

* LDA & Topic Modeling: (Latent Dirichlet Allocation) Continued

↳ Latent Dirichlet Allocation → generative model.

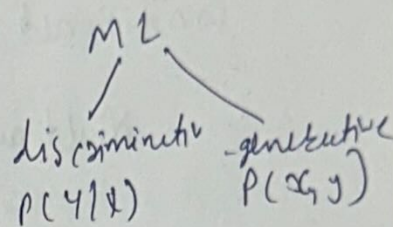
- In ML it is used for the task of Topic Modeling.

Topic Modeling is clustering docs into topics which are latent.

* Some Preliminaries

Multinomial distⁿ

- ~~Generalization~~ Generalization of
binomial distⁿ



- Suppose we have vector x in \mathbb{R}^k such that only one coordinate is 1 & rest are 0. $\hookrightarrow \{x: x \in \{0,1\}^k \text{ & } \|x\|_1 = 1\}$
- If we denote $P(x_k = 1) = \mu_k$ distribution of x is given as

$$P(x|\mu) = \prod_{k=1}^K \mu_k^{x_k}$$

↳ k dimensional vector.
(represent 1st one probab^l,
2nd probab^l :)

where $\mu = \{\mu_1, \dots, \mu_k\}$
 $\mu_k \geq 0$ &
 $\sum_k \mu_k = 1$

ex $x_A = 0, 0, 1$
 $x_B = 0, 1, 0$
 $x_C = 1, 0, 0$
 $\mu_1 = 0.7$
 $\mu_2 = 0.2$
 $\mu_3 = 0.1$

→ Now Consider denote $D \uparrow N$
 independent observation x_1, \dots, x_N
 likelihood

$$P(x_n|\mu) = (0.7)^0 (0.2)^0 (0.1)^1$$

$$P(D|\mu) = \prod_{i=1}^N \prod_{k=1}^K \mu_k^{x_{nk}}$$

$$= \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})}$$

$$= \prod_{k=1}^K \mu_k^{m_k} \quad \left(\because m_k = \sum_n x_{nk} \right)$$

(no. of observation of $x_k = 1$)

↳ $a^x \cdot a^y = a^{x+y}$

Now max^m likelihood solⁿ for μ is

$$\mu_k^{ML} = \frac{m_k}{N}$$

→ we consider joint distⁿ of quantities m_1, \dots, m_k conditioned on μ & N .

$$\text{Multinomial dist}^n(m_1, m_2, \dots, m_k | \mu, N)$$

$$= \binom{N}{m_1, m_2, \dots, m_k} \prod_{k=1}^K \mu_k^{m_k}$$

ex mult(10, 20, ..., 50 | μ ,)

where $\binom{N}{m_1, m_2, \dots, m_k} = \frac{N!}{m_1! m_2! \dots m_k!}$

$$N = \sum_{k=1}^K m_k = N$$

this is multinomial distⁿ

* Dirichlet distⁿ:

Conjugate prior for multinomial distⁿ.

Gamma distⁿ

is conjugate prior for itself

Dirichlet is conjugate prior for Multinomial

$$P(\mu | \alpha) \text{ varies as } \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

where $0 \leq \mu_k \leq 1$

$$\sum_k \mu_k = 1$$

$\alpha_1, \dots, \alpha_K$ are the parameters of distⁿ.

$$\alpha = (\alpha_1, \dots, \alpha_K)$$

↳ (prob. distⁿ over prob. dist)

ex (10 coins ~~flips~~ having 0.7 0.3 ...

... 0.1 & prob. of head

if we need find what is a prob. of head is given prob. distⁿ)

* Dirichlet

listⁿ is given. as

$$\text{Dir}(\mu | d) = \frac{\Gamma(d_0)}{\Gamma(d_1) \dots \Gamma(d_K)} \prod_{k=1}^K \mu_k^{d_k - 1}$$

$$d_0 = \sum_{k=1}^K d_k$$

$$\Gamma(x) = \int_0^\infty u^{x-1} e^{-u} du$$

(extension of
factorial funⁿ
for
reals)

$$\Gamma(x+1) = x \Gamma(x)$$

$$\Gamma(1) = 1$$

multinomial
vector (0 1 0 0 ...)

