

①

P7

score metric for  $(a, b)$  word pair  $S_{a,b}(x, y) = \begin{cases} \cos(\bar{a} - \bar{b}, \bar{x} - \bar{y}) & \text{if } \|\bar{x} - \bar{y}\| \leq \delta \\ 0 & \text{otherwise} \end{cases}$

Analogy:  $\vec{w} \in \mathbb{R}^d$ ,  $\|\vec{w}\| = 1$

seed direction: (she, he)

similarity:  $\uparrow$  in experience

## Direct gender bias

Definition

$$\text{Direct-bias}_C = \frac{1}{|N|} \sum_{w \in N_{\text{set}}} |\cos(\vec{w}, g)|$$

indirect

Bias Definition

$$\beta(w, v) = \left( w \cdot v - \frac{w_{\perp} \cdot v_{\perp}}{\|w_{\perp}\|_2 \|v_{\perp}\|_2} \right) / w \cdot v$$

gender subspace top principle component

strictness

?? confused here.

projection

(Described in next page)

$$w_{\perp} = w - w_g; \quad v_{\perp} = v - v_g; \quad w_g = (w \cdot g) g$$

if 0 means no projection.

$\Downarrow$

$$\beta(w, v) = 0$$

unit vector  $\rightarrow$  orthogonal to each other

B subspace  $\{b_1, \dots, b_k\} \in \mathbb{R}^d \parallel k=1 \Rightarrow$  vector.

original vector

Debiasing Algorithm: Projection direction  $\vec{v}_B = \sum_{d=1}^k (\vec{v} \cdot \vec{b}_d) \vec{b}_d$

$$\beta \text{ value} = \frac{\vec{v}^T \vec{v}_B}{\|\vec{v}_B\|} // \text{orthogonal project: } \vec{v} - \vec{v}_B$$

step 1: identify gender subspace:

Defining sets:  $D_1, \dots, D_n \subset W$  total words

$$\text{mean of } D_i \Rightarrow \mu_i = \sum_{w \in D_i} \vec{w} / |D_i|$$

$\{\vec{w} \in \mathbb{R}^d\}$  word vector  $k \geq 1$

Let Bias subspace

$$C := \sum_{i=1}^b \sum_{w \in D_i} (\vec{w} - \mu_i)(\vec{w} - \mu_i)^T / |D_i| \quad \left[ \begin{array}{l} k \text{ row of SVD} \\ \text{is bias subspace } \mathcal{B} \end{array} \right]$$

(11)

p7

Hard de-biasing:  $\bar{w} := (\bar{w} - \bar{w}_B) / \|\bar{w} - \bar{w}_B\|$  // reembedding definition.   
 new embedding  $\nearrow$  orthogonal projection

$\therefore$  word to Neutralize  $N \subseteq W$   $\uparrow$

family/Equality set  $E = \{E_1, E_2, \dots, E_m\}$  // ?? what we want equidist.   
 got B matrix earlier  $\checkmark$    
 new  $\rightarrow$  step

$$E_i \subseteq W$$

finally all words will have similar component in gender neutral direction

$$\mu := \sum_{w \in E} \frac{w}{|E|}$$

$$v := \mu - \mu_B \quad \text{projection to } B \quad \text{orthogonal}$$

this term varies only for words in E set

$$\text{For } \forall \underline{w \in E}; \bar{w} := \bar{v} + \frac{\bar{w}_B - \mu_B}{\|\bar{w}_B - \mu_B\|}$$

Added for the bias component differences.

output subspace  $B$ , new embedding  $\{\bar{w} \in \mathbb{R}^d\}_{w \in W}$

$$(\bar{v}_B, \bar{w}_{\perp B} = w - w_B)$$

soft bias projection:  $W \in \mathbb{R}^{d \times |\text{vocab}|}$  new -

$T \rightarrow$  transformation  $d \times d$

$$\min_T \left\| (T^T W)^T (T W) - W^T W \right\|_F^2 + \lambda \left\| (T^T N)^T (T B) \right\|_F^2$$

matrix size vocab vocab   
 optimization problem   
 [matrix of the neural embedding words]

Measurement of indirect bias: Between two gender neutral words  $(w, v)$

(measurement)   
 indirect bias  $\beta(\bar{w}, \bar{v}) = \frac{\bar{w}^T \bar{v}}{\bar{w}^T \bar{v}} - \frac{\bar{w}_{\perp B}^T \bar{v}_{\perp B}}{\|\bar{w}_{\perp B}\| \|\bar{v}_{\perp B}\|} \Rightarrow$  [match in gender independent direction]   
 overall match.

①

P.9

simplified version of pagerank:

$$R(u) = c \sum_{v \in B_u} \frac{R(v)}{N_v} \rightarrow \text{Number of links from } u$$

Back link of  $u$  (towards  $u$ )  
 Normalization.  
 Recursive equation.

$A$  matrix  $\rightarrow$  page  $\times$  page big matrix.

$$A_{u,v} = \frac{1}{N_u} \text{ if edge exists.}$$

0 if no edge between  $u, v$

$\rightarrow$  A vector! (bad notation)

$$R = \frac{1}{c} A R$$

$\rightarrow$  eigen vector of  $A$

Eigen vector of  $A \rightarrow$  symmetric matrix.

$\downarrow$   
 All vectors are orthogonal.

Dominant eigenvector

$\rightarrow$  power iteration.

Ranksource modification.

$$R'(u) = \frac{c}{m} \sum_{v \in B_u} \frac{R'(v)}{N_v} + \frac{c}{m} E(u)$$

Source rank vector.

$L_1$  norm,  $\|R'\|_1 = 1$ ,  $c$  is maximized.

if  $E(u) > 0$  the  $u$  is reduced.

Decay factor.

All 1's matrix

$$R' = c (A R' + E) = c (A + E \times \mathbf{1}) R'$$

$\rightarrow$  since  $\|R'\|_1 = 1$   
 Eigenvalue of this one.

(u)

P3 $R_0 \leftarrow s$  // random ints.

loop:

 $\rightarrow R_{i+1} \leftarrow AR_i$  // power iteration (PE)(dominant)  
finding eig vector for A $\rightarrow d \leftarrow \|R_i\|_1 - \|R_{i+1}\|_1$  // constrained. $\rightarrow R_{i+1} \leftarrow R_{i+1} + dE$  // little move from PE $\rightarrow \delta \leftarrow \|R_{i+1} - R_i\|_1$  increases convergence,  
maintain  $\|R\|_1$ ,  
// normalize.while  $\delta > \epsilon$  // convergence

P12: Recommendation problem formulation.

$C \rightarrow$  set of users (userspace,  $\rightarrow$  name, age, demogreaph, ...)

$S \rightarrow$  possible items. (name, title, producers etc)

$u$ , utility function, usefulness between (user, item)

$u: C \times S \rightarrow R$  (rating value  $\rightarrow$  utility)

so objective,  $\forall c \in C$ ,  $s_c' = \arg \max_{s \in S} u(c, s)$

Content based filtering methods:

focus on  $u(c, s_i)$  user already has rated.

