

The problem



Science



Politics



Science



Sports



Sports



Science
Politics



Sports
Science



Sports



Politics



Science



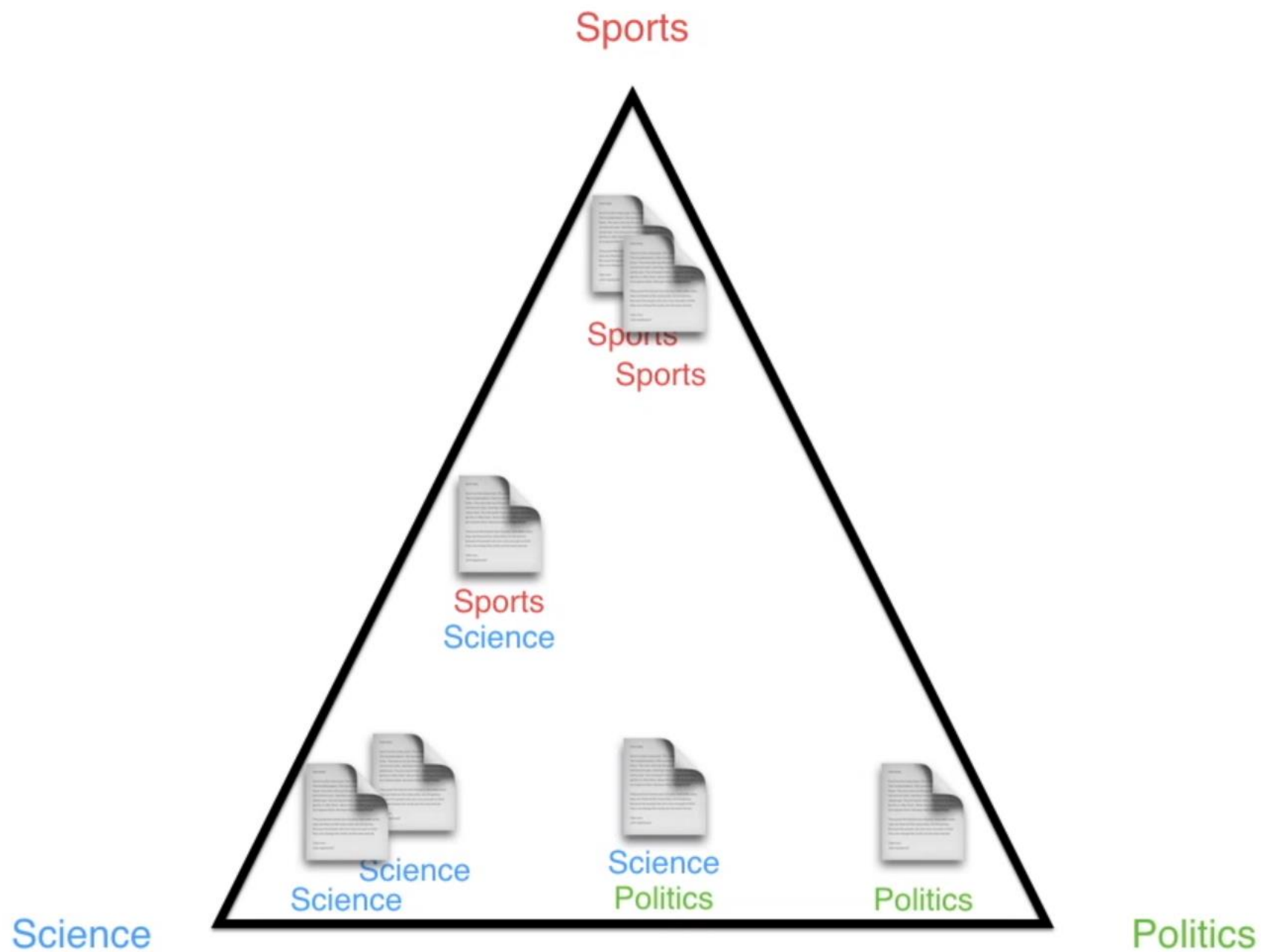
Science

Science

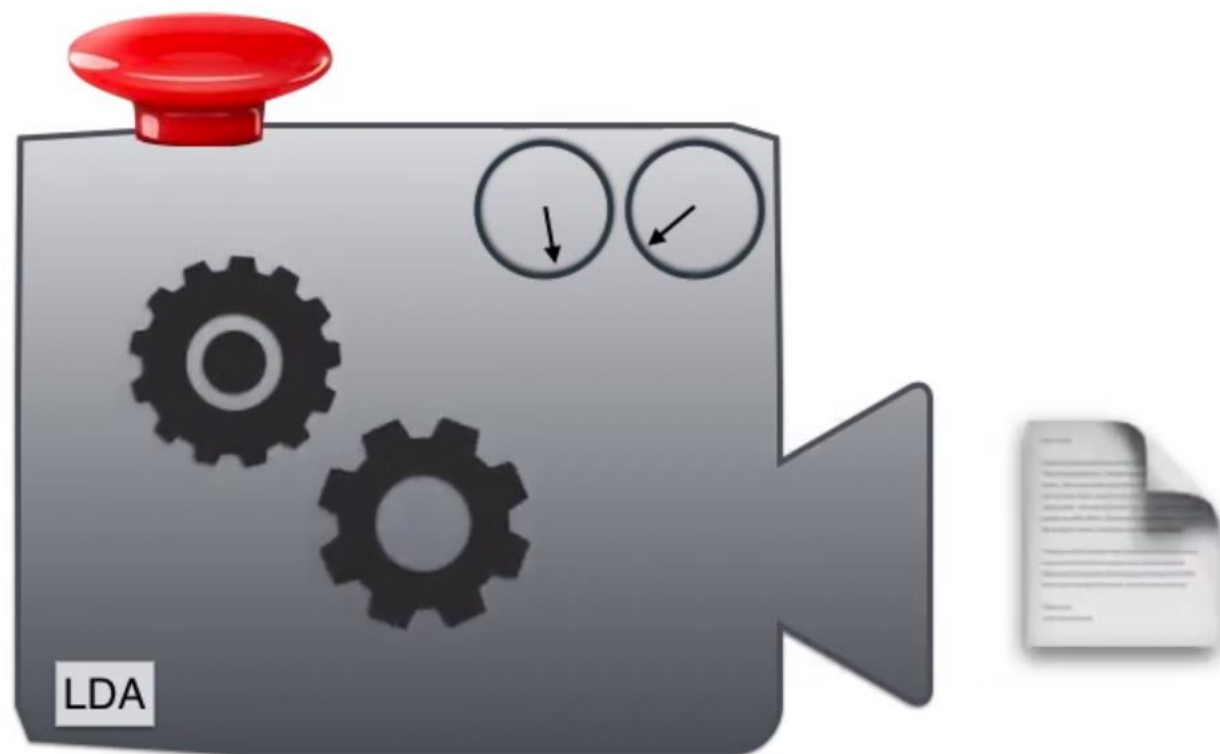
Politics

Sports

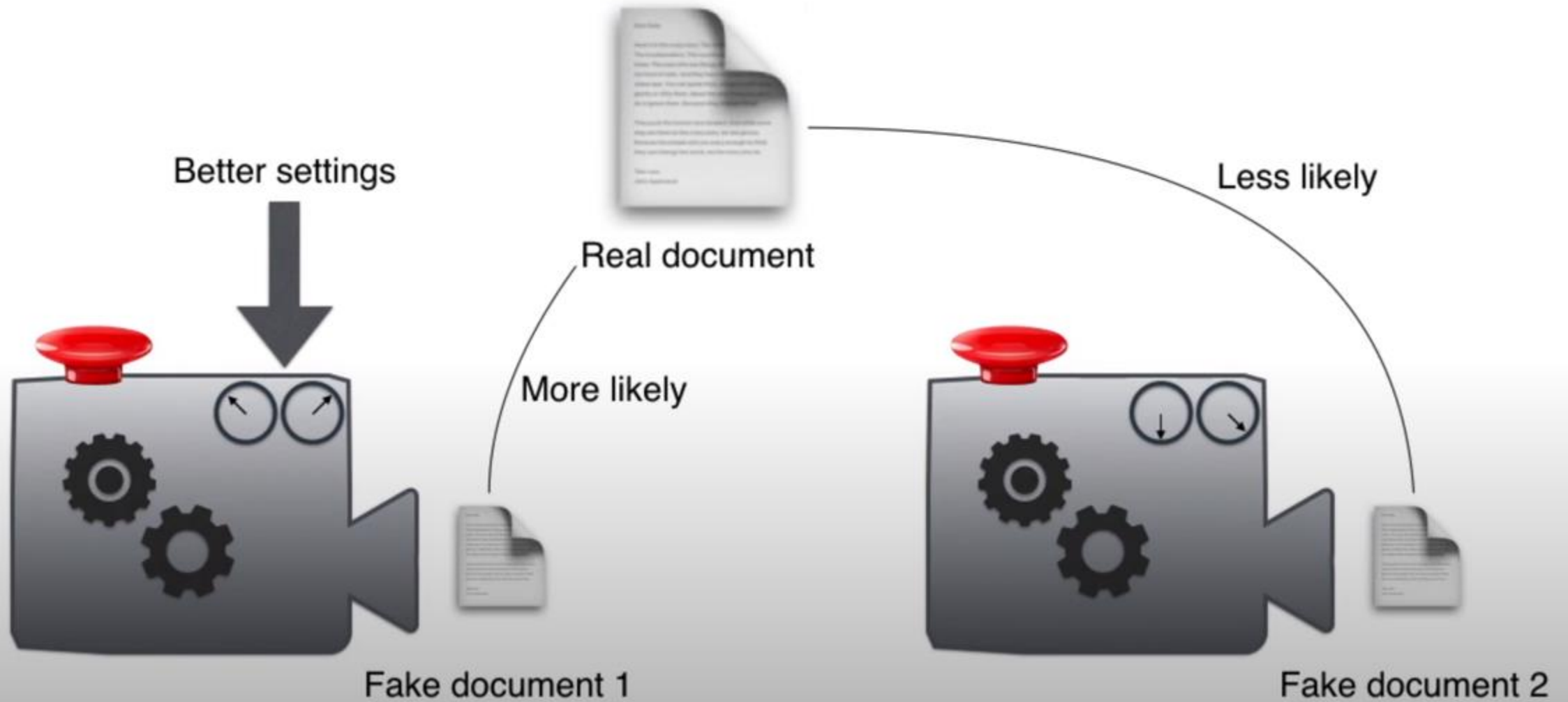
LDA



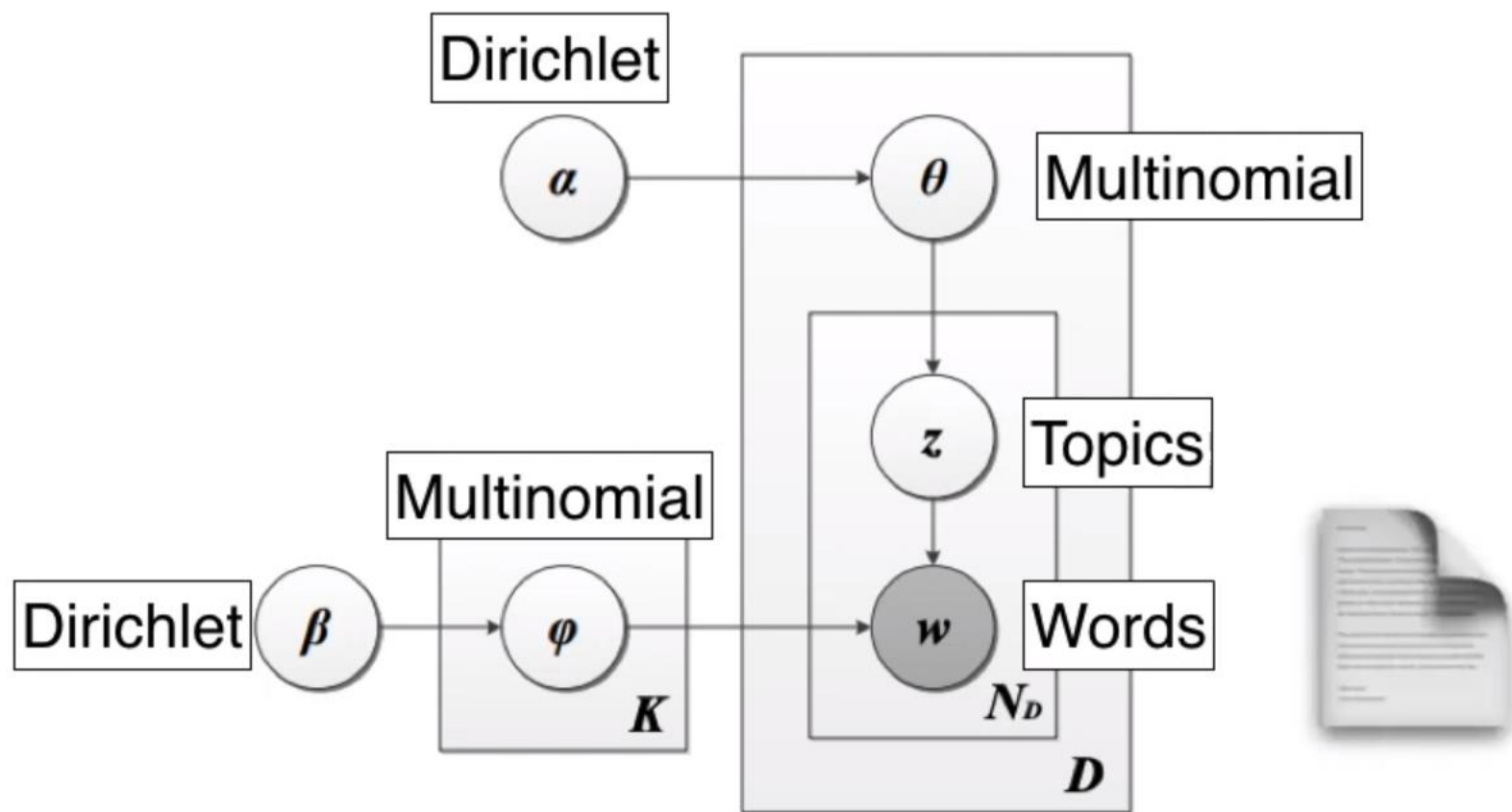
Machine that generates documents



Best settings on the machine

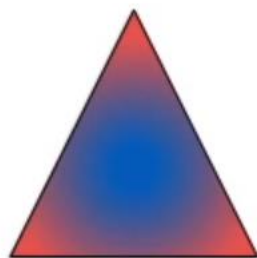


Blueprint for the LDA machine



Probability of a document

$$P(\mathbf{W}, \mathbf{Z}, \boldsymbol{\theta}, \boldsymbol{\varphi}; \alpha, \beta) = \prod_{j=1}^M P(\theta_j; \alpha) \prod_{i=1}^K P(\varphi_i; \beta) \prod_{t=1}^N P(Z_{j,t} | \theta_j) P(W_{j,t} | \varphi_{Z_{j,t}})$$