* Now, $P(\mu | D, \alpha)$

Posterior
also Dirichlet

$$\left(\prod_{k=1}^K \mu_k^{m_k} \right)$$

$$\propto \underbrace{P(D | \mu)}_{\text{Nuts symbol}} \underbrace{P(\mu | \alpha)}_{\text{proportional}}$$

$$\left(\prod_{k=1}^K \mu_k^{d_k - 1} \right)$$

$$\propto \prod_{k=1}^K \mu_k^{d_k + m_k - 1}$$

$$\text{So } (a^x \cdot c^y : a^{x+y})$$

$$\text{Thus } P(\mu | D, \alpha) = \text{Dir}(\mu | d + m)$$

$$= \frac{\Gamma(d_0 + N)}{\Gamma(d_1 + m_1) \dots \Gamma(d_K + m_K)}$$

(must be very clear)

++ Lecture in DSA books

* Lecture in DSA book

The posterior prob. $p(k|x)$ are given as

$$\begin{aligned} \gamma_k^{(n)} &= p(k|x) \\ &= \frac{p(k) p(x|k)}{\sum_k p(k) p(x|k)} \\ &= \frac{w_k N(x|\mu_k, \Sigma_k)}{\sum_k w_k N(x|\mu_k, \Sigma_k)} \end{aligned}$$

Parameters of GMM are w, μ, Σ

$$W = (w_1, \dots, w_k)$$

$$\mu = (\mu_1, \dots, \mu_k)$$

$$\Sigma = [\Sigma_1 \dots \Sigma_k]$$

* How to get values?

Log likelihood is given by $\ln p(x|w, \mu, \Sigma)$

$$= \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K w_k N(x_n | \mu_k, \Sigma_k) \right\}$$

* Another interpretation:

z is K -dim. binary random var. where at a time some z_k is 1 and all other coordinates are 0.

$$\therefore z_k \in \{0, 1\} \text{ and } \sum_k z_k = 1$$

- K possible states of var. z

- Try to define joint dist. $p(x, z)$ in terms of $p(z)$ & $p(x|z)$

$$\rightarrow \text{Now, } p(z_k = 1) = w_k \quad \left[\because 0 \leq w_k \leq 1 \right. \\ \left. \sum_k w_k = 1 \right]$$

$$\rightarrow p(z) = \prod_{k=1}^K w_k^{z_k}$$

(\therefore here there is assumption that every z_k is independent so it's multiplication. 4. one of the ~~vector~~ ^{value} is 1 & others are 0's.)

↑ ^{just like one hot encoding (1 0 0 ...)}

↑ ^{gaussian dist}

$$\text{Also } p(x | z_k = 1) = N(x | \mu_k, \Sigma_k)$$

$$\therefore p(x | z_k) = \prod_{k=1}^K (N(x | \mu_k, \Sigma_k))^{z_k}$$

Now $P(x) = \sum_z P(z) P(x|z)$ (marginal prob.)

$$= \sum_{k=1}^K w_k \cancel{P(x)} N(x|\mu_k, \Sigma_k)$$

\therefore Marginal of x is GMM if we have observations $x_1 \dots x_N$ and as $P(x) = \sum_z P(x, z)$ for every observed data point x_n there is corresponding latent var. z_n .

(Means which come from GMM which is hidden)

Now, $\gamma(z_k)$ i.e. P
(responsibility)

→ Maximizing log likelihood of GMM is more complex problem than the case of single Gaussian, because of summation over K inside log, log for "does not directly on Gaussian".
we

(Matrix cookbook ver)

* Topic Modeling & LDA

- given a corpus of doc. we consider that each word in each doc. has latent assignment into one of K topics. LDA is standard algo. for the task.
- each ~~word~~ doc. is mixture of topics.
- each topic is a P.D. over words

Multi $(m_1, m_2, \dots, m_K / \text{doc.}) : n$
 ex vocabulary size = 10000, probability of each word.

topic 1:	$\begin{bmatrix} 0.01 & 0.05 \end{bmatrix}$	chance
	$w_1 \quad w_2$	
(politics)	(election)	
topic 2:	$\begin{bmatrix} 0.5 & \dots \end{bmatrix}$	
topic K:	$\begin{bmatrix} \dots & \dots \end{bmatrix}$	

each topic 10000 dim. vector

* High level idea of generative model

- each doc. is generated by choosing topic
- each

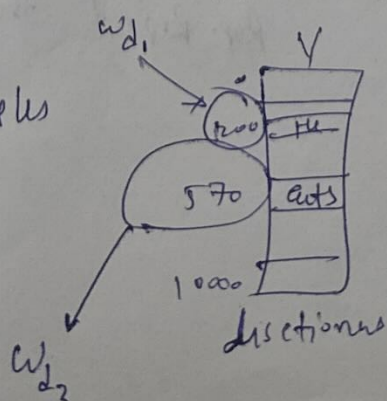
D - documents V - Vocabulary / dictionary

- we represent a doc. d with a vector $w_d \in \mathbb{N}^{n_d}$ where n_d is no. of terms doc.

$$w_d = (w_{d1}, w_{d2}, \dots, w_{dn_d}) \text{ where}$$

w_{di} stands for index i word in position i of doc. in the V .

ex (E) cuts apples
 2 in
 1 topic

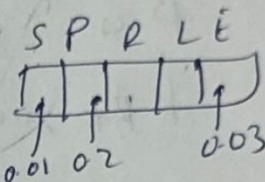


- Let K - num. of latent topics in corpus

- θ_d a vector K -dim. simplex Δ^{K-1} $\theta_d \in [0,1]$
 Shares 1 topic k in doc. d .