$\rightarrow$  similar to previous  $s_i$  will be recommended

Term frequency  $TF_{ij} = \frac{f_{ij}}{\max_k f_{k,j}}$  ;  $f_{ij}$  no time  $k$  word appear in document,  $d_j$

inverse document frequency  $IDF_i = \log \frac{N}{n_i}$  ;  $N \rightarrow$  total documents.  
 $n_i \rightarrow$   $k_i$  appeared in how many documents.

term weight of keyword  $k_i$ , in document  $d_j$

$$w_{ij} = TF_{ij} \times IDF_i$$

for content of document  $d_j$

for all the key words  $k_i$

$$\text{Content}(d_j) = (w_{1j}, w_{2j}, \dots, w_{kj})$$

(11)

P12

content based profile  $(c) = \{w_{c1}, w_{c2}, \dots, w_{ck}\}$  for keyword  $(w)$  in system.

The utility function  $u(c, s) = \text{score}(\text{content based prof } (c), \text{Content}(s))$

$$u(c, s) = \cos(\vec{w}_c, \vec{w}_s) = \frac{\vec{w}_c \cdot \vec{w}_s}{\|\vec{w}_c\|_2 \times \|\vec{w}_s\|_2}$$

collaborative method:

$$u(c, s) \leftarrow u(c', s) ; c' \in C \text{ \& } c' \approx c$$

(similar user group)  
(peer)

① memory based / Heuristic.

rating,  $(r_{cs}) =$  aggregate  $r_{c',s}$  // impute unknown value.  
 not given but estimated  $\uparrow$  gives ratings

Agg. function can be?

$$r_{cs} = \begin{cases} \text{a) } \frac{1}{N} \sum_{c' \in C} r_{c',s} \\ \text{b) } k \sum_{c' \in C} \text{sim}(c, c') \times r_{c',s} \\ \text{c) } \bar{r}_c + k \sum_{c' \in C} \text{sim}(c, c') \times (r_{c',s} - \bar{r}_{c'}) \end{cases}$$

more like collaboration. term.

where,  $k = \frac{1}{\sum_{c' \in C} |\text{sim}(c, c')|}$  (normalizing constant)

$$\bar{r}_c = \left( \frac{1}{|S_c|} \right) \sum_{s \in S_c} r_{cs} \quad \text{where } S_c = \{s \in S \mid r_{c,s} \neq \emptyset\}$$

(Average rating)

(11)

P12

Pearson coefficient based similarity

$$\text{sim}(x, y) = \frac{\sum_{s \in S_{xy}} (\hat{r}_{x,s} - \bar{r}_x) (\hat{r}_{y,s} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (\hat{r}_{x,s} - \bar{r}_x)^2 \sum_{s \in S_{xy}} (\hat{r}_{y,s} - \bar{r}_y)^2}}$$

Iterate over (s) item

Alternatively, Cosine based similarity.  
for each item (s)

(ii) model based Algorithm:

$$\hat{r}_{c,s} = E(\pi_{c,s}) = \sum_{i=0}^n i \times \Pr(\pi_{c,s} = i) \quad \hat{r}_{c,s'}, s' \in S_c$$



①

P 13 selected dimensionality  $f$ :

$q_i \in \mathbb{R}^f$  // item

$p_u \in \mathbb{R}^f$  // user.

interaction between user  $u$  & item  $i$

the approx. rating  $\hat{r}_{ui} = q_i^T p_u$  — ①

how to get it?

SVD? but empty elements??

imputation  $\rightarrow$  bad idea.

So, optimization problem:

$$\min_{p, q} \sum_{(u,i) \in K} (r_{ui} - q_i^T p_u)^2 + \lambda (\|q_i\|^2 + \|p_u\|^2) \quad \text{--- ②}$$

set of given training set (previously observed)  
(explicitly feedback points)

Learning methods

① SGD

$$e_{ui} := r_{ui} - q_i^T p_u$$

$$q_i \leftarrow q_i + \delta (e_{ui} \cdot p_u - \lambda q_i) \quad // \text{Gradient descent}$$

( $\delta$  is changed by internal calculation)

$$p_u \leftarrow p_u + \delta (e_{ui} \cdot q_i - \lambda p_u) \quad + \text{fact}$$

②

Alternate // ALS

to solve nonconvexity  $\rightarrow$  fix one, and solve for the other.



(11)

P 13existence of product/user bias.

modify eq (1) by  $\boxed{b_{ui} = \mu + b_i + b_u}$

So,  $r_{ui} = \underbrace{\mu + b_i + b_u}_{\text{overall rating}} + q_i^T p_u \quad \text{--- (iv)}$

Now the optimization problem changes to,

$$\min_{p, q, b} \sum_{(u,i) \in K} (r_{ui} - \mu - b_i - q_i^T p_u)^2 + \lambda (\|p_u\|^2 + \|q_i\|^2 + b_u^2 + b_i^2) \quad \text{--- (v)}$$

// may bias the model.

Additional input source: cold start overcome. $N(u)$  + implicit preference on items by the users.

$$\boxed{x_i \in \mathbb{R}^f} \quad // \text{ item association}$$

$$\sum_{i \in N(u)} x_i \quad // \text{ sum of implicit preference}$$

Normalization  $\Rightarrow \frac{1}{|N(u)|^{0.5}} \sum_{i \in N(u)} x_i^{4.5} \quad // \text{ empirical.}$

user Attributes  $\rightarrow A(u)$  set  $\Rightarrow \boxed{y_a \in \mathbb{R}^f}$  <sup>associated</sup> factors to the attributes.  
elements.

Now overall:  $r_{ui} = \mu + b_i + b_u + q_i^T \left[ p_u + \frac{1}{|N(u)|^{0.5}} \sum_{i \in N(u)} x_i + \sum_{a \in A(u)} y_a \right] \quad \text{--- (6)}$

two extra terms.

(11)

P13

Temporal dynamics:

including time  $\hat{r}_{ui}(t) = \underbrace{\mu}_{\text{static}} + \underbrace{b_i(t)}_{\text{item bias}} + \underbrace{b_u(t)}_{\text{user bias}} + \underbrace{q_i^T p_u(t)}_{\text{static} \rightarrow \text{human behavior dynamics}}$

input with confidence level:

$$\min_{p, q, b} \sum_{(i, u) \in K} c_{ui} (\underbrace{r_{ui}}_{\text{modified confidence term}} - \mu - b_u - b_i - p_u^T q_i)^2 + \lambda (\|p_u\|^2 + \|q_i\|^2 + b_u^2 + b_i^2)$$

# (Variational Lower bound)

$z \rightarrow x$   
 $p(x) \rightarrow$  Prob dist. over variable  
 $p(x) \rightarrow$  pdf of distribution  $x$   
 Latent observed.

