

Latent Dirichlet Allocation

Satya Pattnaik

I am good

I am great

{ S: 2
am: 2
good: 1
great: 1 }

A₁
A₂
:
:
A_N

2 topics

LDA

T₀ { cats dogs tiger }

T₁ { bat tennis cricket }

Animals

Sports



Normal - μ, σ
Dirichlet - α

Idea

LDA - not a classification algo
— generative algo

$$\left\{ \begin{array}{c} \alpha_1, \alpha_2 \\ \alpha, \beta \end{array} \right\}$$

Two Dirichlet Distributions

Maximum Likelihood Estimation

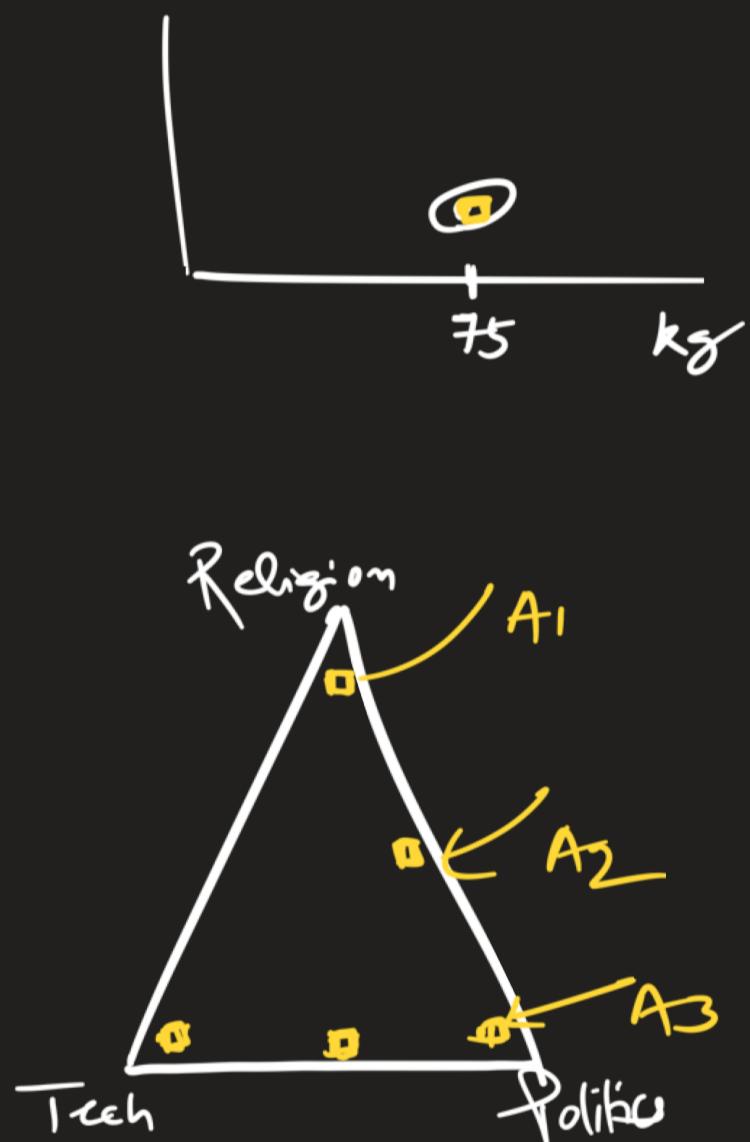
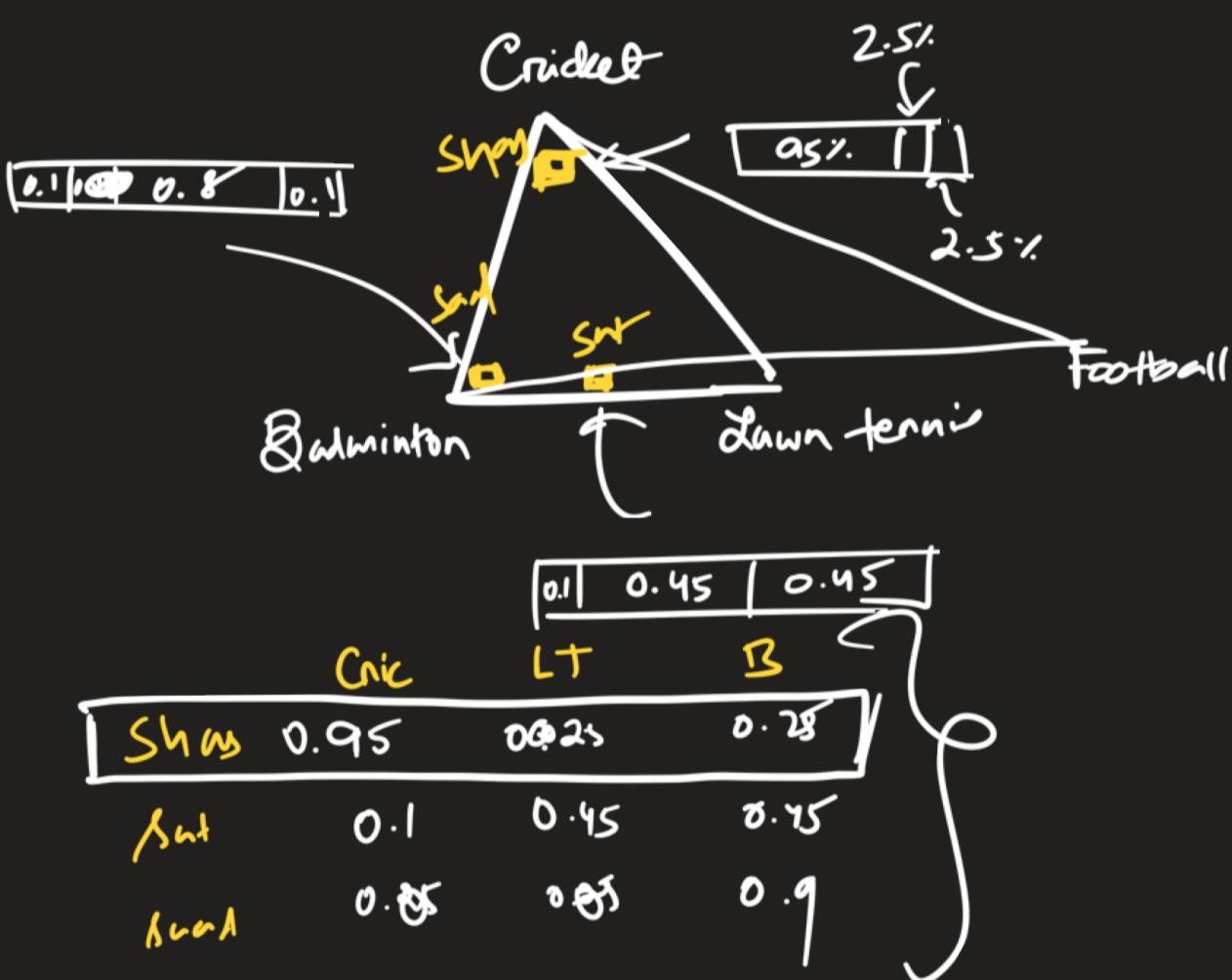
1 5'8"
2 5'9"
3 5'5"
⋮
40