$\beta_k \in \Delta^{V-1}$ vector of term possibilities in topic k .

θ_d



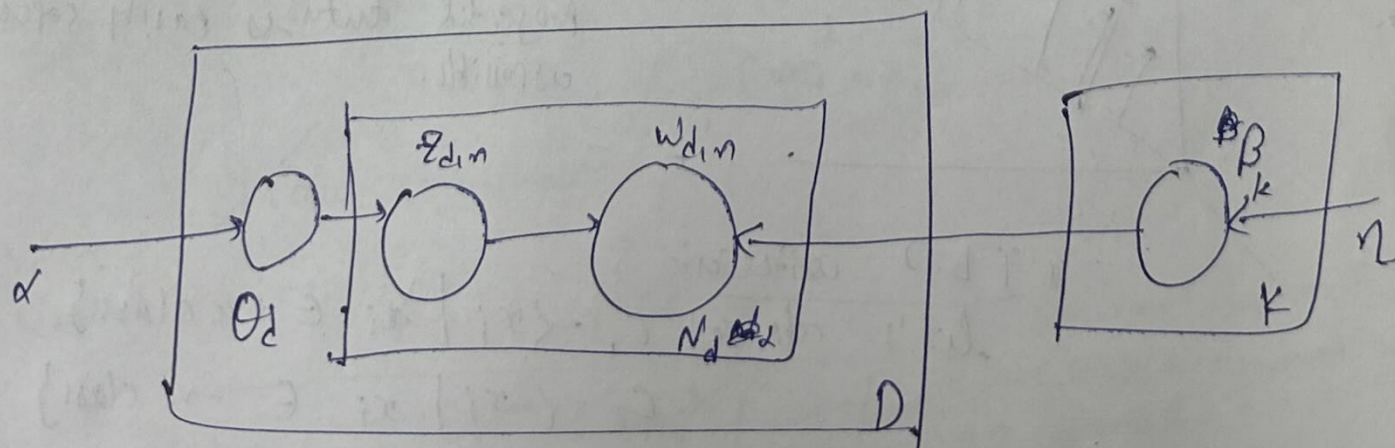
β_{kv} is the prob. of observing word v in topic k

probability of choosing k th topic for generating word

Consider the latent var. $z_{d,n} \in (1, \dots, K)$ is topic assigned to n th word in document d . (pertinence dir. what is proportion into 1 topic)

Algo.

1. Draw K vectors $\beta_k \in \Delta^{V-1}$ from $\text{Dir}(\beta_k | \eta)$
2. Draw D vectors $\theta_d \in \Delta^{K-1}$ from $\text{Dir}(\theta_d | \alpha)$
3. Each word $w_{d,n}$ in doc d is generated in two steps.
 - i. Draw $z_{d,n} \in (1, \dots, K)$ according to topic prob. θ_d .
 - ii. Draw $w_{d,n}$ using term prob. $\beta_{z_{d,n}}$

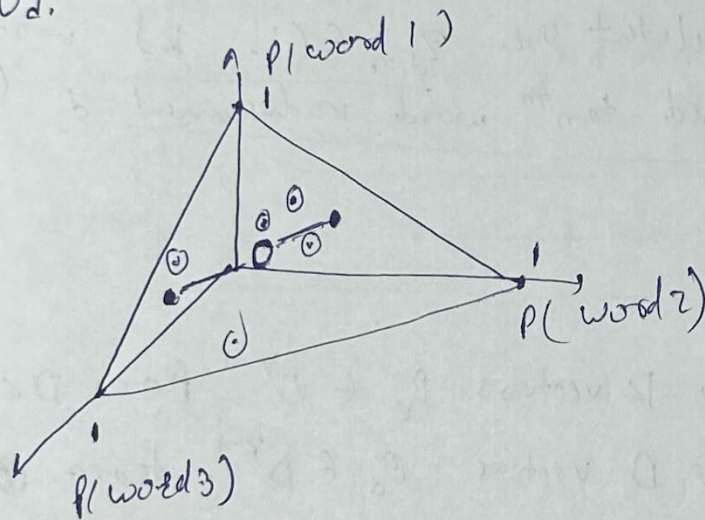


the joint for observed corpus $w_{d,n}$ given η, α :-

$$P(z, \theta, \beta | w, \alpha, \eta) = \left[\prod_{d,n} P(w_{d,n} | z_{d,n}, \beta) \right] \left[\prod_{d,n} P(z_{d,n} | \theta_d) \right]$$