Posterior 
$$P(z|x) = \frac{P(x|z) P(z)}{\int_z P(x|z) P(z)}$$

Derivation 1:  $\log p(x) = \log \int_z p(x, z)$  / (-) of information  $p(x)$  value

(-) or 0 max.  $= \log \int_z p(x, z) \frac{q(z)}{q(z)} dz$

$L = \mathbb{E}_q[\log p(x|z)]$   
 $-KL[q(z)||p(z)]$  ELBO,  $L = \mathbb{E}_q[\log p(x, z)] + H(z)$  // Jensen's inequality

Reorganizing

$\mathbb{E}_q[\log p(x, z)] + H(z)$  // entropy def.

Negative  $> 1$  / is overall (-)  
 Always else  $p(x) = 0$ : no info.

Interpret: (-) of information  $>$  ELBO // reverse it (interesting)  
 No more info than  $-|ELBO|$  in  $p(x)$   
 $\Rightarrow$  more info than  $|ELBO|$

Derivation 2:  $KL(q(z)||p(z|x)) = \int q(z) \log \frac{q(z)}{p(z|x)} dz$

(Backward KL)??

fixed

$= -L + \log p(x)$  // easy  $p(z|x) = \frac{p(x, z)}{p(x)}$

// missed margin of  $q(z)$

$\therefore \log p(x) = L + KL(q(z)||p(z|x)) > 0$

to estimate  $p(z|x)$

(failed) if large gap  $L$   $p(x)$

Interpretation:

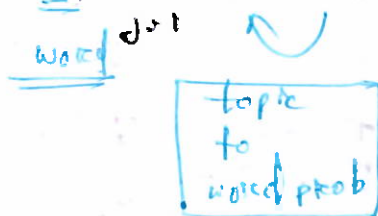
if 0 then  $L = \log p(x)$

By making elbo highest means  $q(z) \approx p(z|x)$

successful posterior estimation.

P15

$$P(w_i) = \sum_{j=1}^T P(w_i | z_i = j) P(z_i = j) \quad // \text{model itself}$$



No. document  $\rightarrow D$   
No. of topic  $\rightarrow T$

(Distribution over Dist.)

(multinomial probs are also from dist)  
↓  
Dirichlet (conjugate prior)

$D$  documents,  $T$  topics  $b$   $W$  unique words.

$$P(w | z = j) = \phi_w^{(j)} \quad // \text{multinomial distribution (T of them)}$$

$$P(z = j) = \theta_j^{(d)} \quad // \text{multinomial (document to topic)}$$

Now the objective

$$\text{Maximize } P(w | \theta, \phi)$$

modified objective for dirichlet distribution.

LDA

$$\max_{\theta, \phi} P(w | \theta, \phi) = \int P(w | \phi, \theta) P(\theta | \alpha) d\alpha \quad // \text{But intractable??}$$

Dirichlet (as conjugate prior for multinomial)  
parameter  $\alpha$   
determined by Variation Bayes / Expectation propagation.

The complete model (Gibbs sampling usage opportunity)

$$\left\{ \begin{array}{l} w_i | z_i, \phi^{z_i} \sim \text{discrete}(\phi^{z_i}) \quad // \text{from earlier} \\ \phi \sim \text{Dir}(\beta) \quad // \text{new [conjugate prior]} \\ z_i | \theta^{d_i} \sim \text{Discrete}(\theta^{d_i}) \quad // \text{earlier} \\ \theta \sim \text{Dir}(\alpha) \quad // \text{new [conjugate prior]} \end{array} \right.$$

$\alpha, \beta$  hyperparameter.



(11)

P15: for single sample: (using just count)

$$\hat{\phi}_j^{(w)} = \frac{n_j^{(w)} + \beta}{n_j^{(w)} + \beta}$$

// for new word  $w$   
and new topic  $z$ .

$$\hat{\theta}_j^{(d)} = \frac{n_j^{(d)} + \alpha}{n_j^{(d)} + T\alpha}$$

Solved by Gibbs Sampling

Alternates: Variational bayes (?), Expectation propagation (?)

Total hyperparameters:  $\alpha, \beta, T \rightarrow$  (varied across)

(fixed it in experiment)  $\downarrow$  Topic Number ??

This is model selection?

Target:  $P(w|T)$   $\uparrow$  topic Number.  
 $\downarrow$  All words  
 Approximated by  $P(w|z, T) \Rightarrow P(z|w, T)$  posterior



P16

Expected value.

$$\tilde{P}(f) = \sum_{x,y} \tilde{P}(x,y) f(x,y) \quad \text{Empirical} \quad \text{①} \quad \frac{1}{n} \times \text{no of times } (x,y) \text{ appears. (training data)}$$

calculate from data.

$$\tilde{P}(x,y) = P(f) \quad \text{indicator function}$$

feature function

we require  $\tilde{P}(f) = P(f)$

train model (training data)

Explicitly,  $\sum_{x,y} \tilde{P}(x,y) P(y|x) f(x,y) = \sum_{x,y} \tilde{P}(x,y) f(x,y)$

model training data

constraint equation

Key goal

$$\mathcal{C} = \{ P \in \mathcal{P} \mid P(f_i) = \tilde{P}(f_i) \text{ for } i = 1, \dots, n \}$$

Entropy:  $H(x) = - \sum_x P(x) \log \left( \frac{1}{P(x)} \right)$

Fig 1

Conditional Entropy,  $H(P) = - \sum_{x,y} \tilde{P}(x) P(y|x) \log P(y|x) \geq 0$

All the 2 possible data prob

model

lower bound

sure model

(set of probabilities distribution)  $H(P)$

Maximum Entropy:  $P_* \in \mathcal{C}$  satisfy  $\tilde{P}(f) = P(f)$

train model

uniform case

upper bound.

uniform case entropy.

log 1/y

cardinality of y

well defined & unique.

$\Rightarrow$  To select a model from a set of prob. distributions, that has maximum entropy.

Parametric Form:

find  $P_*$  -  $\arg \max_{P \in \mathcal{C}} H(P)$  primal optimization.

Lagrangian  $\Lambda(P, \lambda) = H(P) + \sum_i \lambda_i (P(f_i) - \tilde{P}(f_i))$

training (ground truth)

Now,  $P_\lambda = \arg \max_{P \in \mathcal{P}} \Lambda(P, \lambda)$

(model)

$\Psi = \Lambda(P_\lambda, \lambda)$  // max value // dual function.



P 15

solving,  $\underline{p_\lambda(y|x)} = \frac{1}{Z_\lambda(x)} \exp \left( \sum_i \lambda_i f_i(x, y) \right)$

normalization w.r.t.  $y$

unknown (solved later)  $\lambda^*$

model

$$\Psi(\lambda) = - \sum_x \tilde{p}(x) \log Z_\lambda(x) + \sum_i \lambda_i \tilde{p}(f_i)$$

Dual optimization find  $\lambda^* = \underset{\lambda}{\operatorname{argmax}} \Psi(\lambda)$  // maximize

$\lambda^* \rightarrow$  parametric form

Relation to maximum likelihood: Training data.

$$L_P(P) = \log \prod_{x,y} \underbrace{P(y|x)}_{\text{model}}^{\tilde{p}(x,y)} = \sum_{x,y} \tilde{p}(x,y) \log P(y|x)$$

By definition,  $\Psi(\lambda) = L_P(P_\lambda)$  // see earlier section.

