

# Modelling

→ Latent Dirichlet Allocation  
is generative model

↳ used for topic modelling

Topic modeling is clustering  
of documents into topics  
which are latent.

- Some preliminaries:-
- multinomial distribution
  - generalization of binomial distribution
  - suppose we have vector  $x$  into  $\mathbb{R}^K$  such that only one coord. is 1 rest are 0
  - If we denote  $P(X_K = 1) = \mu_K$  distribution of  $x$  is given as

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

where

$$\mu = (\mu_1, \dots, \mu_K) \quad \sum_k \mu_k = 1$$
$$\mu_k \geq 0$$

- now consider dataset D of N independent observations  $x_1, \dots, x_N$

$$\text{Likelihood} = \prod_{i=1}^N \prod_{k=1}^{m_k} u_k$$

$$\Rightarrow \prod_{i=1}^N u_k^{(\sum_n x_{nk})}$$

$$= \prod_{k=1}^K u_k^{m_k} \quad [m_k = \sum_n x_{nk}]$$

[no of observations of  $x_k = 1$ ]

→ Now maximum likelihood solution for  $u$  is

$$\hat{u}_k^{\text{ML}} = m_k / N$$

- we consider joint dist. of quantities  $m_1, \dots, m_K$  conditioned on  $u$  &  $N$

$$\text{mult}(m_1, \dots, m_K | u, N)$$

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

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

$$\sum_{k=1}^K m_k = N \quad (\text{mult. dist.})$$

### • Dirichlet Distribution:-

- Conjugate prior for multinomial distribution

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

$$\text{where } 0 \leq u_k \leq 1$$

$$\sum_k u_k = 1$$

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

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

The dirichlet dist. is given as

$$\text{Dir}(u|\alpha) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_K)} \prod_{k=1}^K u_k^{\alpha_k - 1}$$

$$\alpha_0 = \sum_{k=1}^K \alpha_k$$

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

(extension of factorial func for real)

$$\begin{aligned} - \Gamma(\alpha+1) &= \alpha \Gamma(\alpha) \\ - \Gamma(1) &= 1 \end{aligned}$$

Now  $\alpha$  alpha  
 $P(u|D, \alpha)$

$$\alpha \quad P(D|u) \cdot P(u|d)$$

(varies symbol)

$$\alpha \prod_{k=1}^K \alpha_{ik} + m_k - 1$$

Thus

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

$$= \frac{\Gamma(\alpha_0 + n)}{\Gamma(\alpha_1 + m_1) \dots \Gamma(\alpha_K + m_K)}$$

$$m = (m_1, \dots, m_K)$$

- Topic Modelling & LDA

- Given a corpus of texts we consider that each word

# Ch-1 SUMS

Page No.

Date

① 2 opaque bags  $\rightarrow$  each containing 2 balls

one bag  $\rightarrow$  2 black balls

other  $\rightarrow$  1 black & 1 white

pick a bag at random & pick one of the balls

when you look at the balls

it's black & picked second ball from same bag

prob that ball also black?

$\rightarrow B_1$ : picked ball black

$B_2$ : second picked ball black.

$$P(B_2 | B_1) = \frac{P(B_2 \cap B_1)}{P(B_1)}$$

$$P(B_1) = P(B_1 | \text{Bag 1}) P(\text{Bag 1}) + P(B_1 | \text{Bag 2}) P(\text{Bag 2})$$

$$1 \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2}$$

$$= \frac{3}{4}$$

$\rightarrow P(B_1 \cap B_2)$

$$= P(B_1 \cap B_2 | \text{Bag 1}) P(\text{Bag 1}) + P(B_1 \cap B_2 | \text{Bag 2}) P(\text{Bag 2})$$

$$1 \cdot \frac{1}{2} + 0 \cdot \frac{1}{2} = \frac{1}{2}$$

$$P(B_2 | B_1) = \frac{1/2}{3/4} = \frac{2}{3}$$

(2) consider perceptron in 2D  $h(x) = \text{sign}(\omega^T x)$  where  $\omega = [\omega_0 \omega_1 \omega_2]^T$   
 and  $x = [1 x_1 x_2]^T$

(3) show that regions on plane where  $h(x) = +1$ ,  $h(x) = -1$  are separated by a line

if we express this line by  $x_2 = a x_1 + b$

what are slope 'a' & intercept 'b' in terms of  $\omega_0, \omega_1, \omega_2$ .