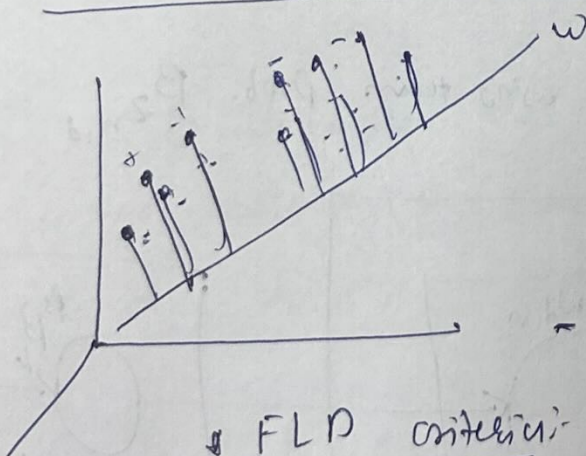
$$= \left[\prod_{d,n} p(\omega_{d,n} | z_{d,n}, \beta) \right] \left[\prod_{d,n} p(z_{d,n} | \theta_d) \right] \cdot \left[\prod_d p(\theta_d | \alpha) \right] \left[\prod_k p(\beta_k | \eta_k) \right]$$

- β is $V \times K$ matrix of term prob. with column k given by β_k
- θ is $K \times D$ matrix of topic membership with column d given by θ_d .



- - topic
- - observed topic
- - generated doc.

* Fischer's Linear Discriminant :



Goal: get a w for which projected data is easily separable as possible.

* FLD criteria:

let's set $C_1 = \{x_i | x_i \in \text{+ve class}\}$

$C_2 = \{x_i | x_i \in \text{-ve class}\}$

$$|C_1| = n_1$$

$$|C_2| = n_2$$

$$n = n_1 + n_2$$

$$m_1 = \frac{1}{n_1} \sum_{x_i \in C_1} x_i, \quad m_2 = \frac{1}{n_2} \sum_{x_i \in C_2} x_i$$

vectors of those classes with d. dim.

point $\in \mathbb{R}^d$

Sign of $\left(\frac{x_i^\perp}{\|x_i^\perp\|} \right) = x_i^T w$ \forall_i
 Perpendicular (Means projected x_i in w)
 which is scalar

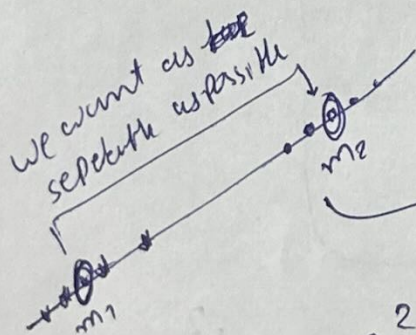
$$\rightarrow \text{if } \frac{1}{n_1} \sum_{x_i \in C_1} x_i^\perp = \frac{1}{n_1} \sum_{x_i \in C_1} x_i^T w$$

$$= \left(\frac{1}{n_1} \sum_{x_i \in C_1} x_i^T \right) w$$

$$= m_1^T w$$

$$= \frac{1}{m_1}$$

Similarly $m_2 = m_2^\perp = m_2^T w$



$$|m_1^\perp - m_2^\perp| = |w^T (m_1 - m_2)|$$

$$S_1^2 = \sum_{x_i \in C_1} (m_1^\perp - x_i^\perp)^2$$

(class 1 variance)

Scatter of class $C_1 \& C_2$

$$\rightarrow S_2^2 = \sum_{x_i \in C_2} (m_2^\perp - x_i^\perp)^2$$

\rightarrow we want $(m_1^\perp - m_2^\perp)^2$ as large & denominator is small.