$P^* \in C$  with maximized entropy is parametric model from  $\underline{P_\lambda(y|x)}$  family that maximize the likelihood of training sample  $\tilde{p}$

Table 1 Summary

Compute the Params: ①  $f_i(x, y) \geq 0$

Algo: Input Iterative Scaling.

input:  $f_1, \dots, f_n \rightarrow$  empirical  $\tilde{p}(x, y)$

output:  $\lambda_i^*, P_\lambda$

1.  $\lambda_i = 0 (i = 1 \dots n)$

2.  $i \in \{1 \dots n\}$

See the algorithm

①

P19

$x_i \in \mathbb{R}^n$  ;  $i = 1, \dots, L$  observations.  $\rightarrow$  ex.  $y_i \in \{1, -1\}$

$f(x, \alpha) \rightarrow$  approximation

Expectation of test error.

parameter

$$R(\alpha) = \int \frac{1}{2} |y - f(x, \alpha)| dP(x, y) \quad // \text{true}$$

$$R_{\text{emp}}(\alpha) = \frac{1}{2L} \sum_{i=1}^L |y_i - f(x_i, \alpha)| \quad // \text{Approximation}$$

vc confidence

Connection between them

Vapnik, 1995

$$R(\alpha) \leq R_{\text{emp}}(\alpha) + \sqrt{\left( \frac{h(\log(2L/h) + 1) - \log(n/4)}{L} \right)}$$

Here,  $0 \leq \eta \leq 1$  ; with prob  $(1 - \eta)$  holds.

$L \rightarrow$  example

$h \rightarrow$  vc dimension.

$n \uparrow$   
 $h \downarrow$  we want  
 $L \uparrow$

set to minimize this bound

[vc confidence lower  $\rightarrow$  the better] (may overfit)

The vc dimension: increases function capacity  $\uparrow$  confidence boundary (higher is Bad)

$\infty$  infinite vc dimension  $f(x, \alpha) = \theta(\sin(\alpha x))$ ,  $x, \alpha \in \mathbb{R}$

$$x_i = w^{-i}$$

$y_i$  -- Assign anything.

$$\alpha = \pi \left( 1 + \sum_{i=1}^L \frac{(1 - y_i) w^i}{2} \right)$$

true Approximate

vc dimension  $\neq \infty$

shattering depends on choice of points.

choose points that can be shattered.

## Separable Case

① satisfying hyperplane:  $\underline{w} \cdot \underline{x} + b = 0$  (HP)  
 & Linear  $\nearrow$   
 normal to HP.  
 projection to HP =  $-\underline{b}$   
 bias  
 (All points projected as  $-b$  amount)  
 (nice)

Separable cases  
 $\left\{ \begin{array}{l} x_i \cdot \underline{w} + b \geq 1 ; \text{ overshoot in projection } y_i = +1 \\ x_i \cdot \underline{w} + b \leq -1 ; y_i = -1 \end{array} \right.$

$$\Rightarrow y_i (x_i \cdot \underline{w} + b) - 1 \geq 0 \quad \forall$$

introducing lagrangian  $\alpha_i, i=1 \dots d$ ;  $\alpha_i \geq 0$   
 multiply & add all.

$$L_P = \frac{1}{2} \|\underline{w}\|^2 - \sum_{i=1}^d \alpha_i y_i (x_i \cdot \underline{w} + b) + \sum_{i=1}^d \alpha_i$$

primal problem.

$$\left\{ \begin{array}{l} \underline{w} = \sum_i \alpha_i y_i x_i \quad // \text{ solution.} \\ \sum_i \alpha_i y_i = 0 \end{array} \right.$$

Now the dual problem.

putting values  $\hookrightarrow$

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\underline{x}_i \cdot \underline{x}_j)$$

// dual formulation  
kernel idea.

use  $K(\underline{x}_i, \underline{x}_j)$  in nonlinear case

KKT condition:

$$\frac{\partial}{\partial \underline{w}_0} L_P = \underline{w}_0 - \sum_i \alpha_i y_i x_i = 0 \quad ; i=1 \dots d$$

$$\frac{\partial}{\partial b} L_P = y_i (\underline{x}_i \cdot \underline{w} + b) - 1 \geq 0 \quad ; \sum_i \alpha_i y_i = 0$$

$$y_i (\underline{w} \cdot \underline{x}_i + b) - 1 \geq 0 \quad i=1 \dots d$$

$$\alpha_i \geq 0 \quad \forall i$$

$$\alpha_i (y_i (\underline{w} \cdot \underline{x}_i + b) - 1) = 0 \quad \forall i$$

(if separable)

P19 Non linear SVM:

$\phi: \mathbb{R}^d \rightarrow \mathcal{H}$  Higher dimension Projection.

$$K(\underline{x}_i, \underline{x}_j) = \phi(\underline{x}_i) \cdot \phi(\underline{x}_j)$$

need  $\mathcal{H}$  of  $|\mathcal{H}|$  dimensional??

How to use the kernel??

→ we just need dot product.

test phase →  $f(x) = \sum_{i=1}^{N_s} \alpha_i y_i \phi(\underline{s}_i) \cdot \phi(x) = \sum_{i=1}^{N_s} \alpha_i y_i \underbrace{K(\underline{s}_i, x)}_{\text{only need this}} + b$

using the train data  
No of support vector

mercer's Condition:-  $(\mathcal{H}, \phi)(d \rightarrow \mathcal{H})$

$\checkmark K(\underline{x}, \underline{y}) = \sum_i \phi(\underline{x})_i \phi(\underline{y})_i$  mapping exists

if  $\forall g(x) \geq 0$  that satisfy.

$\int g(x)^2 dx$  is finite.

then

$\int K(\underline{x}, \underline{y}) g(\underline{x}) g(\underline{y}) d\underline{x} d\underline{y} \geq 0 \Rightarrow \text{Positive Semidefinite (PSD)}$

open Question: How to formulate  $\phi$ ?

since VC dimension is  $|\mathcal{H}| + 1$  // in this case

so  $|\mathcal{H}| \uparrow$  bad generalize.

Radial basis kernel:

$K(\underline{x}, \underline{y}) = e^{-\|\underline{x} - \underline{y}\|^2 / 2\sigma^2}$  // maybe infinite VC Dimension.

→ two layer sigmoid NN.

(generalize  
n vs 1 classifier)

1. Layer I  $N_s$  weights each  $d_1$  dimensional

2. Layer II  $N_s$  weights ( $\alpha_i$ )

finally Sigmoid.

