(1)

P7

$$\left\{\begin{array}{l}\text{score metric} \\ \text{for } (x,y) \text{ word} \\ \text{pair}\end{array}\right. \quad \underbrace{S_{(a,b)}(x,y)}_{\text{baseline}} = \begin{cases} \cos(\vec{a}-\vec{b}, \vec{x}-\vec{y}) & \text{if } \|\vec{x}-\vec{y}\| \le \delta \\ 0 \end{cases}$$

$\leftarrow$ 7th experime

→ seed direction (she, he) $\downarrow$ similarity th

-Analogy

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

### Direct gender bias

Gender neutral word

strictness

Definition $\boxed{\text{Direct-bias}_c = \frac{1}{|N|} \sum_{w \in N, \text{set}} |\cos(\vec{w}, g)|^c}$

?? confused here. projection

→ [gender subspace top principle component]

Indirect Bias Definition
(Described in next page)

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

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

if 0 means no projection.
$\Downarrow$
$\beta(w,v)=0$ orthogonal to each other

if , $w_g = w \Rightarrow w_\perp = 0$

then $\Rightarrow \beta(w,v) = 1$

unit vector

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

original vector

### Debiasing Algorithm:

Projection Direction $v_B = \sum_{d=1}^{k} (v \cdot b_j) b_j$

S value $= v^T_B B$ //orthogonal project: $v - v_B$

step 1: Identify gender subspace:

Defining sets. total words
$D_1 \dots D_n \subseteq W$

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

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

Let Bias subspace

$$C := \sum_{i=1}^{b} \sum_{w \in D_i} (\vec{w}-\mu_i)^T(\vec{w}-\mu_i)/|D_i| \quad \begin{bmatrix}\text{k row of SVD} \\ \text{is bias subspace } \textcircled{B}\end{bmatrix}$$

P7

Hard de-biasing : $\underset{\text{new embedding}}{\overline{w}} := (\overline{w} - \overline{w}_B)/\|\overline{w} - \overline{w}_B\|$ // re-embedding definition.
← orthogonal projection

∴ word to Neutralize $N \subseteq W$ ↑ → got B matrix/earlier ✓
here $\overset{2.2}{\longrightarrow}$ step → projection

family/Equality set $E = \{E_1', E_2', \dots, E_m'\}$ //?? what we want equidist.

$E_i' \subseteq W$ → finally all words will have similar
component in gender neutral Direction

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

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

this term varies only for words in E set

For $\forall \; w \in E; \quad \overline{w} := \overline{v} + \sqrt{1 - \|v\|^2} \dfrac{\overline{w}_B - \mu_B}{\|\overline{w}_B - \mu_B\|}$

Added for the bias component differences.

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

$$\left( v_B, \quad w_{\perp B} = w - w_B \right)$$

___

s-oft bias Correction : $W \in \mathbb{R}^{d \times |vocab|}$ new

$T \to$ transformation $\underline{d \times d}$

$$\min_{T} \| (T W)^T (T W) - \overline{w}^T w \|_F^2 + \lambda \| (T N)^T (T B) \|_F^2$$

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, $\beta(\overline{w}, \overline{v}) = \dfrac{\overline{w}^T \overline{v} - \dfrac{w_{\perp B}^T v_{\perp B}}{\|w_{\perp B}\| \|v_{\perp B}\|}}{\overline{w}^T \overline{v}}$ ⟹ [match in gender Independent direction]

overall match.

P3

simplified version of pagerank:

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

→ Number of links from u

Back link of u (towards u)

→ Normalization.

→ Recursive equation.

A matrix → page × page big matrix.

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

eigen value

0 if no edge between u, v

A vector! (bad notation)

$$R = c \cdot A \cdot R$$

→ eigen vector of A

→ Eigen vector of A → symmetric matrix.

→ All vectors are orthogonal.

→ Dominant eigenvector

→ 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 of rank.
vector.

L1 Norm, $\|R'\|_1 = 1$, c is maximized.

if $E(u) \geq 0$ the c is reduced.

Decay factor.    [All 1 → matrix]

⇒ $R' = c(AR' + E) = c(A + E \times 1) R'$  // since $\|R'\|_1 = 1$

Eigen value of this one.