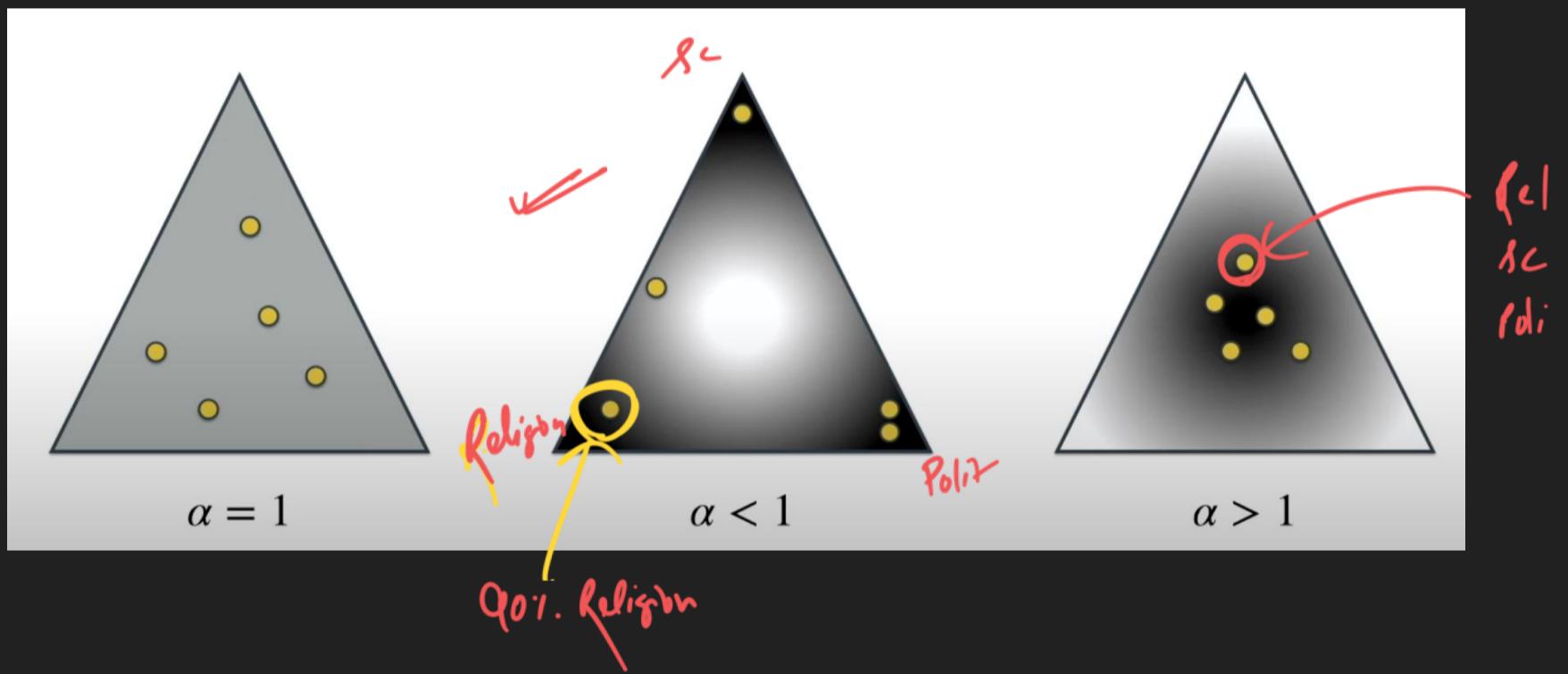
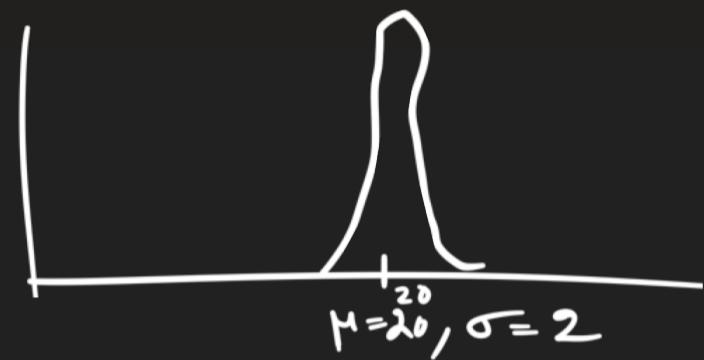
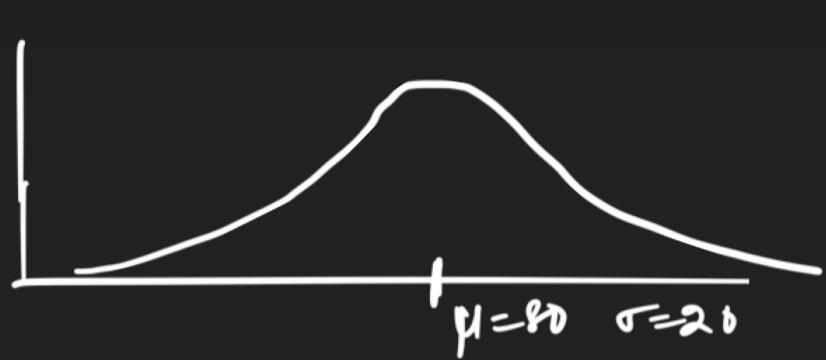
Data $\rightarrow \mu, \sigma$