(iii)

F19

Nonseparable Case:

may cause

error

$$\underline{x}_i^T \underline{w} + b \geq 1 - \underline{\epsilon}_i \quad ; \quad y_i = +1$$

$$\underline{x}_i^T \underline{w} + b \leq -1 + \underline{\epsilon}_i \quad ; \quad y_i = -1$$

$$\epsilon_i \geq 0 \quad \forall i$$

if any  $\epsilon_i > 1$  error occurs.

$\sum_i \epsilon_i$  = upper bound of training error.

Dual problem)  $L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \underline{x}_i^T \underline{x}_j$

Subject to:  $0 \leq \alpha_i \leq C$  → use parameter.  
 $\sum \alpha_i y_i = 0$  (higher penalty to error)

Solution is  $\underline{w} = \sum_{i=1}^{N_S} \alpha_i y_i \underline{x}_i$

skip gram model maximize  $\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-c \leq j \leq c \\ j \neq 0}} \log P(w_{t+j} | w_t)$   
 sequence of words  $\{w_1, w_2, \dots, w_T\}$

So, maximizing  $P(w_{t+j} | w_t)$  is

defined as  $P(w_o | w_I) = \frac{\exp(v_{w_o}^T v_{w_I})}{\sum_w \exp(v_w^T v_{w_I})}$  
 $\begin{cases} c \uparrow \\ \text{training time} \uparrow \\ \text{Accuracy} \uparrow \end{cases}$

$\begin{cases} v_{w_o} = \text{output vector reps.} \\ v_{w_o} = \text{input vector reps} \end{cases}$   $w=1$   $\rightarrow$  All the words ?? [huge computation]  
 $v_{w_o}$  - vector representation of  $w_o$  (via Network)

$w \rightarrow$  huge size !!

Each word has two reps  
 $\rightarrow$  input  $v_w$   
 $\rightarrow$  output  $v_w'$

Hierarchical Softmax: Need  $\log w$  nodes.

$$P(w | w_I) = \prod_{j=1}^{L(w)-1} \sigma \left( \left[ n(w_j, d) = \begin{cases} 1 & \text{if true, else 0} \end{cases} \right] v_{n(w_j, d)}^T v_{w_I} \right)$$

$\rightarrow$  care about input/output representation. [computation  $\propto L(w)$ ]

Negative Sampling should be high (interesting)

NEG objective:  $\log \sigma(v_{w_o}^T v_{w_I}) + \sum_{i=1}^K E_{w_i \sim P_n(w)} [\log \sigma(-v_{w_i}^T v_{w_I})]$

$\Downarrow$  row of matrix positive  
 $\uparrow$  different choices  
 $\rightarrow$  negative word how (neg)  
 $\Downarrow$  should be low & negative

$$\sigma(x) = \frac{1}{1 + \exp(-x)}$$

// modified NCE



subsampling of freq words.

Discarding prob

$$P(w_i) = 1 - \sqrt{\frac{t}{f(w_i)}}$$

word frequency.  
// Discarding prob.

✓  $f(w_i) \uparrow$   $P(w_i) \uparrow$

✓  $f(w_i) \downarrow$   $P(w_i) \downarrow$   
(High freq words discarded more)

Balance between rare & frequent words.

Bigram skip:

(Phrase selection)

score( $w_i, w_j$ ) =

$$\frac{\text{count}(w_i, w_j) - \delta}{\underbrace{\text{count}(w_i)}_{\text{unigram } w_i} \times \underbrace{\text{count}(w_j)}_{\text{unigram } w_j}}$$

bigram



①

P21

The probabilistic model:

$$P(\{s_t, y_t\}) = P(s_1) P(y_1 | s_1) \prod_{t=2}^T P(s_t | s_{t-1}) P(y_t | s_t)$$

$\downarrow$  observed, D  
 $\downarrow$  Hidden state, k

$\underbrace{\prod_{t=2}^T P(s_t | s_{t-1}) P(y_t | s_t)}_{\text{separable. Conditional Independence.}}$

Let, k states,

 $P(s_t | s_{t-1}) \Rightarrow k \times k$  matrix (Transition matrix) <sup>state</sup>
 $P(y_t | s_t) \Rightarrow k \times D$  observation matrix

 $\rightarrow$  modeled by GMM/Neural Networks.

what if:  $s_t = s_t^{(1)}, s_t^{(2)}, \dots, s_t^{(m)}$  } factorial HMM !!

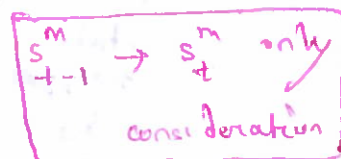
$\downarrow$   $\downarrow$   $\downarrow$   
 $k^{(m)}$  possible values. each

 $k^{(m)} = k$  // simplicity
Then,  $k^m$  states  $k^m \times k^m$  state transition matrix!!
 $\rightarrow$  impossible to work with

 $\rightarrow$  Requires constraint on state tx mat.
factorial HMM  $\rightarrow$  underlying state tx is constrained

$$P(s_t | s_{t-1}) = \prod_{m=1}^M P(s_t^{(m)} | s_{t-1}^{(m)})$$

Decoupled.



p21

observation:

→  $D \times D$  variance matrix

$$P(Y_t | s_t) = |C|^{1/2} (2\pi)^{-D/2} \exp \left\{ -\frac{1}{2} (Y_t - \mu_t)' C^{-1} (Y_t - \mu_t) \right\}$$

for all  $m$ th contribution

$$\mu_t = \sum_{m=1}^M w_t^{(m)} s_t^{(m)}$$

probabilistic

columns are contribution of each states →  $D \times 1$  finally.

state variable  $s_t^{(m)} \Rightarrow K \times 1$  vectors. → only one 1 (one hot encoding)

→ depends on  $s_t^{(m)}$  value

Learning & Inference:Expectation maximization: param learning

$$Q(\phi^{new} | \phi) = E \left\{ \log P(s_t, Y_t | \phi^{new}) \mid \phi, \{Y_t\} \right\} \quad (6)$$

↓  
current param  
new params.

$$p_t^{(m)} = P(s_t^{(m)} | s_{t-1}^{(m)})$$

Factorial HMM:  $\phi = \{ \underline{w}^{(m)}, \underline{\pi}^{(m)}, p^{(m)}, \underline{c} \}$  // find all of these parameters.

↓  
 $p(s_t^{(m)})$

E step:compute  $Q$ : Expand  $s$  by using forward equations.

→ can be expressed as Expectation of

$$E \{ \cdot | \phi, Y_t \} \Rightarrow \underbrace{\langle s_t^{(m)} \rangle}_{\substack{\text{state} \\ \text{occupation} \\ Y_t}} ; \underbrace{\langle s_t^{(m)} \cdot s_t^{(n)} \rangle}_{\substack{K \times 1 \text{ vec.} \\ \text{two states jointly}}} ; \underbrace{\langle s_{t-1}^{(m)} s_t^{(m)} \rangle}_{\substack{\text{state} \\ \text{transition} \\ \sum_t K \times K \text{ mat}}}$$

M step: maximize  $Q$  using Jensen's inequality.

solved by: weighted linear regression.

