Revise (NLP) → Topic modelling (unsupervise)   - Intuitive way

(NLP) → (NLU)  { clustering diff articles
                 topic and them
                 one topic }   - what is
                                happen

process

[words | token]

google news   (NLS)

                        ↳ Global
                          sports    (NLU)    - what is
                                              outcome
text → numeric

⇒ count          Elon ✓
⇒ tfidf          Rupee ✓           Business
⇒ encoding       Google ✓   →      ~~name of topic~~
                 Nifty ✓
DTM

TOC
→ Algebra  (A/M/s/p  Matrix)  } mimic today } machine

→ Trig
→ calculus → dej
              → integ
→ Steels

Deepali are

---

( unigram ), ( bigram ) , trigram , n-gram , feature

DTM    | I | am | I am | is | - toke  indepen — [ I am ]        tf.idf

I am
[I ✗]  | I am | Taufique is | Aman/Deepali | Deepali are | Aman & Deepali are
                                                            studying NLP

→ Taufique is
( ✓ am )

# Usage of Topic Modelling

$\underline{1}$ cells $\longrightarrow$ $\underline{1}$ tissues $\rightarrow$ $\underline{1}$ organ $\longrightarrow$ $\underline{1}$ human body

$\underline{1}$ words $\rightarrow$ 1 sentence $\rightarrow$ 1 para $\rightarrow$ 1 doc

topic modelling

---

Assumption $\rightarrow$ start    (Doc) $\longrightarrow$ (Topic) $\rightarrow$ (term)     (1 topic)

$\underline{DTM}$

(1 doc)

(1 doc $\rightarrow$ 1 topic $\rightarrow$ 1 term) find out

n topic $\rightarrow$ 1 term

topic model

1 doc $\rightarrow$ 1 topic

$\downarrow$

n terms

$doc \longrightarrow$ mixture topic

$topic \longrightarrow$ mixture terms $\Big\}$ topic modelling

$Doc \rightarrow topic \rightarrow words$

assmp

1 do | Sourav Ganguly & Greg Chappell |

Sports , Controversy

$\hookrightarrow SG \perp GC$

SG
GC

Latent Dirichlet Allocation DTM Latent
topic

Doc $\rightarrow$ topic $\rightarrow$ words

$100°C \rightarrow 100°C$

Maths ~~Airic~~ → intuition

topics → edge of shape

doc → (Topic) → words

(GC)

words

balsm

Messi

Foot
contoo

SG

LHS → RHS

doc → mul topi

topic → n words

topics → edge of shape

(LDA) food (T_1) (doc)

DD

(T_2) (T_3)

Song Drink

$T_1$ equally distant

$\ell$

$\ell$

$\ell$

$\ell$

$\ell/2$

$\ell$

$\ell$

$\ell$

$L$

3 length = $\ell$

# The problem

Goal

Science  Politics  Science  Sports  Sports  Science Politics  Sports Science

**Doc 1**
- ball
- ball
- ball
- planet
- galaxy

**Doc 2**
- referendum
- planet
- planet
- referendum
- referendum

**Doc 3**
- planet
- planet
- galaxy
- planet
- ball

**Doc 4**
- planet
- galaxy
- referendum
- planet
- ball

Sports
Science

## LDA

Science  Politics  Science

Sports  Sports  Science Politics

Sports Science

know

Sports

Science  Politics

Sports
Science

Science  Science Science  Science Politics  Politics  Politics

# Best settings on the machine



Better settings

Real document

More likely

Less likely

winning

$1 \Rightarrow$

Fake document 1

Fake document 2

gibbers

$D_1 = 100 \text{ words}$

(LDA) → module→dw?

Perplexity
└→ SS

"$\alpha$"  "$\beta$"