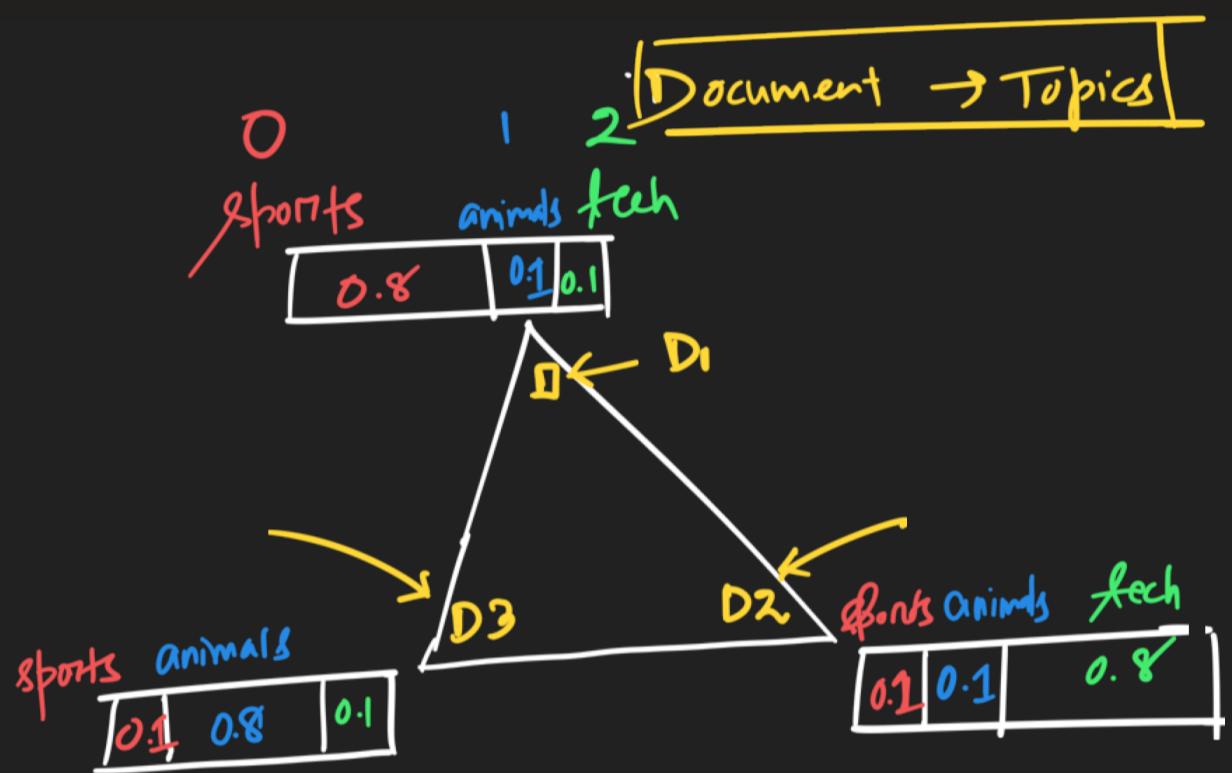
Andreas
Rainer

$$\left\{ \begin{array}{c} \alpha_1, \alpha_2 \\ \alpha, \beta \end{array} \right\}$$