P ?

$$R_0 \leftarrow s \quad // \text{ Random ints.}$$

**loop:**

finding Eig vector for $\underline{A}$ (Dominant)

$\rightarrow R_{i+1} \leftarrow A R_i \quad // \text{ power iteration (PE)}$

$\rightarrow d \leftarrow \|R_i\|_1 - \|R_{i+1}\|_1 \quad // \text{ constrained }.$

$\rightarrow R_{i+1} \leftarrow R_{i+1} + \boxed{d} E \quad // \text{ little move from PE}$

$\rightarrow \delta \leftarrow \|R_{i+1} - R_i\|_1 \longrightarrow$ increases convergence.

maintain $\|R\|$,

$// \text{Normalize}$,

**while** $\delta > e \quad // \text{convergence}$

P12: Recommendation problem formulation.

$C \to$ set of users (user space, $\to$ name, age, demograph,...)

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

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

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

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

## Content based filtering methods:

focus on $\boxed{u(c, \widehat{s_i})}$ · user already has rated.

$\quad\quad\quad \hookrightarrow$ similar to previous $s_i$ will be recommended

Term frequency $\to TF = \dfrac{f_{i,j}}{\underset{z}{Max}\ f_{z,j}}$ ; $f_{i,j} \triangleq$ no time $k_i$ word appear in document, $d_j$

inverse Document frequency $\to IDF_i = \log \dfrac{N}{n_i}$ ; $N \to$ total documents.

$\quad\quad\quad\quad n_i \to \textcircled{$k_i$}$ appeared in how many documents.

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

$$w_{i,j} = TF_{i,j} \times IDF_i$$

for content of document $d_j$ \qquad for all the key words $k_s$.

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

P12

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

→ for user. $c \in C$

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

$$\Rightarrow \boxed{u(c,s) = cos\left(\vec{w_c}, \vec{w_s}\right) = \frac{\bar{w_c} \cdot \bar{w_s}}{\|\bar{w_c}\|_2 \times \|\bar{w_s}\|_2}}$$

## Collaboreative method :

$$u(c,s) \longleftarrow u\left(\boxed{c_v^j}, s\right) \; ; \; c_j \in C \; \text{\&} \; c_j \approx c$$

( similar user group)
(peer)

### ① memory based / Heuristic.

rating , $\boxed{r_{cs}}$ = $\boxed{\text{aggregate}}$ $r_{c',s}$ // impute unknown value.

$c' \in C$

not given but estimated     given ratings

⇓ Agg. function can be ?

$$r_{c,s} = \begin{cases} ⓐ \; \frac{1}{N} \sum\limits_{c' \in C} r_{c',s} & \text{more like collaboration.} \\ & \text{terom.} \\ ⓑ \; k \sum\limits_{c' \in C} sim(c,c') \times r_{c',s} \\ ⓒ \; \bar{r_c} + k \sum\limits_{c' \in C} sim(c,c') \times \left(r_{c',s} - \bar{r_{c'}}\right) \end{cases}$$

where, $k = \dfrac{1}{\sum\limits_{c' \in C} |sim(c,c')|}$ (Normalizing constant)

$$\bar{r_c} = \left(\frac{1}{|S_c|}\right) \sum\limits_{s \in S_c} r_{c,s} \quad \text{where} \; S_c = \{s \in S | r_{c,s} \neq \phi\}$$

(Average rating)

P 12

Pearson coefficient based similarity:

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

// Iterate over (s) item

Alternatively, Cosine based Similarity.

for each item $s$

(ii) model based Algorithm:

$$c_s = E(r_{c,s}) = \sum_{i=0}^{n} i \times Pr\left(r_{c,s} = i \mid r_{c,s'}, \; s' \in S_c\right)$$

**P 13**    selected dimentionality $f$ :

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

$$P_u \in \mathbb{R}^f \quad // \text{ users.}$$

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 → bad idea.

So, optimization problem:

$$\min_{P, q} \sum_{(u,v) \in k} \left( r_{ui} - q_i^T P_u \right) + \lambda \left( \| q_i \|^2 + \| P_u \|^2 \right) \quad ——②$$

given training Set (previously observed)
(Explicitly feed back points)

**Learning methods**

①SGD

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

$$q_i \leftarrow q_i + \gamma \left( e_{ui} \cdot P_u - \lambda q_i \right) \quad // \text{ Gradinet descent}$$

$(\gamma)$ is changed by internal calculation

$$P_u \leftarrow P_u + \gamma \left( e_{ui} \cdot q_i - \lambda P_u \right) \quad \pm \text{ fast}$$

②

**Alternate // ALS**

to solve nonconvexity → fix one, and solve for the other.

P 13

Existance of product/user bias.