$$J(w) = \frac{|m_1^\perp - m_2^\perp|^2}{S_1^2 + S_2^2}$$

scaler

one of the is larger known as FLD criterion

$$\begin{aligned} \tilde{m}^+ &= w^T m^+ - m_1 \\ \tilde{m}^- &= w^T m^- - m_2 \end{aligned}$$

$$\therefore \text{maximize } |\tilde{m}^+ - \tilde{m}^-|$$

$$= |w^T (m^+ - m^-)|$$

$$\therefore |\tilde{m}^+ - \tilde{m}^-|^2 = [w^T (m^+ - m^-)]^2$$

$$= w^T (m^+ - m^-) \cdot w^T (m^+ - m^-)$$

$$= \underbrace{w^T (m^+ - m^-) (m^+ - m^-)^T w}_{S_B \text{ (Rank 1)}}$$

S_B (Rank 1)
between class scatter
(Inter class) matrix

$$\therefore |\tilde{m}^+ - \tilde{m}^-|^2 = w^T S_B w$$

Scatter for projected points for +ve class

$$\begin{aligned} \tilde{S}^{+2} &= \sum_{x_i \in C^+} (\tilde{x}_i - \tilde{m}^+)^2 \\ \tilde{S}^{-2} &= \sum_{x_i \in C^-} (\tilde{x}_i - \tilde{m}^-)^2 \end{aligned}$$

$$\tilde{S}^{+2} = \sum_{x_i \in C^+} (w^T (x_i - m^+))^2$$

$$= \sum_{x_i \in C^+} \{w^T (x_i - m^+) (x_i - m^+)^T w\}$$

Scatter matrix for +ve class.

$$= w^T \left\{ \sum_{x_i \in C^+} (x_i - m^+) (x_i - m^+)^T \right\} w$$

Summation is rank 1 matrix

$$\tilde{S}^{+2} = w^T S^+ w$$

$$\tilde{S}^{-2} = \omega^T \tilde{S} \omega$$

$$\tilde{S}^{+2} + \tilde{S}^{-2} = \omega^T (S^+ + S^-) \omega$$

$$= \omega^T S_\omega \omega \Rightarrow S_\omega = S^+ + S^-$$

(within class scatter class $S^+ + S^-$)
(intra)

* FLD Criteria:

$$J(\omega) = \frac{|\tilde{m}^+ - \tilde{m}^-|^2}{\tilde{S}^{+2} + \tilde{S}^{-2}}$$

$$= \frac{\omega^T S_B \omega}{\omega^T S_\omega \omega}$$

(S_B, S_ω all known)

$$\omega_{\text{fisher}} = \underset{\omega}{\text{argmax}} J(\omega)$$

$$J(\omega): \mathbb{R}^d \rightarrow \mathbb{R}$$

Fisher hyperplan.

* Maximizing the $J(\omega)$

$$\frac{\partial J(\omega)}{\partial \omega} = \frac{(\omega^T S_\omega \omega)(2\omega^T S_B) - (\omega^T S_B \omega)(2\omega^T S_\omega)}{(\omega^T S_\omega \omega)^2}$$

$$= \frac{2\omega^T (\omega^T S_\omega \omega S_B - \omega^T S_B \omega S_\omega)}{(\omega^T S_\omega \omega)^2}$$

$$= \frac{(\omega^T S_\omega \omega)(S_B \omega) - (\omega^T S_B \omega)(S_\omega \omega)}{(\omega^T S_\omega \omega)^2}$$

$$\frac{\partial J(\omega)}{\partial \omega} \Big|_{\omega = \omega^*} = 0$$

$$(\omega^{*T} S_{\omega} \omega^*) (S_B \omega^*) = (\omega^{*T} S_B \omega^*) (S_{\omega} \omega^*)$$

$$S_B \omega^* = \frac{(\omega^{*T} S_B \omega^*) (S_{\omega} \omega^*)}{(\omega^{*T} S_{\omega} \omega^*)}$$

$$S_B \omega^* = J(\omega^*) (S_{\omega} \omega^*)$$

2

$$AV = \lambda BV$$

$$AV = \lambda IV$$

generalized Eigen
value Problems

$$S_B \omega^* = (m^+ - m^-) \underbrace{(m^+ - m^-)^T}_{\text{Scalar}} \omega^*$$

This is some vector along with direction $(m^+ - m^-)$

$$\therefore S_B \omega^* = J(\omega^*) \cdot (S_{\omega} \omega^*)$$

$$J(\omega^*) S_B \omega^* = \alpha (m^+ - m^-)$$

$$J(\omega^*) S_{\omega} \omega^* = (m^+ - m^-) \underbrace{(m^+ - m^-)^T \omega^*}_{\text{Scalar}}$$

$$S_{\omega} \omega^* = \alpha (m^+ - m^-)$$

$$\left(\alpha = \frac{(m^+ - m^-)^T \omega^*}{J(\omega^*)} \right)$$

$$\therefore \omega^* = \alpha S_{\omega}^{-1} (m^+ - m^-)$$

direction ω^*

It is invertible