Q: 15

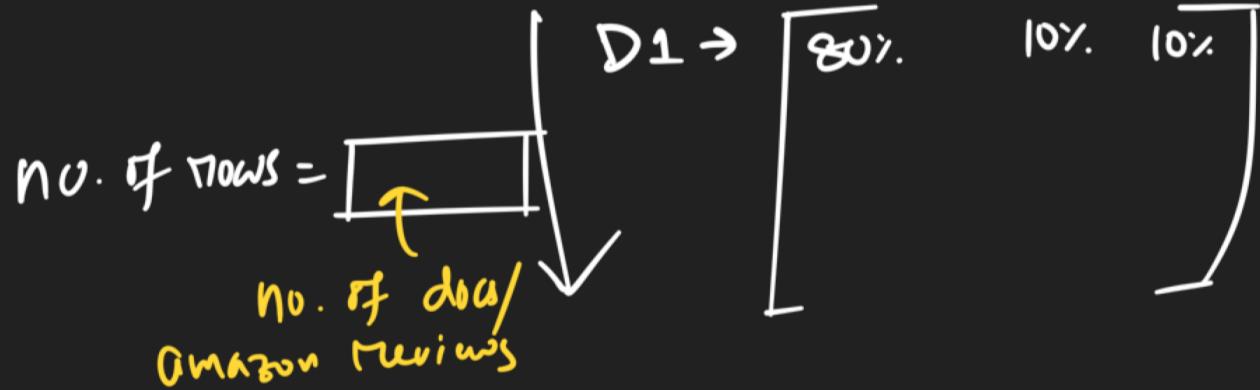




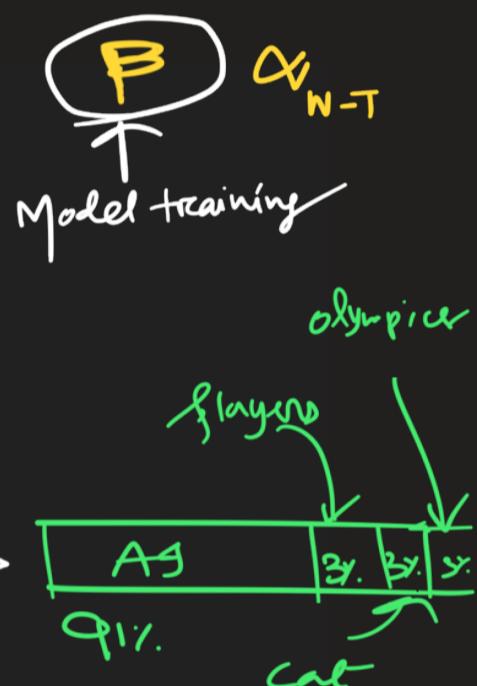
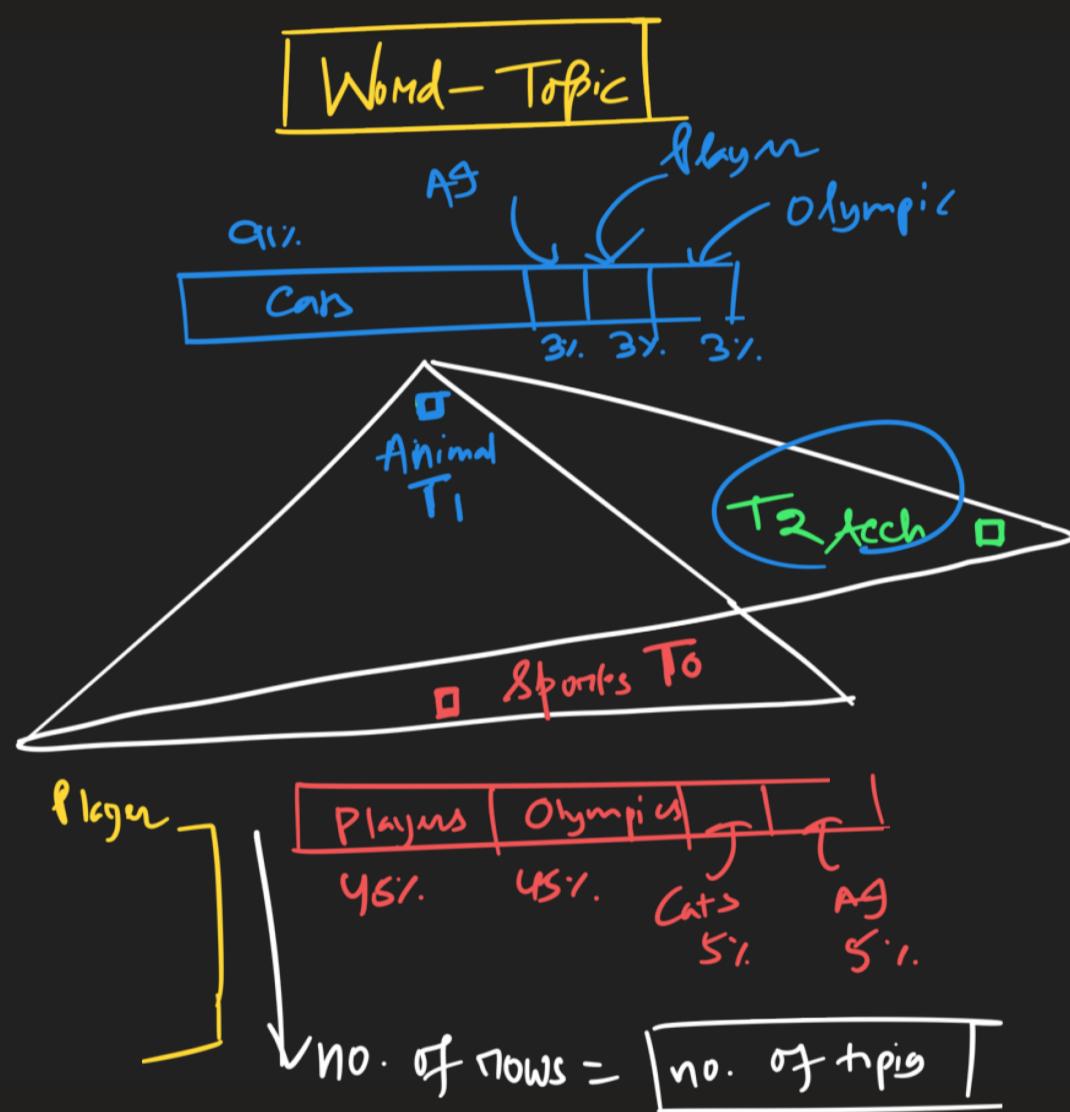


α α_{D-T}

↑
Model Training



$\frac{\text{Vocab} \text{ (1 million)}}{\rightarrow \text{Cats}}$
 $\rightarrow \text{AG}$
 $\rightarrow \text{Olympics}$
 $\rightarrow \text{Players}$
 $\underline{\text{no. of cols} = \text{Vocab size}}$



Notations

k → no. of topics (user defined)

V → no. of unique words (vocab)

M → no. of docs/reviews

N → no. of words of in doc

M'

| | | |
|-------|-------|--|
| M_1 | T_1 | |
| 60 | 27 | |
| 50 | 10 | |
| 70 | 10 | |
| 65 | 20 | |

$\left\{ \begin{bmatrix} 80\text{kg}, 77\text{kg}, 67\text{kg}, \dots \end{bmatrix} \right\}$