Topics



Words



Dirichlet
Distributions



Topics

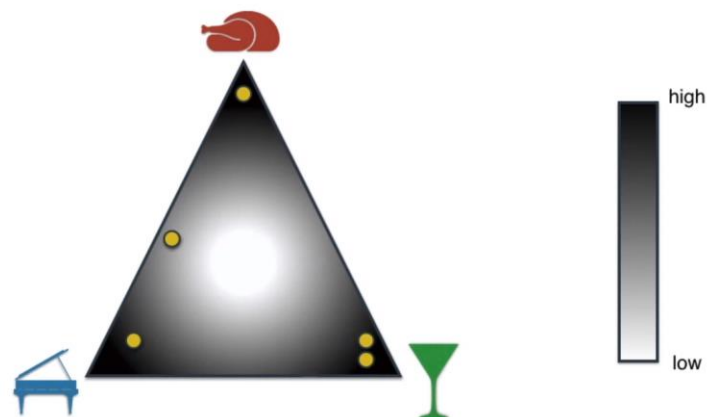


Words

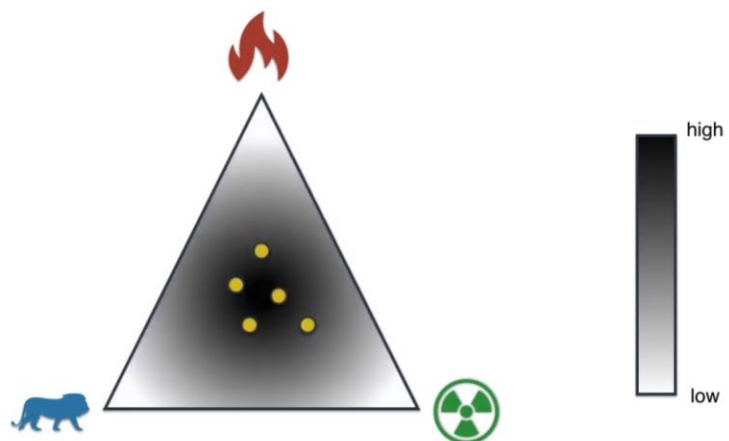


Multinomial
Distributions

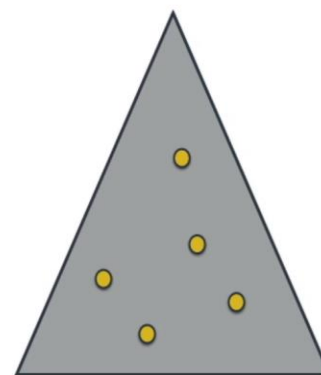
At a party



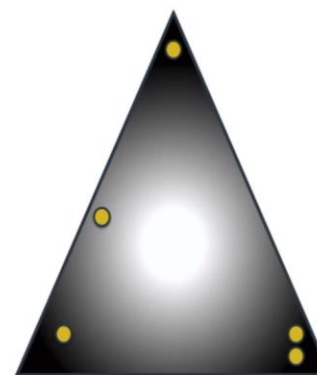
At a party



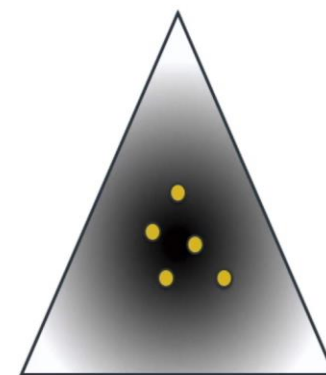
Dirichlet Distributions



$\alpha = 1$



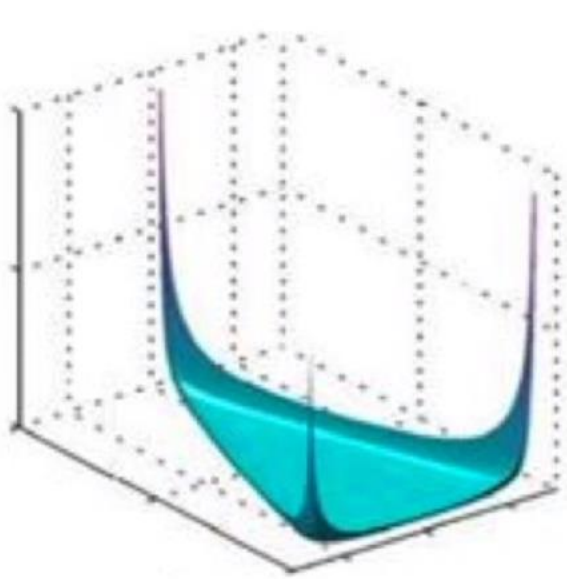
$\alpha < 1$



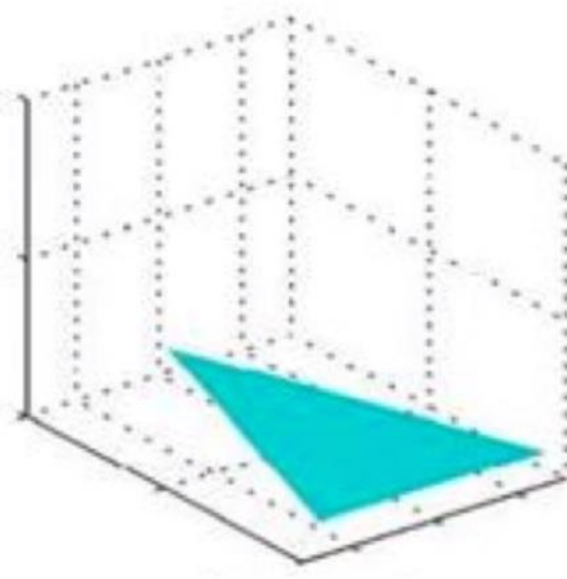
$\alpha > 1$

Dirichlet Distributions

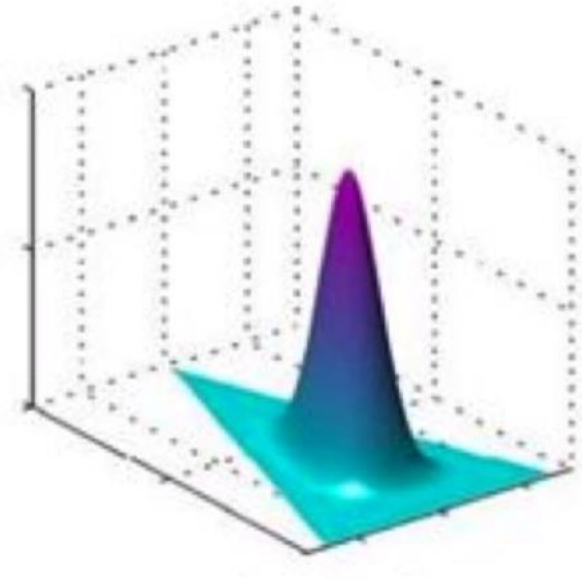
$$f(x_1, \dots, x_K; \alpha_1, \dots, \alpha_K) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{i=1}^K x_i^{\alpha_i - 1}$$



0.7, 0.7, 0.7

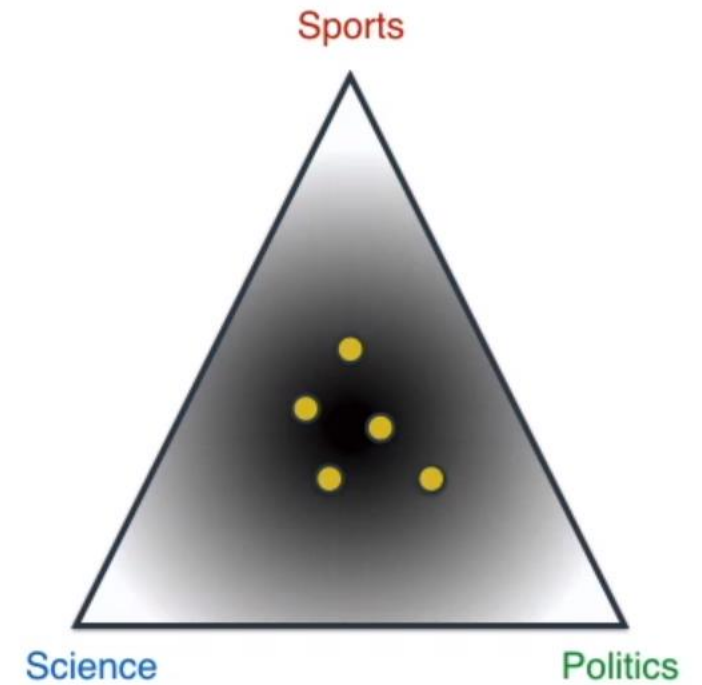
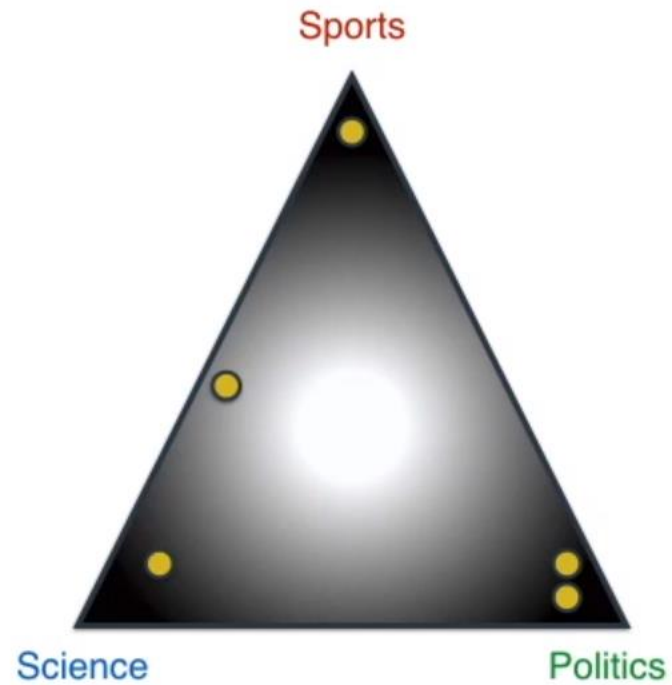
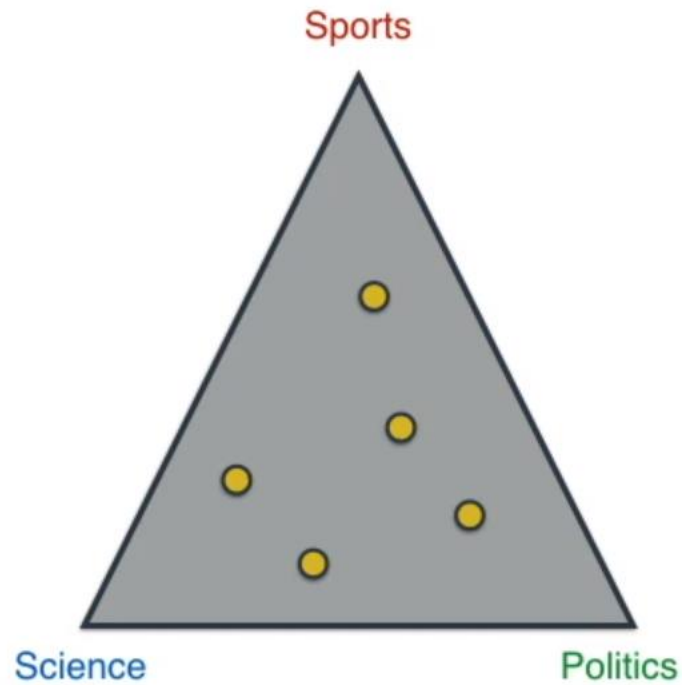