(b) case  $\omega = [1, 2, 3]^T$ ,  $\omega = -[1, 2, 3]^T$

→ Soln: -  $h(x) = \text{sign}(\omega^T x)$ .

(4)  $h(x) = +1$ ,  $\omega^T x > 0$

↳ so separation b/w 2 regions is line whose eq. is  $\omega^T x = 0$

$$\omega_0 + \omega_1 x_1 + \omega_2 x_2 = 0$$

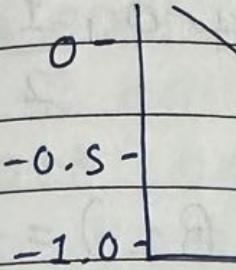
↳ also can write  $x_2 = a x_1 + b$

$$\text{so, } a = \frac{-\omega_1}{\omega_2}, \quad b = \frac{-\omega_0}{\omega_2}$$

$$(b) \omega = (1, 2, 3)^T$$

$$\omega = -(1, 2, 3)^T$$

(Identical line  
both case)



(b)

$$\omega^T c(t) \omega$$

conclude

$$\hookrightarrow \omega^T c(t) \omega^* =$$

$$= \omega^T c t$$

$$\geq \omega^T c(t) \cdot$$

(3)

prove that PLA eventually converges to linear separator for separable data.

- $\omega^*$  be an optimal set of value weights,
- proof shows PLA weights  $\omega(t)$  get more aligned with  $\omega^*$  with each iteration.  $\omega(0) = 0$

(a)

Let  $\rho = \min_{1 \leq n \leq N} y_n (\omega^* \cdot x_n)$   
show that  $\rho > 0$

as every  $x_n$  well classified by  $\omega^*$  for all  $n = 1, \dots, N$

$$y_n = \text{sign}(\omega^* \cdot x_n)$$

$$y_n (\omega^* \cdot x_n) > 0$$

- for all  $n = 1, \dots, N$

$$\rho = \min_n y_n (\omega^* \cdot x_n) > 0$$

$\hookrightarrow$  Ind  
obvious

$$(b) \omega^T(t) \omega^* \geq \omega^T(t-1) \omega^* + p$$

conclude that  $\omega^T(t) \omega^* \geq t_p$

$$\hookrightarrow \omega^T(t) \omega^* = [\omega^T(t-1) + y(t-1)x^T(t-1)]\omega^*$$

$$= \omega^T(t-1) \omega^* + y(t-1) \omega^{*T} x(t-1)$$

$$\geq \omega^T(t-1) \omega^* + p$$

$\hookrightarrow$  Induction, If  $t=0$  we  
obviously get  $0 \cdot \omega^* \geq 0$

(6) Sample of 10 marbles from bin that holds red & white marbles. Prob of red marble = 0.8 for  $\mu = 0.05$ ,  $\lambda = 0.5$ . Compute prob of getting no red marbles ( $V=0$ ) in following case.

- (a) we draw only one such sample compute prob  $V=0$
- (b) we draw 1000 independent samples, compute prob that one of sample has  $V=0$
- (c) repeat for 100000 ind. samples

Soln :-

(a) for one sample  
 $P(V=0) = (1-\mu)^{10}$

$$\mu = 0.05 \rightarrow P(V=0) = 0.59873$$

$$\mu = 0.5 \rightarrow P(V=0) = 9.7656 \times 10^{-6}$$

$$\mu = 0.8 \rightarrow P(V=0) = 1.024 \times 10^{-7}$$

(b) for 1000 independent samples

$P(\text{At least one sample has } V=0)$

$$= 1 - P(\cup_{i=1}^{1000} V_i > 0)$$

$$= 1 - \prod_{i=1}^{1000} (P(V_i > 0))$$

$$\begin{aligned} &= 1 - (1 - 0.59873)^{1000} \\ &= 1 - e^{-1000 \cdot 0.59873} \\ &= 1 - e^{-598.73} \end{aligned}$$