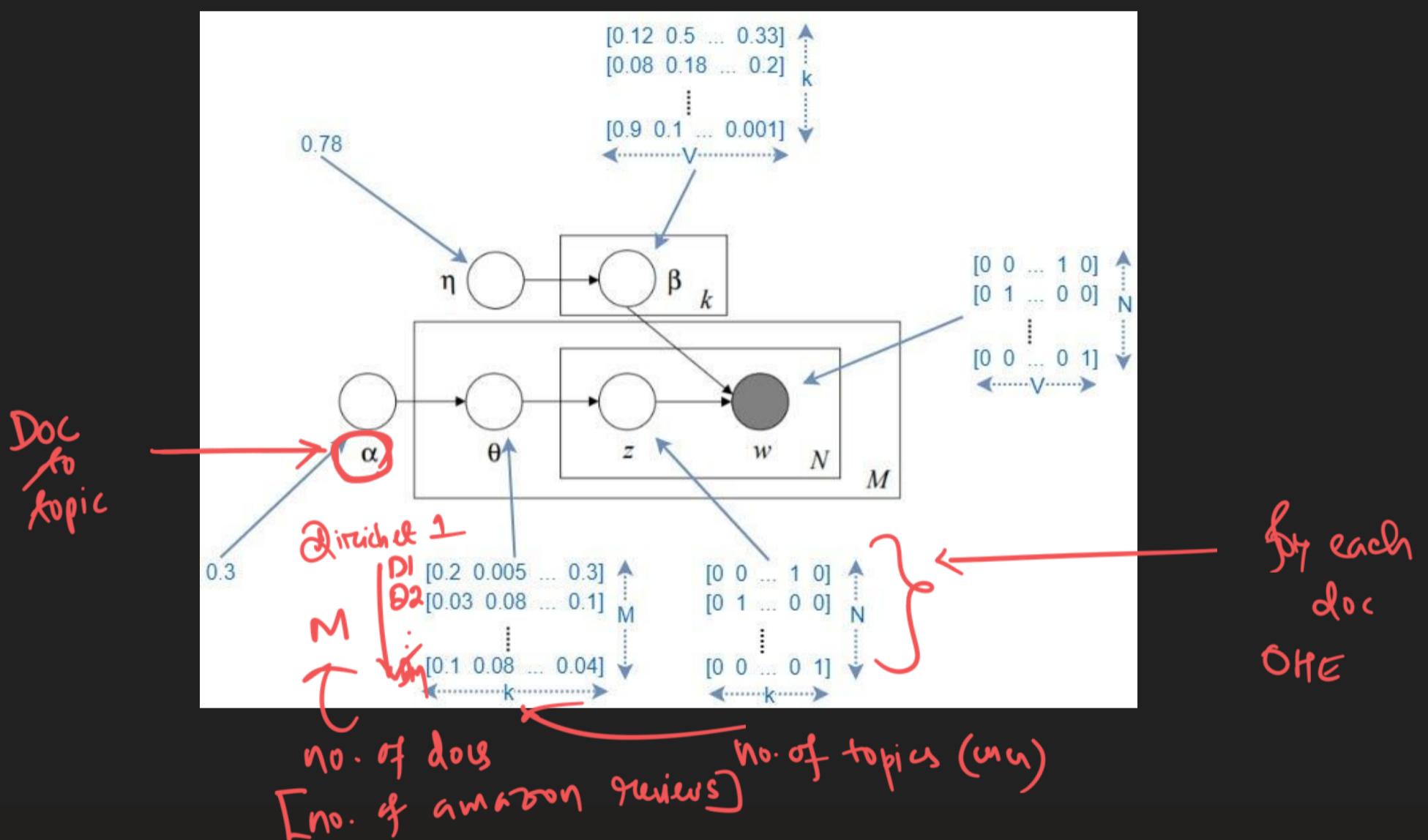
Calculus - $\alpha_{\text{optimal}}, \beta_{\text{optimal}}$

$$\boxed{\frac{\partial L}{\partial w} = 0} \quad w_{\text{optimal}}$$

| D-T | W-T |
|---------------|-------------|
| Dirichlet D+1 | Dirichlet 2 |
| α_1 | β_1 |
| α_2 | β_2 |
| α_3 | β_3 |
| . | . |
| α_n | β_n |

Generated docs
 Doc Set₁ / All Reviews₁
 Doc Set₂ $\xleftrightarrow{\text{most similar}}$ All docs / All Review₂
 Doc Set₃ $\xleftrightarrow{\text{most similar}}$ All docs / All Review₃
 .
 .
 .
 Doc Set_n

Doc-set [all the docs in main doc]



$N=5$
 \Rightarrow each doc
 is of
 5 words

| | | | |
|----|---|---|---|
| W1 | 1 | 0 | 0 |
| W2 | 1 | 0 | 0 |
| W3 | 1 | 0 | 0 |
| W4 | 1 | 0 | 0 |
| W5 | 0 | 0 | 1 |