1, 1, 1

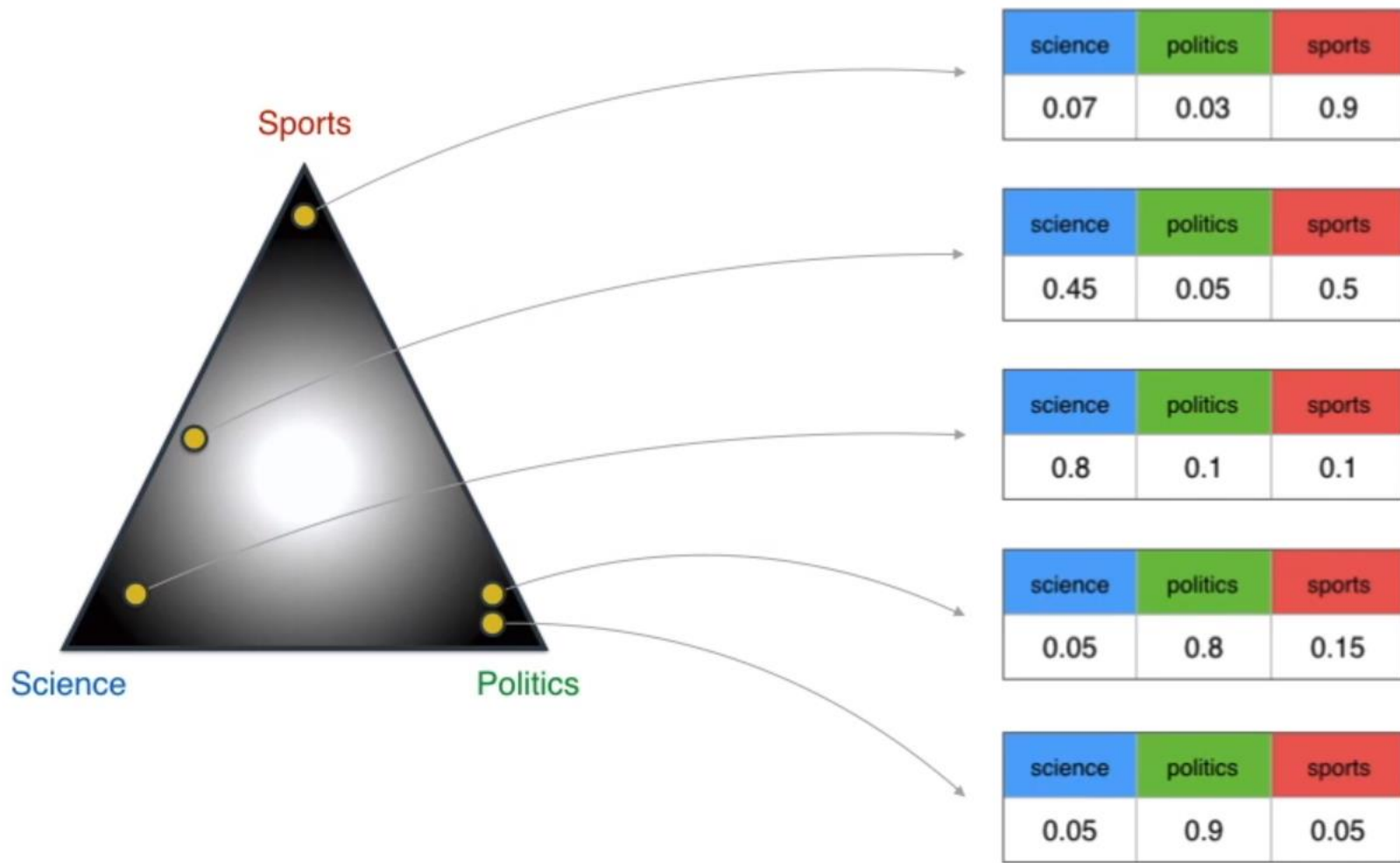


5, 5, 5

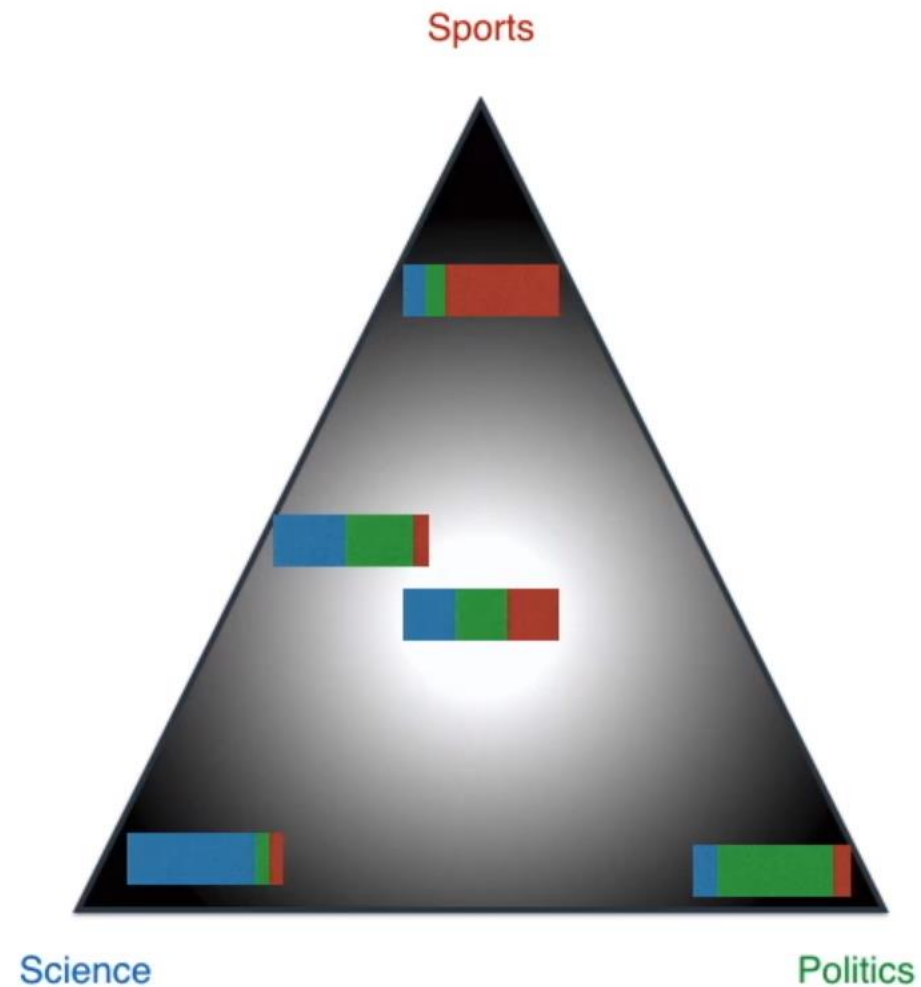
Quiz: Which one for topics?