(c) for 1000 samples

$P(\text{At least one sample has } V=0)$

(7) sample of tossing coin. Assume  $\mu$  generate

for coin tosses generate

$P(K/N)$

training

① assume sample size ( $N$ ) = 10

if all coins have  $\mu = 0.05$

compute prob at least one coin will have  $V=0$  for case of 1 coin  
1000 coin, 1000000 coins.

repeat for  $\mu = 0.8$

$$\downarrow \\ \text{So } h^n \rightarrow \mu = 0.05$$

$P(\text{at least one coin has } V=0)$

$$= (1 - 0.05)^{10}$$

$$= 0.05987$$

→ for 1000 coins =

$$1 - [1 - (1 - 0.05)^{10}]^{1000} = 1;$$

→ for 1000000 coins

$$1 - [1 - (1 - 0.05)^{10}]^{1000000}$$

②  $N=6, 2$  coins with both coins plot prob  $P[\max |V_i - \mu_i| > \epsilon]$

$$L \leq G \in [0, 1]$$

hoeffding bound

$$\downarrow \\ P[|\bar{V} - \mu| > \epsilon] \leq 2e^{-2N\epsilon^2}$$

$\downarrow$  So  $h^n, N=6, \mu=0.5 \rightarrow$  Hoeffding inequality  
2 coins

$$P\left(\max_i |\bar{V}_i - \mu_i| > \epsilon\right)$$

$$P(|V_1 - u_1| > \epsilon \text{ or } |V_2 - u_2| > \epsilon)$$

$$= P(|V_1 - u_1| > \epsilon) + P(|V_2 - u_2| > \epsilon)$$

$$- P(|V_1 - u_1| > \epsilon \text{ and } |V_2 - u_2| > \epsilon)$$

$$= P(|V_1 - u_1| > \epsilon) + P(|V_2 - u_2| > \epsilon)$$

$$- P(|V_1 - u_1| > \epsilon) P(|V_2 - u_2| > \epsilon)$$

= 1

coins  
if  $|V_i| > \epsilon$

$$\begin{aligned}
 &= 1 - \prod_{i=1}^{1000} [1 - P(V_i=0)] \\
 &= 1 - \prod_{i=1}^{1000} [1 - (1-u)^{10}]^{1000} \\
 &= 1 - [1 - (1-u)^{10}]^{1000} \\
 u = 0.05 \rightarrow \text{ans} = 1
 \end{aligned}$$

(c) for 100000 independant sample.

$$P(\text{At least one sample has } V=0) = 1 - (1 - (1-u)^{10})^{100000}$$

(7) sample of head & tail created by tossing coin a no. of times independently. Assume we have no. of coins that generate diff. samples independently

→ for coin → prob of head =  $u$   
 prob of obtaining  $k$  heads in  $N$  tosses given by binomial dist.

$$P[k/N, u] = \binom{N}{k} u^k (1-u)^{N-k}$$

training error  $v = k/N$

Page No. \_\_\_\_\_  
Date \_\_\_\_\_

2.1 eq:  $E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}$

in this eq. set  $\delta = 0.09$

let  $E(M, N, \delta) = \sqrt{\frac{1}{2N} \ln \frac{2M}{\delta}}$

- ⑥ for  $M=1$ , how many examples do we need to make  $\epsilon \leq 0.05$ ?

$M=1, \delta=0.03, \epsilon \leq 0.05$

$$N \geq \frac{1}{2\epsilon^2} \ln \frac{2M}{\delta}$$

$$\geq \frac{1}{2 \cdot (0.05)^2} \ln \frac{2}{0.03}$$

$$= 839.94$$

- ② Show that for learning model for tve rectangles

$$m_H(4) = 2^4, m_H(5) < 2^5$$

Hence give a bound for  $m_{H(N)}$

for  $N=4$ , if we consider 4 non-aligned call dichotomis generat H shatters these points so we have  $m_H(4) = 2^4$

→ for  $N=5$ , no matter how you place 5 points, some points inside rectangle

- not able to generate all dichotomies  
 $\therefore m_H(5) < 2^5$

L for rectangle  $dvc = 4$

$$m_{H(N)} = N^4 + 1$$

③ compute maximum no. of dichotomies  $m_H(N)$  for these learning models. and consequently compute  $dvc$ , the VC dim.