Gibbs Sampling: Inference:

$$s_t^{(m)} \sim p(s_t^{(m)} | \underbrace{\{s_t^{(n)} : n \neq m\}}_{\text{Neigh}}, \underbrace{s_{t-1}^{(m)}}_{\text{past}}, \underbrace{s_{t+1}^{(m)}}_{\text{future}}, \underbrace{y_t}_{\text{observation}})$$

$$\propto p(s_t^{(m)} | s_{t-1}^{(m)}) p(s_{t+1}^{(m)} | s_t^{(m)}) p(y_t | s_t^{(m)}, \dots, s_t^{(m)}, \dots, s_t^{(m)})$$

↑ state transition
↑ state transition
graphical model design itself

↓ markovian

completely factorized Variational Inference:

$$\begin{aligned} \textcircled{1} \quad \log p(\{y_t\}) &= \log \sum_{\{s_t\}} p(\{s_t, y_t\}) \\ &= \log \sum_{\{s_t\}} q(\{s_t\}) \frac{p(\{s_t, y_t\})}{q(\{s_t\})} \\ &\geq \sum_{\{s_t\}} q(\{s_t\}) \log \left[ \frac{p(\{s_t, y_t\})}{q(\{s_t\})} \right] \end{aligned}$$

②

The difference between ① & ② is  $|\textcircled{1} - \textcircled{2}|$  // simple math.

$$KL(q||p) = \sum_{\{s_t\}} q(\{s_t\}) \log \left[ \frac{q(\{s_t\})}{p(\{s_t | y_t\})} \right]$$

→ change parameters of  $q(\{s_t\})$  to minimize:

p > 1

(iv)

$$Q(\{s_t | \theta\}) = \prod_{t=1}^T \prod_{m=1}^M Q(s_t^{(m)} | \theta_t^{(m)})$$

time step
possible steps at each t.

vector itself.  
 $\theta_t^{(m)} = \begin{bmatrix} \theta_{t,1}^{(m)} \\ \theta_{t,2}^{(m)} \\ \vdots \end{bmatrix}$

$$Q(s_t^{(m)} | \theta_t^{(m)}) = \prod_{k=1}^K \left( \theta_{t,k}^{(m)} \right)^{s_{t,k}^{(m)}}; \quad s_{t,k}^{(m)} \in \{0, 1\}$$

multiply. All the params. with math markovian chain state  $k$ , at time  $t$

$\sum_{k=1}^K s_{t,k}^{(m)} = 1$   
 only one is 1 else 0

$s_t^{(m)} = \begin{bmatrix} s_{t,1}^{(m)} \\ s_{t,2}^{(m)} \\ \vdots \end{bmatrix}$

$\theta_t^{(m)}$  → state occupation prob. with multinomial over  $s_t^{(m)}$   
under distribution  $Q$

vector softmax elementwise

$$\theta_t^{(m) \text{ New}} = \phi \left\{ W^{(m)'} C^{-1} \tilde{y}_t^{(m)} - \frac{1}{2} \Delta^{(m)} + (\log P^{(m)}) \theta_{t-1}^{(m)} + (\log P^{(m)})' \theta_{t-1}^{(m)} \right\}$$

residual error  $\tilde{y}_t^{(m)} = y_t - \sum_{l \neq m} W^{(l)} \theta_t^{(l)}$

vector of diagonal elements  $W^{(m)'} C^{-1} W^{(m)}$

Structured Variational Inference:

$$Q(\{s_t\} | \theta) = \frac{1}{Z_Q} \prod_{m=1}^M Q(s_1^{(m)} | \theta) \prod_{t=1}^T Q(s_t^{(m)} | s_{t-1}^{(m)}, \theta)$$

$Z_Q$  → normalized

$$Q(s_1^{(m)} | \theta) = \prod_{k=1}^K \left( h_{1,k}^{(m)} \pi_k^{(m)} \right)^{s_{1,k}^{(m)}}$$

$$Q(s_t^{(m)} | s_{t-1}^{(m)}, \theta) = \prod_{k=1}^K \left( h_{t,k}^{(m)} \sum_{j=1}^K P_{k,j}^{(m)} s_{t-1,j}^{(m)} \right)^{s_{t,k}^{(m)}}$$

P21

(v)

$$Q(s_t^{(m)} | s_{t-1}^{(m)}, \theta) = \prod_{k=1}^K \left( h_{t,k}^{(m)} \prod_{j=1}^K (p_{kj}^{(m)})^{s_{t-1,j}^{(m)}} \right)$$

one hot vectors.

$$\theta = \{ \pi^{(m)}, p^{(m)}, h_t^{(m)} \}$$

$K \times 1 \rightarrow$  prob of observation  $P(y_t | s_t)$

for each  $K$  setting  $s_t^{(m)}$

$$Q(s_{t,j}^{(m)} = 1 | \theta) = h_{t,j}^{(m)} P(s_{t,j}^{(m)} = 1 | \phi)$$

$\Rightarrow$  having an observation at  $t=1$ , under  $s_{t,j}^{(m)} = 1$   
has prob of  $h_{t,j}^{(m)}$

Can be proved that,  $KL(Q||P)$  is minimized.

$$h_t^{(m) \text{ new}} = \exp \left\{ w^{(m)'} e^{-1} \tilde{y}_t^{(m)} - \frac{1}{2} \Delta^{(m)} \right\}$$

recs  $\rightarrow y_t = \sum_{k \neq m} w^{(k)} \langle s_t^{(k)} \rangle$

connected to ELBO bound

$$F(q, \phi) = \mathbb{E}_Q \left\{ \log P(y, s | \phi) \right\} - \mathbb{E}_Q \log \{ Q(s) \} \leq \log P(y)$$

(-) value.  $\rightarrow$  reverse it (interesting)  
No more info than ELBO

①

P22

Problem formulation:

$\{A_1, A_2, \dots, A_m\} \rightarrow m$  smart phone.

$A_i = \{A_{i1}, \dots, A_{iL}\}$  // complete trace for  $A$ .

$A_{ij} = \{t_{ij}, x_{ij}, y_{ij}\}$  // temporal/spatial information.

Query  $Q = \{Q_1, \dots, Q_f\}$ ;  $f \ll L$  time.

Targets  $k$  relevant trajectories of  $Q$  from  $A$  site

Trajectory comparison function  $LESS(Q, A_i)$   
 compare their trajectory.