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

→ overall reating

So, $r_{ui} = \underline{\mu} + \underline{b_i} + \underline{b_u} + q_i^T p_u$ —— ⓘⱽ

Now the optimization problem changes to,

$$\min_{P,q,b} \sum_{(\mu,i)\in k} \left(r_{ui} - \mu - b_u - b_i - q_i^T p_i\right)^2 + \lambda\left(\|p_u\|^2 + \|k_i\|^2 + b_u^2 + b_i^2\right) —— ⓥ$$

// may bias the model.

Additional Input source: cold start overcome.

$N(u)$ → implicit preference on items by the users.

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

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

Normalization ⇒ $|N(u)|^{-0.5} \sum\limits_{i \in N(u)} x_i^{+.5}$ // empirical.

user Attributes →$A(u)$ ⇒ $\boxed{y_a \in \mathbb{R}^f}$ Associated factor to the Attributes.
                    set        elements.

Now overall: $r_{ui} = \mu + b_i + b_u + q_i^T\left[p_u + |N(u)|^{-0.5} \sum\limits_{i \in N(u)} x_i \right.$

$$\left. + \sum\limits_{a \in A(u)} y_a\right] —— ⑥$$

two extra terms.

P13

Temporal dynamics:

Including time $\quad \hat{r}_{ui}(t) = \mu + b_i(t) + b_u(t) + q_i^\top p_u(t)$

Dynamic

static → item bias

user bias

static → human behaviour dynamics.

Input with confidence level:

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

→ modified confidence term.

(Variational Lower bound)

$z \rightarrow x$   $P(x) \rightarrow$ Prob dist. over variable
Latent observed.   $p(x) \rightarrow$ Pdf of Distribution $x$

Posterior   $P(z|x) = \dfrac{P(x|z)\, P(z)}{\displaystyle\int_z P(x|z)\, P(z)}$

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

$(-)$ or 0 max,   $= \log \displaystyle\int_z P(x,z)\, \dfrac{q(z)}{q(z)}\, dz$

Reorganizing

$L = \underset{q}{E}\big[\log P(x|z)\big]$   $\geq \underset{q}{E}\Big[\log \dfrac{P(x,z)}{q(z)}\Big]$   // Jensen's inequality

$- KL\big[q(z)\|p(z)\big]\underline{\ ELBO}$,   $L = \underset{q}{E}\big[\log P(x,z)\big] + \underset{q}{H}(z)$   // entropy def.

$\underbrace{\phantom{E[\log P(x,z)]}}$   $\underbrace{\phantom{H(z)}}$

$|\text{Negative}| > |+| $  //so overall $(-)$

$\overline{|\text{Always else } P(x) = \emptyset \, 1|}$ : no info.

Interprete: $(-)$ of information $>$ $\underline{ELBO}$   // reverse it $-$(interesting)

**No more info than** $-|ELBO|$ in $P(x)$
$\Rightarrow$ more info than $|ELBO|$

**Derivation 2:**   $KL\big(q(z)\|P(z|x)\big) = \displaystyle\int q(z)\, \log \dfrac{q(z)}{P(z|x)}\, dz$

(Backward KL)??

fixed   $= -L + \log P(x)$   // easy $P(z|x) = \dfrac{P(x,z)}{P(z)}$

// missed margin of $q(z)$

$\therefore \boxed{\log P(x) = L + \boxed{KL\big(q(z)\|P(z|x)\big)}}\,{\geq 0}$   to estimate $p(z|x)$

(failed) if large

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

**interpretation:**
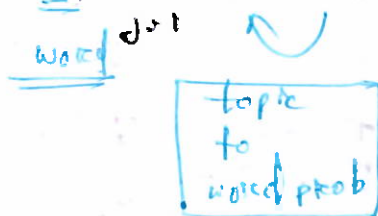
By making elbo highest means $q(z) \approx P(z|x)$
successful posterior estimation.

P 15

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

word — topic to word prob — topic prob.

No. Document → D
No. of topic → T

(Distribution over Dist.)

(multinomial probs are also from Dist) → Dirichlet (conjugate prior)

D documents, T topics & W unique words.

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

event (count) (topic to word)

(k face dice n time, count) | what events?
(w words total W counts) | how many time!

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

(d topics total D times)

Now the objective | Maximize $P(w | \theta, \phi)$

Modified objective for dirichlet distribution.

LDA

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

with dirichlet parameter($\alpha$) (as conjugate prior for multinomial)

determined by → Variation bayes / Expectation propagation.

