aracly Muehligian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



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

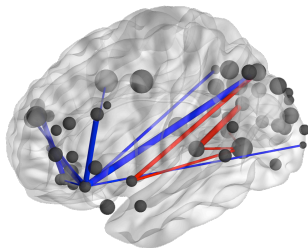
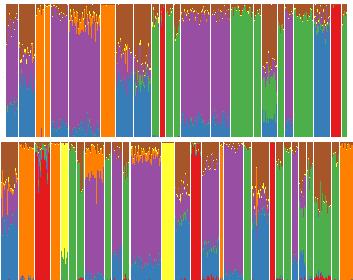
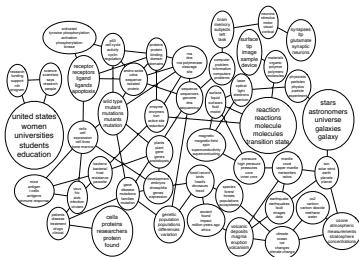
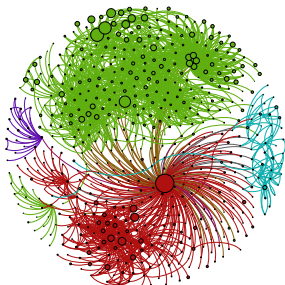
Topic proportions and assignments

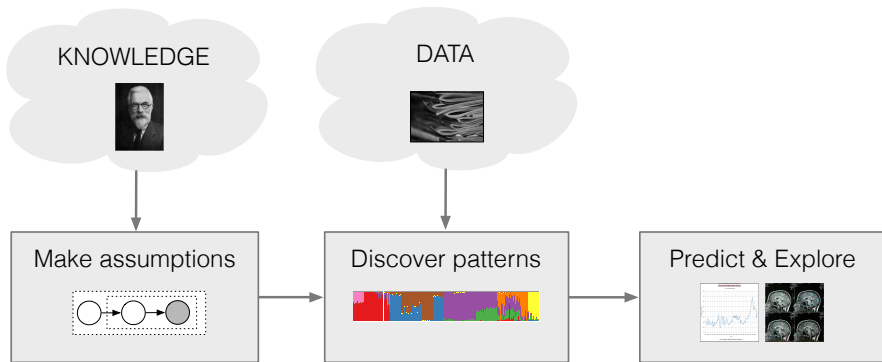


II. Discover those patterns from data

$$\nu^* = \arg \max_{\nu} \mathbb{E}_q [\log p(x, z, \beta \mid \alpha)] + \mathbb{H} [q(z, \beta \mid \nu)]$$

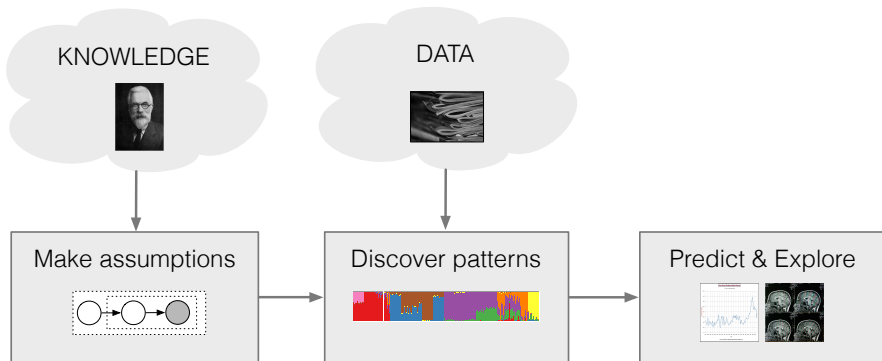
III. Use the discovered patterns to predict about and explore the data





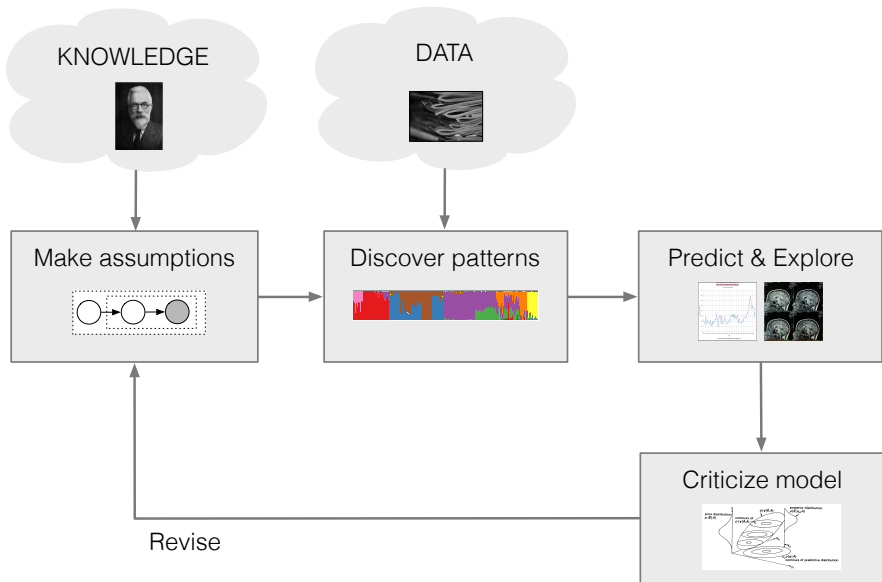
Our perspective:

- ▶ Customized data analysis is important to many fields.
- ▶ This pipeline separates assumptions, computation, application.
- ▶ It facilitates solving data science problems.



What we need in probabilistic ML:

- ▶ **Flexible** and **expressive** components for building models
- ▶ **Scalable** and **generic** inference algorithms
- ▶ **Easy to use** software to stretch probabilistic modeling into new areas





We should seek out unfamiliar summaries of observational material, and establish their useful properties... And still more novelty can come from finding, and evading, still deeper lying constraints.

(John Tukey, *The Future of Data Analysis*, 1962)