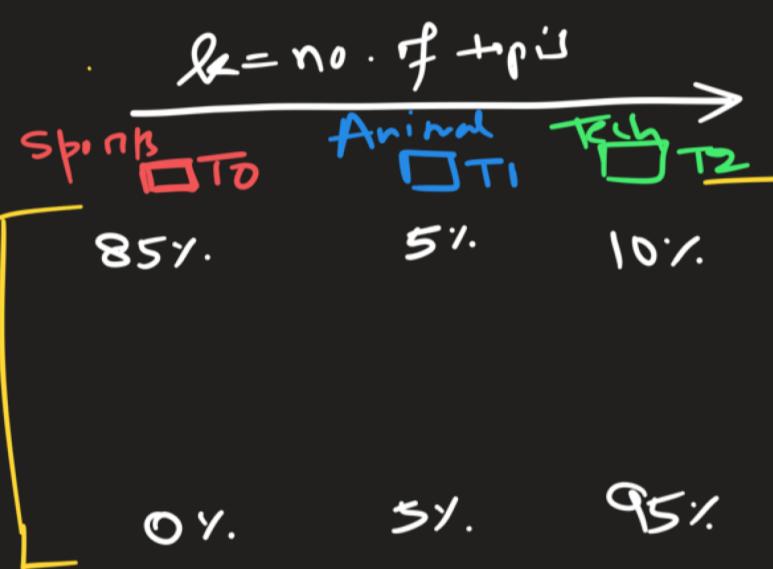


Dirichlet - 1
 Dist

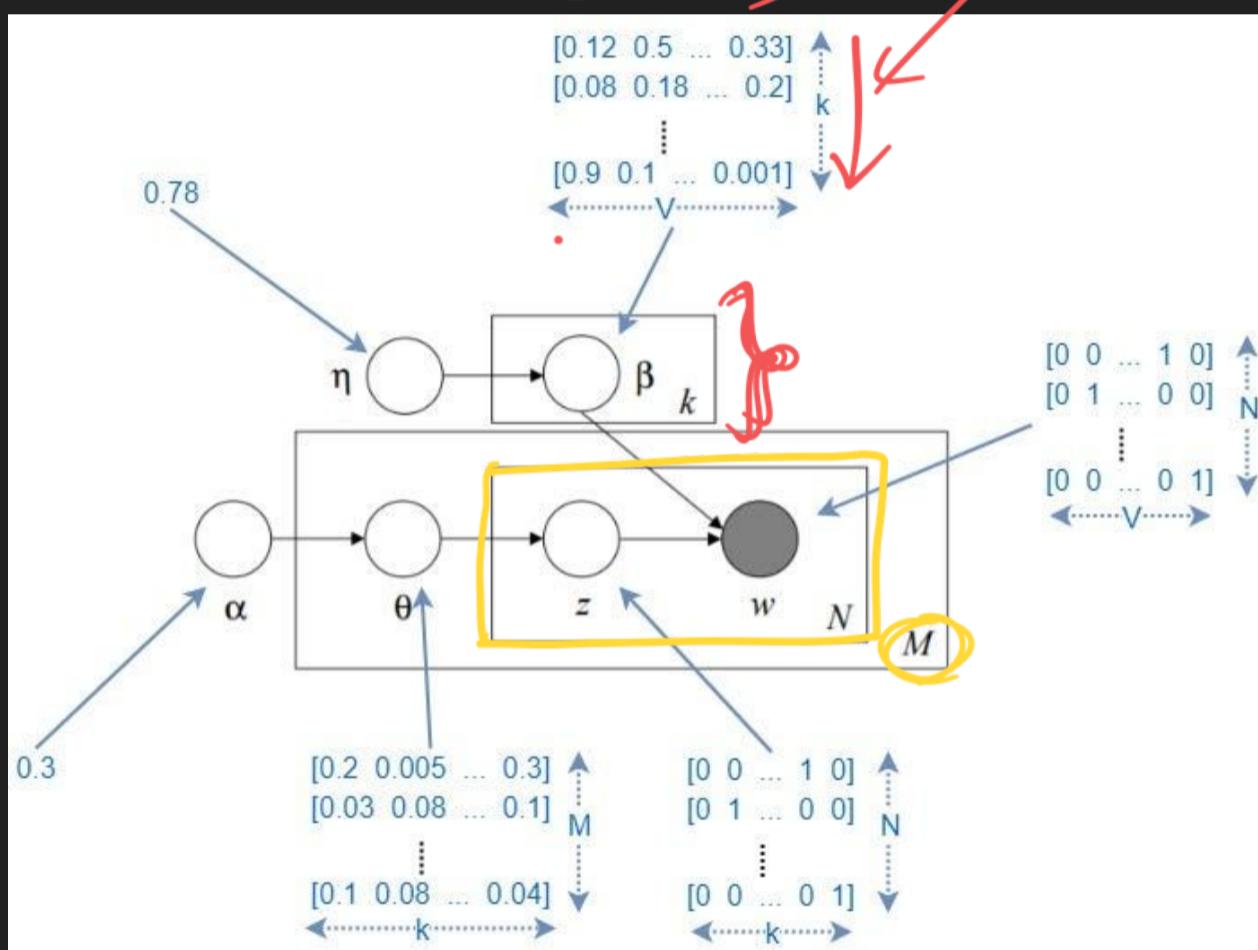
Doc1
 Doc2
 ...

↓
 DocN

| | | |
|----|---|---|
| W1 | 0 | 0 |
| W2 | 0 | 0 |
| W3 | 0 | 0 |
| W4 | 0 | 0 |
| W5 | 0 | 0 |



D2 $\xrightarrow{\text{no. of words}}$ $k = \text{top } i$



Sample 5 words

| | |
|------------|------------|
| players(1) | olympic(2) |
| playing(3) | olympic(4) |
| | cat(5) |

| | Cats | AG | Players | Olympic |
|----|------|----|---------|---------|
| W1 | 0 | 0 | 1 | 0 |
| W2 | 0 | 0 | 0 | 1 |
| W3 | 0 | 0 | 1 | 0 |
| W4 | 0 | 0 | 0 | 1 |
| W5 | 1 | 0 | 0 | 0 |