→ ① positive or negative ray

$H$  contains func. which are  $+1$  on  $[a, \infty)$  for some  $a$ . together with those that are  $+1$  on  $(-\infty, q]$

Soln:- we already know that growth function for the rays is equal to  $N+1$

- if we enumerate dichotomies added by negative rays, we get  $N-1$  new dichotomies.

- you got opposite of ones from positive rays and you have to subtract 2 dichotomies where all points are  $+1$  &  $-1$ .

$$\therefore m_H(N) = 2N$$

- As the largest value of  $N$  for which  $m_H(N) = 2^N$   
 (as  $m_H(3) = 6$ )  
 so  $dVC = \underline{2}$

b) positive / neg. interval

- $H$  contains func. which are +1 on an interval  $[a, b]$  & -1 elsewhere on an interval  $[a, b]$ , +1 elsewhere

→ growth func for the interval  
 $= N^2/2 + N/2 + 1$

If we add new dichotomies generated by -ve interval we get  $N-2$  new ones.

$N=1 \rightarrow$  already generate 2 dichotomies with the interval

$$m_H(N) = \frac{N^2}{2} + \frac{3N}{2} - 1 \quad N \geq 1$$

$$\times \quad m_H(N) = 2 \text{ if } N=1$$

largest value of  $N$  for which  
 $m_H(N) = 2^N \rightarrow ③$   
 $dVC = 3$

① 2 concentric spheres :-

↳  $H$  contains func. which are  
+1 for  $a \leq \sqrt{x_1^2 + \dots + x_d^2} \leq b$

↳ growth func of concentric circle  
map.

$$R^d \rightarrow [0, +\infty]$$

$$\phi(x_1, \dots, x_d) \rightarrow r = \sqrt{x_1^2 + \dots + x_d^2}$$

Same as the interval

$$m_H(N) = \frac{N^2}{2} + \frac{N}{2} H,$$

independent of  $d$ .

As the largest value of  $N$  for  
which  $m_H(N) = 2^N$ .

$$\hookrightarrow dvc = 2$$

\* Show that  $m_H(2N) \leq m_H(N^2)$   
hence obtain generalization  
bound. which involves  $m_H(N)$

→ ex: we have 3 ways to dichotomize  
2 points  $x_1, x_2 \rightarrow$  & 2 ways to  
dichotomize another 2 points  
here we have  $3 \times 2 = 6$  ways to  
dichotomize all 4 points.

let us say  $m_{H(N)} = p$   
 now if we partition any set of  $2N$  points into set of  $N$  points each, each of these 2 partition produce  $p$  dichotomies at most.

If we combine  $\rightarrow$  max. dicho. would be prod. of  $p$  by  $p$   
 $m_H(2N) = m_{H(N)}^2$

\* Suppose  $m_{H(N)} = N+1$  so  $dvc=1$   
 you have 100 training ex. Use generalization bound to give a bound for  $E_{out}(g)$  with 90% confidence. Repeat for  $N=1$

$$N = 100$$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{100} \ln \frac{4(2 \cdot 100+1)}{0.1}}$$

$$= E_{in}(g) + 0.8481596$$

$\rightarrow$  with prob at least 90%.  $N=10000$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{8}{10000} \ln \frac{4(2 \cdot 10000+1)}{0.1}}$$

\* let  $H = \{$   
 prove that  
 $dvc(H) = \dots$   
 $m_H(d) = 2$

$$2^d = m_{H(d)}$$

\* let  $H = \{h_1, \dots, h_m\}$  with finite  $M$ .  
 prove that  $\text{dvc}(H) \leq \log_2 M$

$\text{dvc}(H) = d$ , by definition we have  
 $m_H(d) = 2^d$

$$\begin{aligned} 2^d &= m_H(d) = \max_{x_1, \dots, x_d} |H(x_1, \dots, x_d)| \\ &= \max_{x_1, \dots, x_d} |\{h(x_1), \dots, h(x_d)\} : h \in H\}| \\ &\leq |H| = M \end{aligned}$$

$$\text{so } d \leq \log_2(M)$$

\* for hypothesis set  $H_1, \dots, H_K$  with finite VC dim  $\text{dvc}(H_k)$  derive and prove tightest upper & lower bound that you can get on  $\text{dvc}$