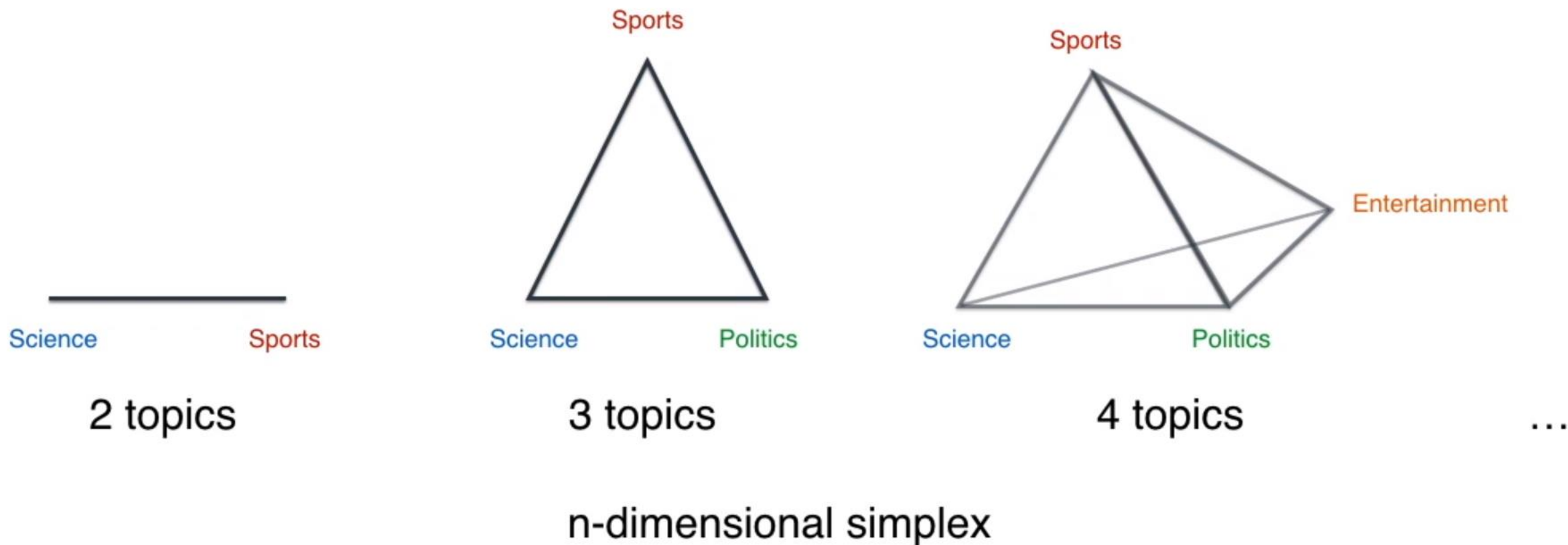
Quiz: Which one for topics?



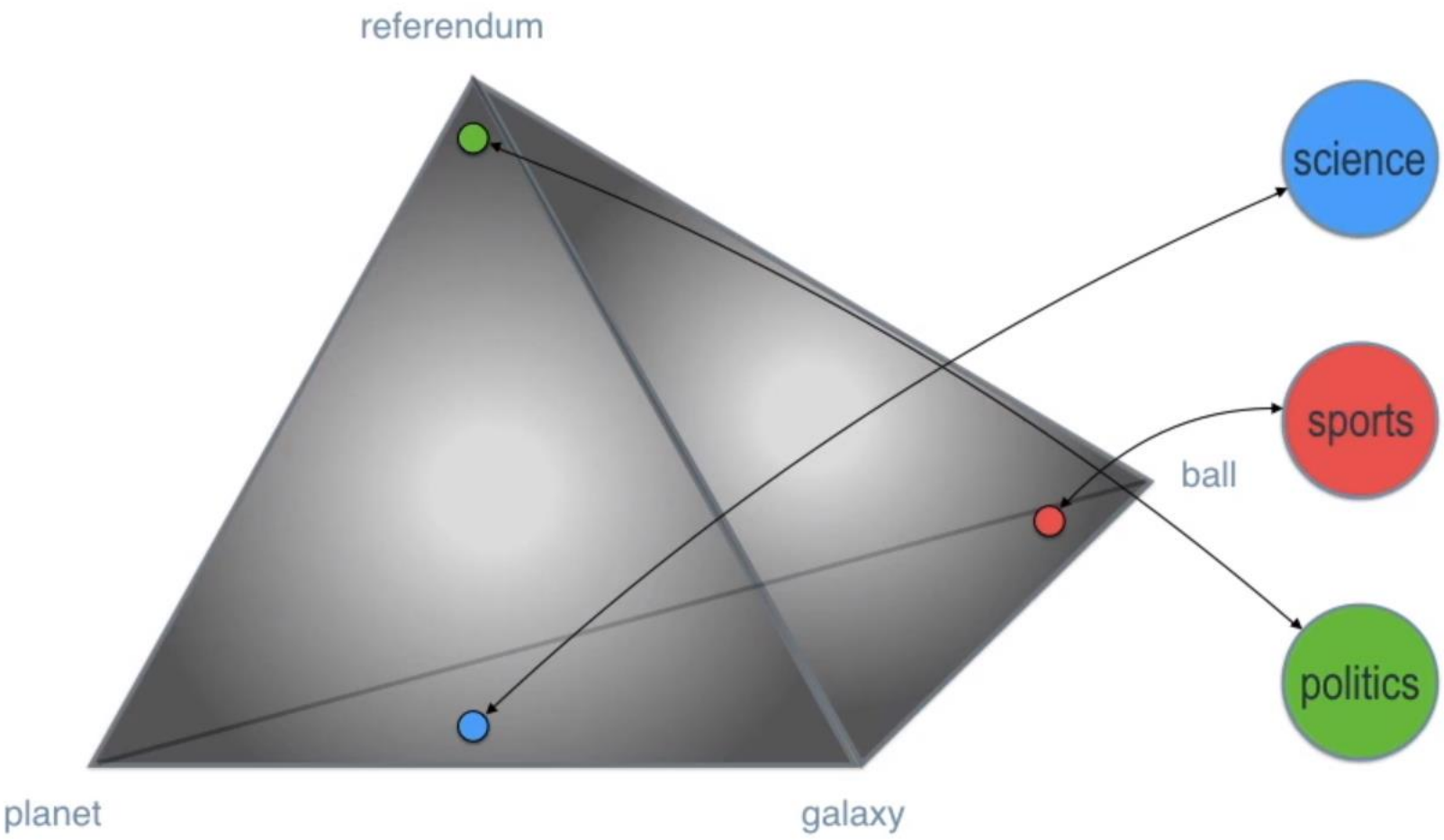
A distribution of distributions



More topics? More dimensions



Quiz: Where to put the topics?