Longest Common SubSequence (LCSS)

By definition:

$$LESS_{\delta, \epsilon}(A, B) = \begin{cases} 0, & A \text{ OR } B = \emptyset \\ 1 + LESS_{\delta, \epsilon}(\text{Head}(A), \text{Head}(B)) & \text{if: } |a_{x:L_1} - b_{x:L_2}| < \epsilon \text{ (time)} \\ & |a_{y:L_1} - b_{y:L_2}| < \epsilon \text{ (x coordinate)} \\ & |L_1 - L_2| < \delta \text{ (time, y coordinate)} \\ \max(LESS_{\delta, \epsilon}(\text{Head}(A), B), LESS_{\delta, \epsilon}(A, \text{Head}(B))) & \text{otherwise} \end{cases}$$

time matching window  $\delta, \epsilon$   
 spatial matching window  
 Both are application Specific.

[iterative Algorithm]

$(x, y \text{ at time } t)$

$\text{Head}(A) = ((a_{x:1}, a_{y:1}), \dots, (a_{x:L-1}, a_{y:L-1}))$

P22

Bounding above LCSS: easier computation.

$$LCSS(MBE_Q, A_i) = \sum_{j=1}^{|A_i|} \begin{cases} 1, & \text{if } A_i[j] \text{ within envelop.} \\ 0, & \text{otherwise.} \end{cases}$$

$MBE_Q$ : Minimum Bounding Envelop of Query Q.

$MBE_Q$  is the area between high envelop &  $EnvHigh[i]$   
Low envelop.  $EnvLow[i]$

$$EnvHigh[i] = \max(Q[j] + \epsilon); |i-j| \leq \delta$$

$$EnvLow[i] = \min(Q[j] - \epsilon); |i-j| \leq \delta$$



unique solution

$\left\{ \begin{array}{l} G \rightarrow \text{recovers the training data distribution} \\ D \rightarrow \frac{1}{2} \text{ everywhere} \end{array} \right.$

Adversarial Nets:

gen  $\sim p_g$

noise vector  $z \rightarrow \text{map } G(z, \theta)$

$\Downarrow$

differentiable function (MLP) multi layer perceptron

$D(x, \theta_D)$  / differentiable MLP.

$\rightarrow$  Discrimination  $\rightarrow x$  from data /  $G$  ??

$$\min_{G_D} \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left\{ \log[D(x)] \right\} + \mathbb{E}_{z \sim p_z(z)} \left\{ \log[1 - D(G(z))] \right\}$$

$\Rightarrow$  iterative numerical approach:

$\left\{ \begin{array}{l} k \text{ step of } D \rightarrow \text{to keep near optimal (inner loop)} \\ 1 \text{ step of } G \rightarrow \text{changes slowly enough.} \end{array} \right. \quad \left\{ \begin{array}{l} \text{ML/PCD way??} \end{array} \right.$

Theory: Algorithm 1  $\Rightarrow$  crack of jack.

understanding sequence: when  $D$  is optimal  $\Rightarrow$  ??  $\checkmark$   
 How/when  $G$  is optimal |  $D$  is optimal  $\checkmark$   
 what happens when  $D$  is optimal ??  $\checkmark$   
 what happens when both are optimal

for no of train

mnest  
loop  
keep k  
near  
optimal

for  $k$  in range ( $k$ )sample  $z^1, \dots, z^m \rightarrow P_g(z)$ sample  $x^1, \dots, x^m \rightarrow P_{data}$ update  $\theta_D$  by ascending ~~grad~~ stochastic grad.

$$\nabla_{\theta_D} \frac{1}{m} \sum_{i=1}^m [\log D(x_i) + \log (1 - D(G(z_i)))] \quad // \text{maximize}$$

→ end. for

sample  $m$  noise  $\{z^1, \dots, z^m\}$ Update the gen.  $\theta_g$  descending gradient

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z_i))) \quad // \text{minimize.}$$

end. for.

global optimality:  $P_g = P_{data}$ .loop 1: For  $G$  fixed,

$$D_{opt}(x) = \frac{P_{data}(x)}{P_{data}(x) + P_g(x)} \quad // \text{optimal Discriminator (maximize)}$$

training criterion:

$$V(G, D) = \int_x P_{data}(x) \log(D(x)) dx + \int_z P_g(z) \log(1 - D(G(z))) dz$$

$$= \int_x P_{data}(x) \log(D(x)) dx + \int_z P_g(z) \log(1 - D(G(z))) dz$$

for any function  $y \rightarrow a \log y + b \log(1-y)$  achieves its  
maximum in  $[0, 1]$  at  $\frac{a}{a+b} = y$

$$\underline{C(G)} = \max_D V(G, D)$$

$$= E_{x \sim P_{\text{data}}} [\log D_G^*(x)] + E_{z \sim P_z} [\log (1 - D_G^*(G(z)))]$$

$$= E_{x \sim P_{\text{data}}} [\log D_G^*(x)] + E_{x \sim P_g} [\log (1 - D_G^*(x))] \quad // \text{from earlier argument } (y = \frac{a}{a+b}) \text{ maximized}$$

$$= E_{x \sim P_{\text{data}}} \left[ \log \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \right] + E_{x \sim P_g} \left[ \log \frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)} \right]$$

Now D is set to tune  $P_g(x) / G(x)$  to minimize this — (1)

The global minima for  $C(G)$  is if  $P_g(x) = P_{\text{data}}(x)$

Now D fails to comprehend

in that case,

$$E_{x \sim P_{\text{data}}} \left[ \log \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \right] + E_{x \sim P_g} \left[ \log \frac{P_g(x)}{P_{\text{data}}(x) + P_g(x)} \right] =$$

$$= -\log 2 - \log 2 = -2 \log 2$$

Re organizing the equation (1) we get  $P_g$  may not be optimal

$$C(G) = -\log 4 + E_{x \sim P_{\text{data}}} \left[ \log \frac{2 P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \right] + E_{x \sim P_g} \left[ \log \frac{P_g(x)}{\frac{P_{\text{data}}(x) + P_g(x)}{2}} \right]$$

$$= -\log 4 + \underbrace{KL \left( P_{\text{data}} \parallel \frac{P_{\text{data}} + P_g}{2} \right) + KL \left( P_g \parallel \frac{P_{\text{data}} + P_g}{2} \right)}$$

$$= -\log 4 + 2 JS(P_{\text{data}} \parallel P_g) \quad // \text{Jensen-Shannon Divergence}$$

for optimal  $G^* \Rightarrow C(G) = -\log 4$  !! if  $P_g$  is optimal if only  
minimum for optimal D ( $P_g = P_{\text{data}}$ )

Because  $G$  is too good  
 Not because D is bad  
 D is still OPTIMAL

P 29

convergence of Algo: <sup>(1)</sup> 1st discrimination reaches optimal <sup>(2)</sup>  $D^*(x) = P_{data} / (P_g + P_{data})$  and  $P_g$  is updated to improve.