$D \; D$ → relate → sample

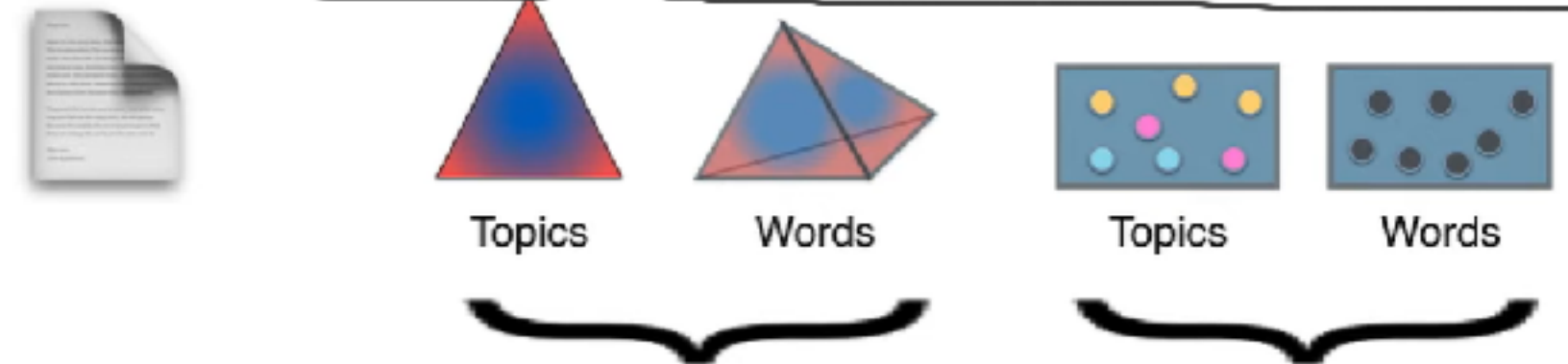## Best settings on the machine

P(same)

Best settings

# Blueprint for the LDA machine



## Probability of a document

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

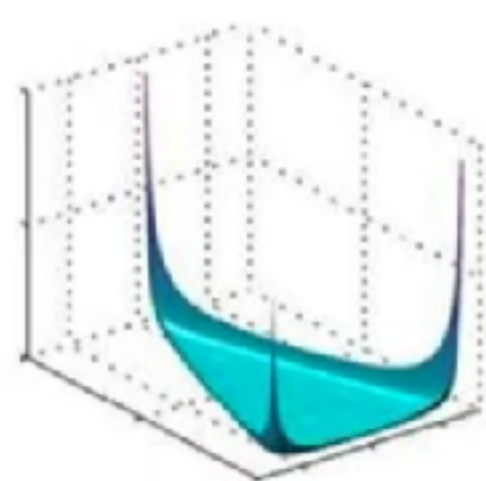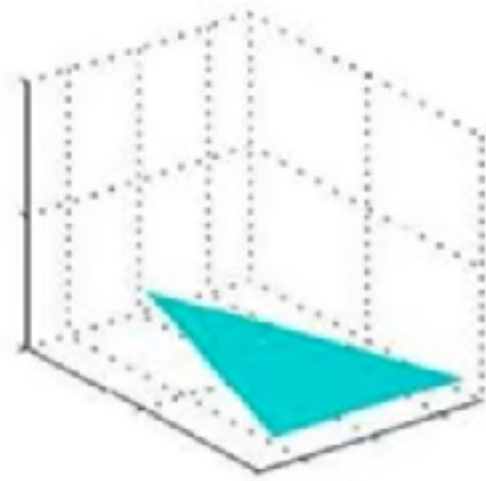Topics    Words    Topics    Words

Steps | Intuition-2

(0, '0.024*"ban" + 0.017*"order" + 0.015*"refugee" + 0.015*"law" + 0.013*"trump"
'+ 0.011*"kill" + 0.011*"country" + 0.010*"attack" + 0.009*"state" + '
'0.009*"immigration"')   topic → R ∪ w          topic 1

(1, '0.020*"student" + 0.020*"work" + 0.019*"great" + 0.017*"learn" + '
'0.017*"school" + 0.015*"talk" + 0.014*"support" + 0.012*"community" + '
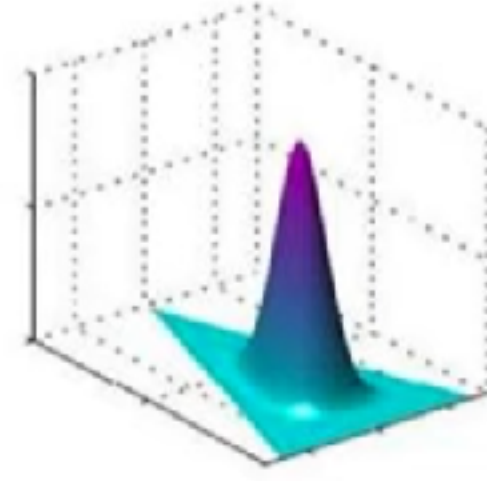'0.010*"share" + 0.009*"event")   topic → ? → Arma

# Dirichlet Distributions

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



| 0.7, 0.7, 0.7 | 1, 1, 1 | 5, 5, 5 |

```
ia.vate {index , ( )}

    for col is col-names

    sum & np-sum (_____,axis

    sum (is), =
```

Eucl (SS) → PD

unsupervise

P → PD

SS

Perplexity

Number of Topics

Smoothed Unigram
Smoothed Mixt. Unigrams
LDA
Fold in pLSI

1/5(elbow method