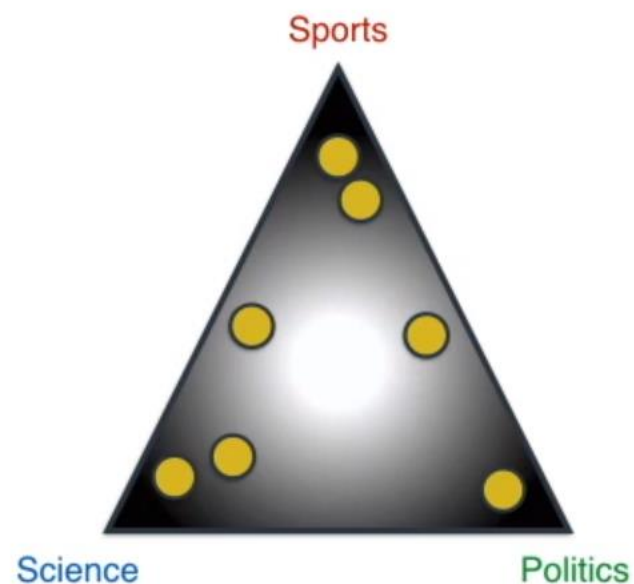


Galaxy	Planet	Ball	Referendum
0.4	0.4	0.1	0.1

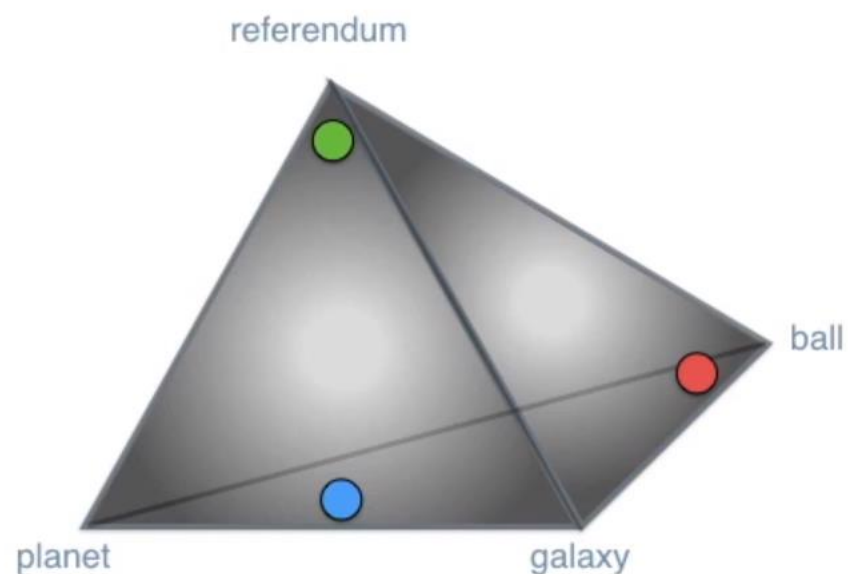
Galaxy	Planet	Ball	Referendum
0.3	0.1	0.5	0.1

Galaxy	Planet	Ball	Referendum
0.1	0.1	0.1	0.7

Two Dirichlet distributions

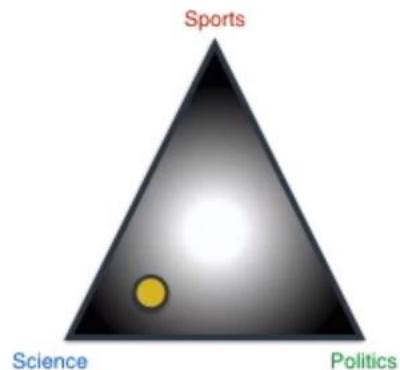


Documents-Topics



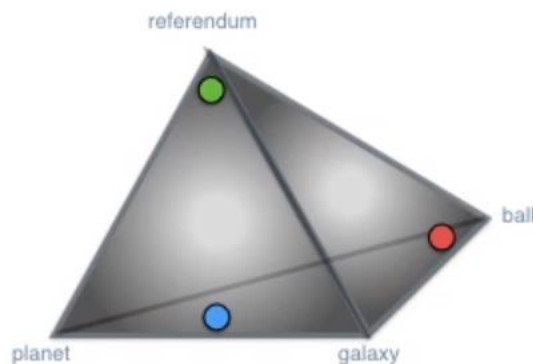
Topics-Words

$$\prod_{j=1}^M P(\theta_j; \alpha)$$



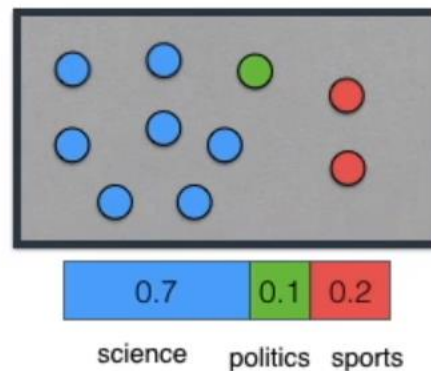
science	politics	sports
0.7	0.1	0.2

$$\prod_{i=1}^K P(\varphi_i; \beta)$$




Galaxy	Planet	Ball	Referendum
0.4	0.4	0.1	0.1
Galaxy	Planet	Ball	Referendum
0.1	0.1	0.1	0.7
Galaxy	Planet	Ball	Referendum
0.3	0.1	0.5	0.1

$$\prod_{t=1}^N P(Z_{j,t} | \theta_j)$$



$$P(W_{j,t} | \varphi_{Z_{j,t}})$$



galaxy	galaxy	planet
galaxy	planet	ball
galaxy	planet	planet
		referendum
planet	ball	referendum
referendum		referendum
galaxy	referendum	
referendum		referendum
galaxy	ball	ball
planet	referendum	galaxy

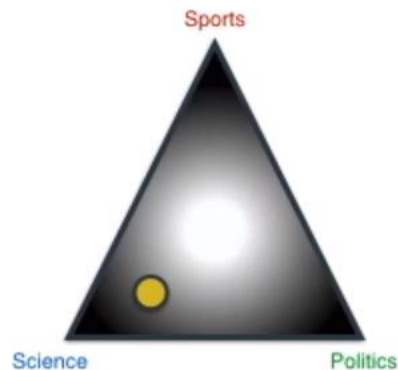
Topics

science
science
sports
science
science
politics
sports
sports
science

Words

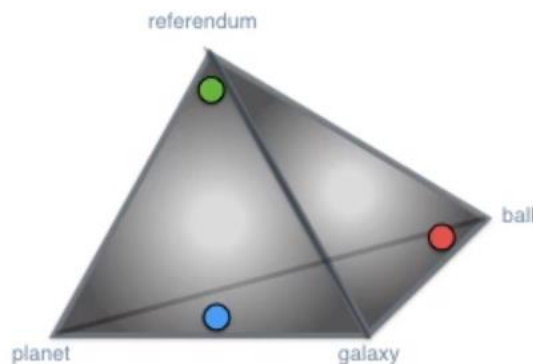
planet
galaxy
ball
planet
galaxy
referendum
galaxy
ball
referendum

$$\prod_{j=1}^M P(\theta_j; \alpha)$$



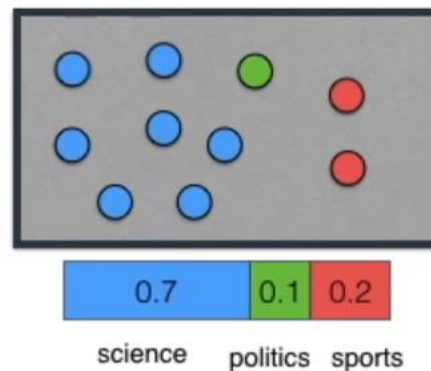
science	politics	sports
0.7	0.1	0.2

$$\prod_{i=1}^K P(\varphi_i; \beta)$$



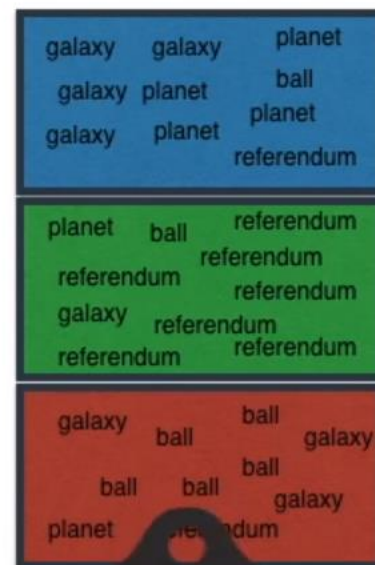
Galaxy	Planet	Ball	Referendum
0.4	0.4	0.1	0.1
Galaxy	Planet	Ball	Referendum
0.1	0.1	0.1	0.7
Galaxy	Planet	Ball	Referendum
0.3	0.1	0.5	0.1

$$\prod_{t=1}^N P(Z_{j,t} | \theta_j)$$



Topics
 science
 science
 sports
 science
 science
 politics
 sports
 sports
 science

$$P(W_{j,t} | \varphi_{Z_{j,t}})$$

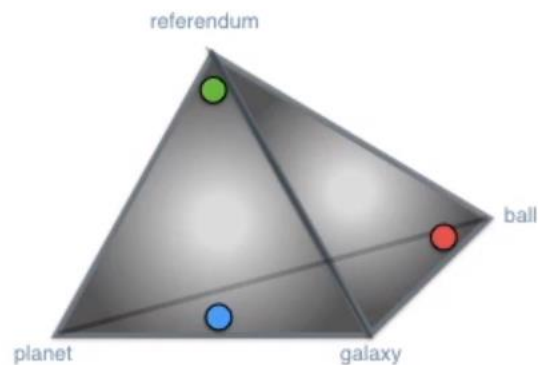
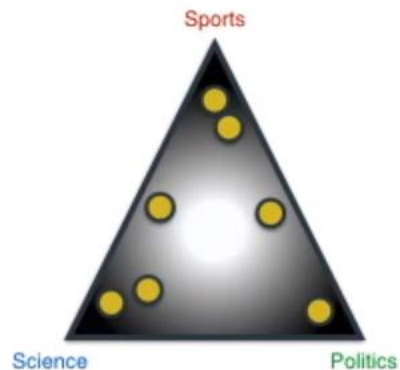