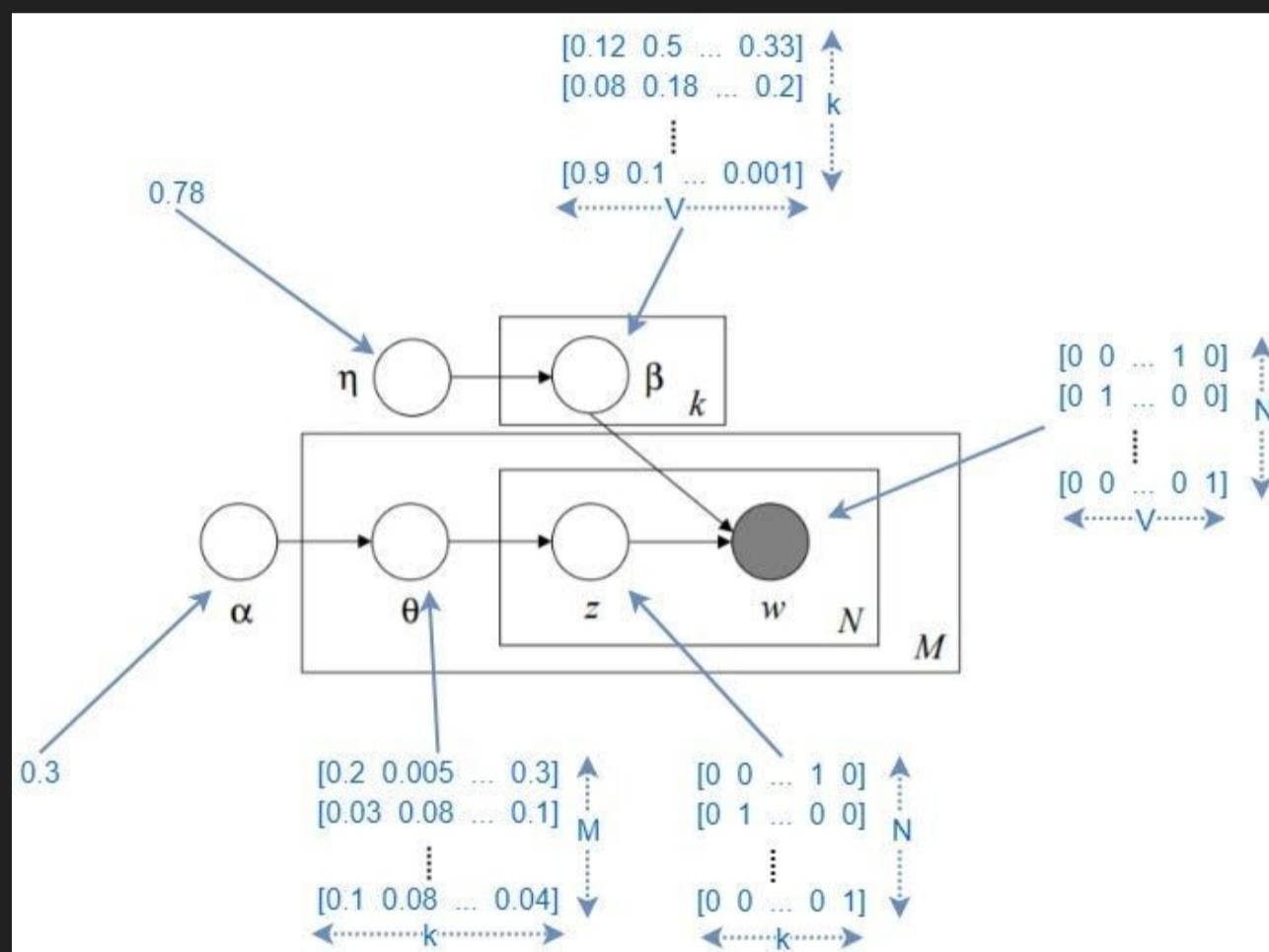
✓ Document 1

| | | |
|---|---|---|
| 0 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 1 | 0 | 0 |
| 0 | 0 | 1 |

| | Cats | AG | Players | Olympic |
|-----------|------|-------|---------|---------|
| T1 sports | 5% | 5% | 45% | 45% |
| T2 Animal | 90% | 8.33% | 3.3% | 3..33% |
| T3 Tech | 3.3% | 90% | 3-3% | 3.3% |

original doc

| DOC1 | W1 | players |
|------|----|---------|
| | W2 | olympic |
| | W3 | Players |
| | W4 | olympic |
| | W5 | cat |



$$\left\{ \begin{array}{l} 0 : \text{phone} \\ 1 : \text{guitar} \\ \quad \swarrow \\ \quad 3 \text{ times} \end{array} \right\}$$

all possible word in training data

$$\left\{ \begin{array}{l} \text{inter} \\ Q_1 : 3 \end{array} \right\}$$