<sup>(1)</sup> 
$$E_{x \sim P_{data}} [\log D_g^*(x)] + E_{x \sim P_g} [\log (1 - D_g^*(x))]$$

<sup>(2)</sup>  $P_g$  converge to  $P_{data}$ .  
 $\Downarrow$   
 MLP via function  $g(z; \theta_g)$  <sup>optimize this instead of  $P_g$</sup>   
 (multiple critical point) <sup>(multiple local optima)</sup>  
 theory  $\rightarrow$  (not yet)

D vs Ad:  $G$  must not be trained too much!! (D must be mode collapse!! convergence??) (sync with  $G$ )

Ad: only backpropagation  
 Large distribution learning.

$\rightarrow$  if  $D$  is too strong to learn nothing??  $\leftarrow$

Joint dist  $P(x, h | \eta) = P(x | h) P(h | \eta)$

$\uparrow$  hidden  
 $\uparrow$  obs  $\downarrow$  param

$$P(h | x, \eta) \propto P(h, x | \eta)$$

predictive,  $P(x_{\text{new}} | x) = \int P(x_{\text{new}} | h) P(h | x, \eta) dh$

mixture (dir)

$$P(\mu_{1:k}, \theta, z_{1:n} | x_{1:n})$$

$\downarrow$  mean  $\downarrow$  class  $\downarrow$  observ

Generative probabi. mod:

mixture prior  $\left\{ \begin{array}{l} \theta \sim \text{Dirichlet}(\alpha) \\ \mu_k \sim N(0, \sigma_0^2) \end{array} \right.$

global hid var

① mixture assignment  $z_n | \theta \sim \text{Dirichlet}(\theta)$  ?

② Data point  $x_n | z_n, \mu \sim N(\mu_{z_n}, 1)$

easy posy:- sample  $\theta \rightarrow$  corresponding  $\mu_{1:k}$  for each data point estimate  $z_n \leftarrow 1 \dots k \rightarrow$  given the distribution.

Paper 30

②

joint distribution:

$$p(\theta, \mu, z, x | \sigma_0^2, \alpha) = p(\theta | \alpha) \prod_{k=1}^K p(\mu_k | \sigma_0^2) \prod_{i=1}^N p(z_i | \theta) p(x_i | z_i, \mu)$$

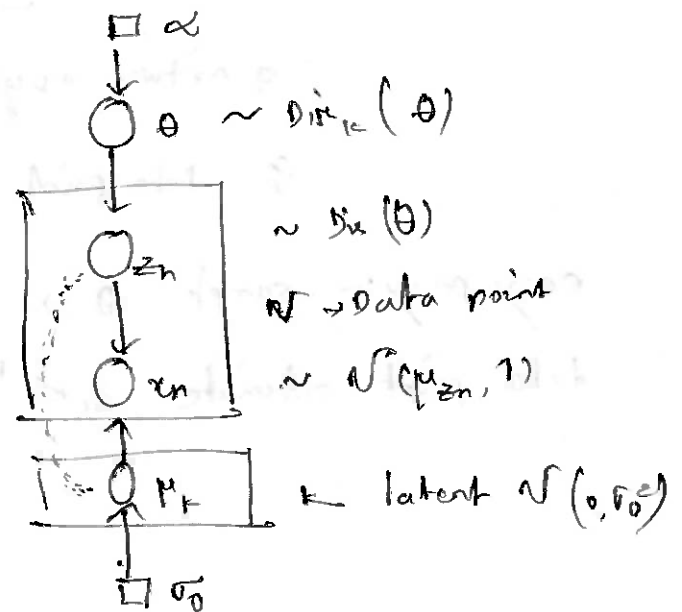
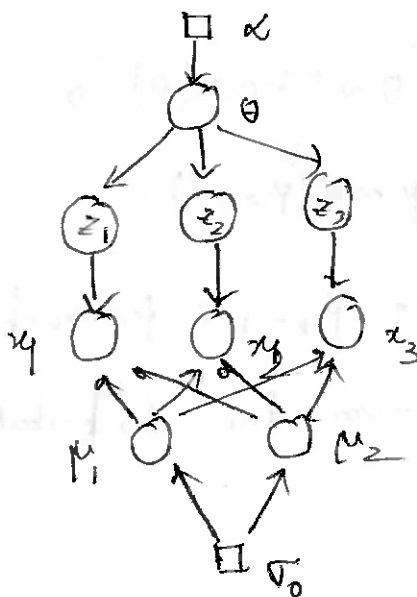
Posterior

$$p(\theta, \mu, z | x, \sigma_0^2, \alpha) = \frac{p(\theta, \mu, z, x | \sigma_0^2, \alpha)}{\int p(x | \sigma_0^2, \alpha)}$$

Predictive Distribution.

$$p(x_{\text{new}} | x, \sigma_0^2, \alpha) = \int \left( \sum_{z_{\text{new}}} p(z_{\text{new}} | \theta) p(x_{\text{new}} | z_{\text{new}}, \mu, \sigma_0^2) p(\theta, \mu | x, \sigma_0^2, \alpha) \right) d\theta \cdot d\mu$$

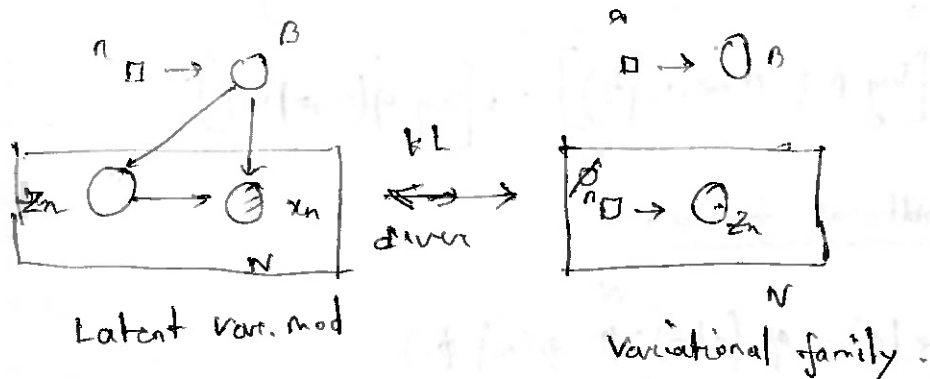
The graphical model:



### Example model:

- (i) linear factor model: PCA, factor model. (graph 3)
- (ii) mixed membership model:
- (iii) Matrix factorization model:
- (iv) Time Series
  - ~~Time Series models~~ Hidden Markov model
  - Kalman filter

### Posterior Inference      With mean field : Variational method.



### Conditional Conjugate model:

$$p(\beta, z, x | \eta) = p(\beta | \eta) \prod_{n=1}^N p(z_n | \beta) p(x_n | z_n, \beta)$$

local latent  $\uparrow$   
 global latent  $\downarrow$       obs  $\downarrow$   
 Always fixed mixture prop  $\uparrow$   
 changes so local  $\downarrow$

Posterior:

$$p(\beta, z | x) = \frac{p(\beta, z, x)}{\int p(\beta, z, x) d\beta, dz} \rightarrow \text{now? problem!}$$