$$\prod_{j=1}^M P(\theta_j; \alpha)$$

$$\prod_{i=1}^K P(\varphi_i; \beta)$$

$$\prod_{t=1}^N P(Z_{j,t} \mid \theta_j)$$

$$P(W_{j,t} \mid \varphi_{Z_{j,t}})$$



planet
galaxy
ball
planet
galaxy
referendum
galaxy
ball

ball
planet
galaxy
galaxy
ball
ball
referendum
ball

referendum
referendum
referendum
referendum
planet
referendum
galaxy

referendum
referendum
planet
galaxy
referendum
planet
galaxy
ball

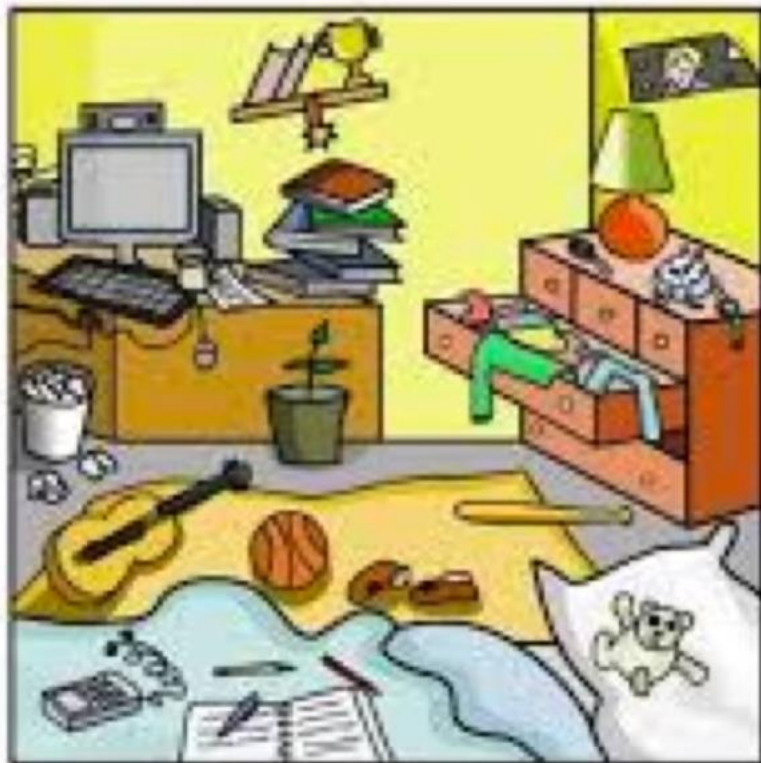
galaxy
planet
galaxy
ball
ball
planet
planet
galaxy

galaxy
planet
referendum
referendum
planet
ball
galaxy
planet

referendum
galaxy
ball
ball
galaxy
referendum
planet
galaxy

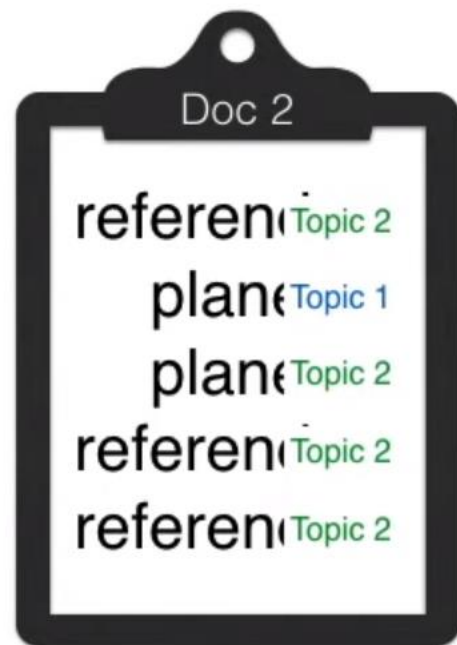


Gibbs Sampling



Gibbs sampling





Topic 1

Topic 2

Topic 3



80% Topic 3
20% Topic 1

Topic 1



80% Topic 2
20% Topic 1

Topic 2



80% Topic 1
20% Topic 3

Topic 3



60% Topic 1
20% Topic 2
20% Topic 3



Property 1:
Articles are as monochromatic as possible

Property 2: Words are as monochromatic as possible

Words

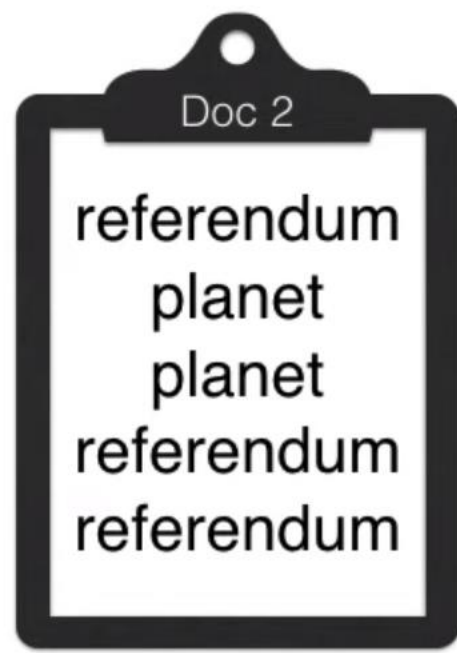


planet
planet
planet
planet
planet
planet
planet
planet

referendum
referendum
referendum
referendum

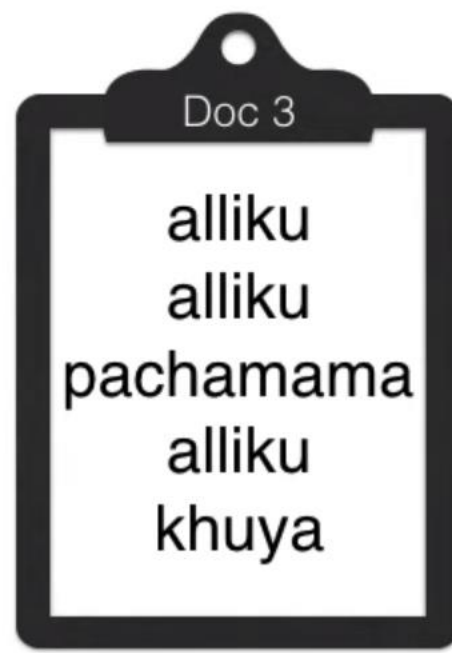
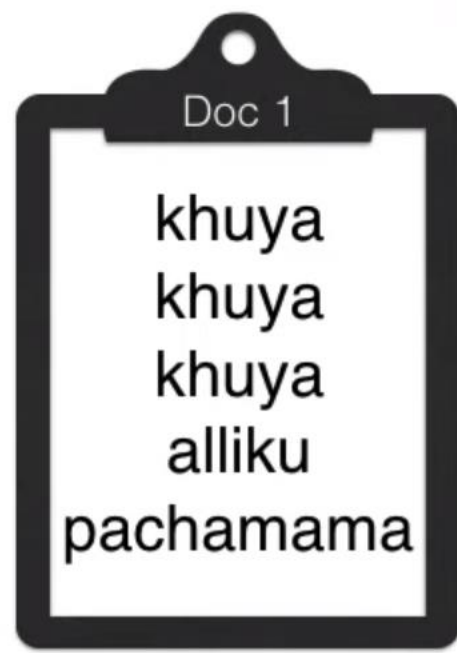
ball
ball
ball
ball
ball

galaxy
galaxy
galaxy



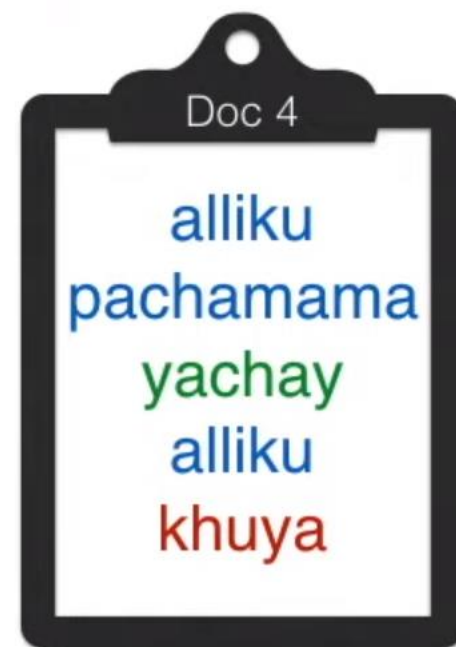
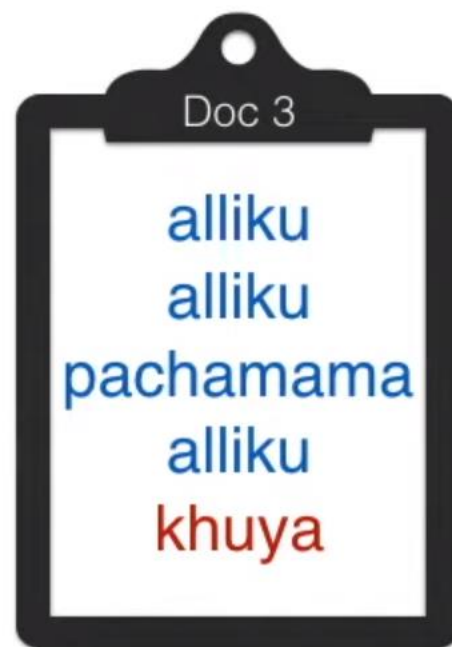
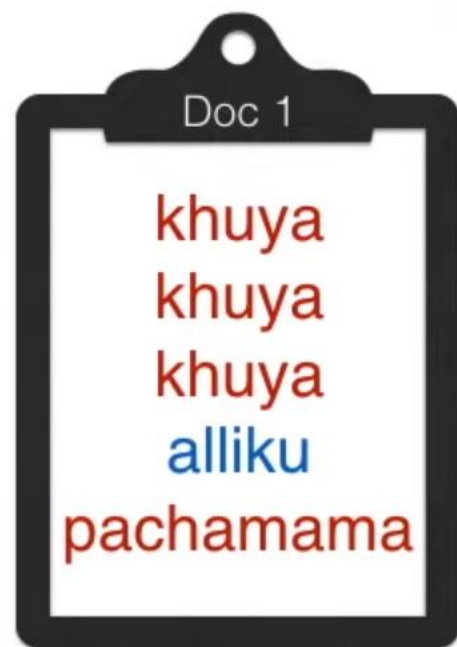
Goal: Color each word with **blue**, **green**, **red**

1. Each article is as monochromatic as possible
2. Each word is as monochromatic as possible



Goal: Color each word with **blue**, **green**, **red**

1. Each article is as monochromatic as possible
2. Each word is as monochromatic as possible



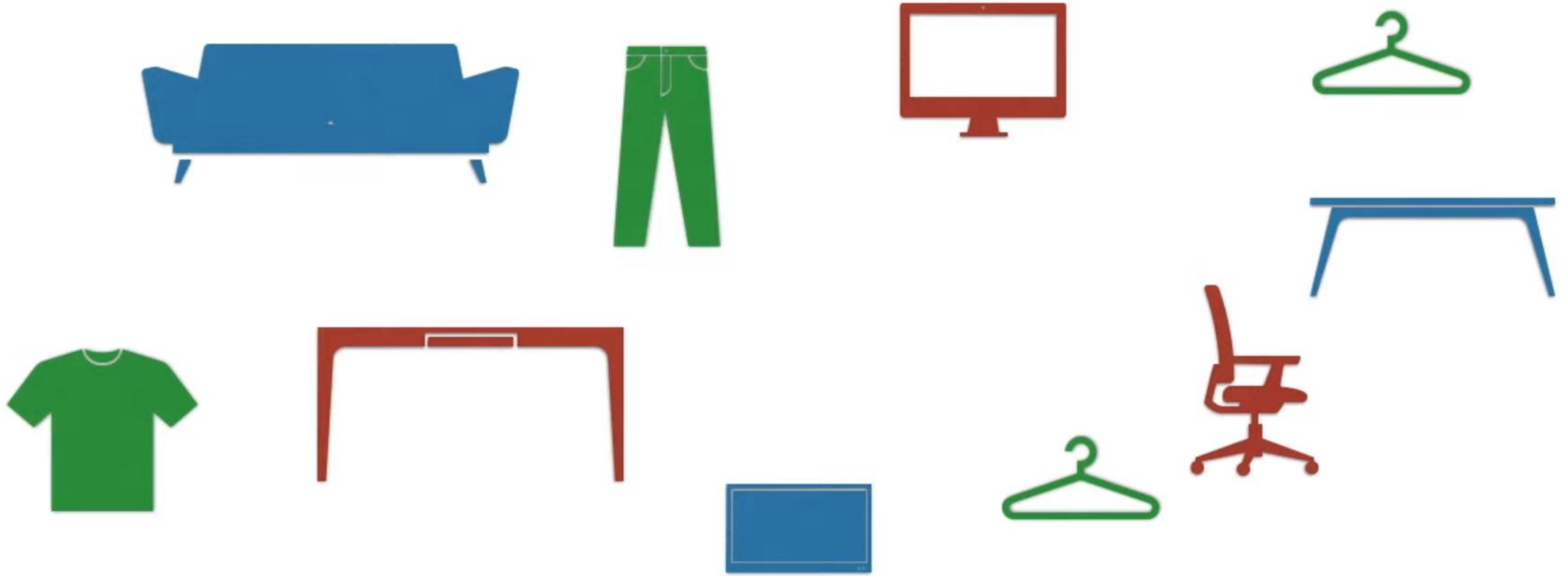
alliku
alliku
alliku
alliku
alliku
alliku
alliku

yachay
yachay
yachay
yachay

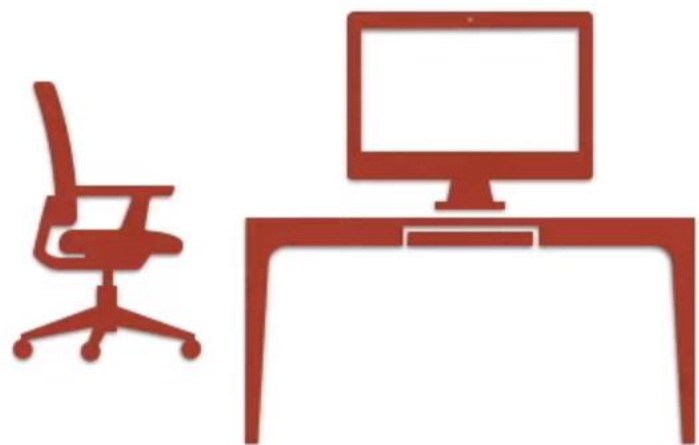
khuya
khuya
khuya
khuya
khuya

pachamama
pachamama
pachamama

Gibbs sampling



Gibbs sampling





Goal: Color each word with **blue**, **green**, **red**

1. Each article is as monochromatic as possible
2. Each word is as monochromatic as possible



Topic 1

How much is Topic 1 in Doc 1?

$$2 + \alpha$$

How much is 'ball' in Topic 1?

$$0 + \beta$$



Topic 2

How much is Topic 2 in Doc 1?

$$0 + \alpha$$

How much is 'ball' in Topic 2?

$$1 + \beta$$



Topic 3

How much is Topic 3 in Doc 1?

$$2 + \alpha$$

How much is 'ball' in Topic 3?

$$3 + \beta$$



80% Topic 3
20% Topic 1

Science

Topic 1
planet (7)
galaxy (2)



80% Topic 2
20% Topic 1

Politics

Topic 2
referendum (4)
planet (1)



80% Topic 1
20% Topic 3



60% Topic 1
20% Topic 2
20% Topic 3

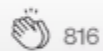
Sports

Topic 3
ball (5)
galaxy (1)

**Shashank
Kapadia**

Data Science
Manager @Monster
Building scalable
and operationalized
ML solutions for
data-driven
products. My articles
on Medium don't
represent my
employer.

Follow



816



22



Eye Balling Models

- Top N words
- Topics / Documents

Intrinsic Evaluation Metrics

- Capturing model semantics
- Topics interpretability

Human Judgements

- What is a topic

Extrinsic Evaluation Metrics/Evaluation at task

- Is model good at performing predefined tasks, such as classification