The Complete model (Gibbs Sampling usage opportunity)

$$\begin{cases} w_i | z_i, \phi^{z_i} \sim \text{Discrete}(\phi^{z_i}) & \leftarrow \text{from earlier} \\ \phi \sim \text{Dir}(\beta) & // \text{new [conjugate prior)} \\ z_i | \theta^{d_i} \sim \text{Discrete}(\theta^{d_i}) & \leftarrow \text{earlier} \\ \theta \sim \text{Dir}(\alpha) & // \text{new [conjugate prior]} \end{cases}$$

$\alpha, \beta$ hyperparameter.

P 15

By integrating w.r.t. $\theta, \phi, \Rightarrow P(w, z)$

$P(w, z) = P(w|z) P(z)$

No document here

word $w \leftarrow$ j topic

no. time (Assign)

where, 
$$P(w|z) = \left(\frac{\Gamma(W\beta)}{\Gamma(\beta)^W}\right)^T \prod_{j=1}^{T} \frac{\prod_w \Gamma(n_j^w + \beta)}{\Gamma(n_j^{(\cdot)} + W\beta)}$$

//see wiki
multinomial
Approx.
—②

j topics

similar to

Gramma function.

$$\Gamma(z) = \int_0^\infty x^{z-1} \bar{e}^{x} dx$$

similarly,
$$P(z) = \left(\frac{\Gamma(T\alpha)}{\Gamma(\alpha)^T}\right)^D \prod_{d=1}^{D} \frac{\prod_j \Gamma(n_j^{(d)} + \alpha)}{\Gamma(n^{(d)} + T\alpha)}$$

no of time
word w from
docu d
j topic
district from early.
—③

$\Rightarrow \quad P(z|w) = \dfrac{P(w, z)}{\sum\limits_{z} P(w, z)}$

//{Again Intractable $T^n$ terms
so require Approximation

$\downarrow$

Solve it By MCMC, (Gibbs Sampling).
$\downarrow$
require $P(z_i | z_{-i}, w)$

using the earlier ② and ③ we get. (by cancellation)

$$P(z_i | z_{-i}, w) \propto \frac{n_{-i,j}^{(w_i)} + \beta}{n_{-i,j}^{(\cdot)} + W\beta} \cdot \frac{n_{-i,j}^{(d)} + \alpha}{n_{-i,j}^{(\cdot)} + T\alpha}$$

(prob of $w_i$ under topic $j$)

prob of topic $j$ in document $d$

//just need counter

$n_{-i}^{(\cdot)} \rightarrow$ not include current assignment $(z_i)$

P15: for single Sample: (using just count)

$$\hat{\phi}_j^{(w)} = \frac{n_j^w + \beta}{n_j^{()} + \beta}$$    // for new word w
                                                              and new topic z.

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

Solved by Gibbs Sampling ↙

Alternates: Variational bayes ⑦ , Expectation propagot ②

Total hyperparameter : $\alpha$, $\beta$, Ⓣ → (Varied Across)

(Fixed it    ↓ Topic Number ??
in experiment) ↓
                This is model selection?

Target :    $P(w | T)$   ┌→ topic
                          Number.
            ↓
         ~All words  ↘ Approximated
                       by      posterior
                    $P(w | z, T)$ ⇒ $P(z | w, T)$

Expected value.

empirical ①

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

→ $\frac{1}{N}$ ✗ No. of time $(x,y)$ appears. (training data) (model)

training data ↓ calculate from data.

↓ model given ↓ indicator function

$P(x) \sim x$ (empirical) of x

feature function

we requires $\boxed{\tilde{P}(f) = P(f)}$

train model (training data)

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

model training data.

constraint equation

$C \equiv \{ P \in P \mid P(f_i) = \tilde{P}(f_i) \text{ for } i = 1 \cdots n \}$

Key goal

Entropy: $H(x) = \sum_{x} P(x) \log\left[\frac{1}{P(x)}\right]$

Fig 1

All the 2 *possible Data prob

lower ↓ bound

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

model ↑ sure model

(set of probabilities distribution) $H(P) \leq$ uniform case entropy.

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

train model upper bound. $= \log |y|$

$P_* = \arg\max_{P \in C} H(P)$ // uniform case

$\parallel$ cardinality of y

well defined & unique.

⇒ To select a model from a set of prob. distributions, that has maximum entropy.

Parametric Form:

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

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

training (ground truth)

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

→ (model)

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

Solving, $\boxed{P_\lambda (y|x)} = \frac{1}{Z_\lambda(x)}$ exp $\left( \sum_i \boxed{\lambda_i} f_i(x,y) \right)$, $\boxed{\lambda^*}$ — unknown (solved later)

model

$\sqcup$ normalization — w.r.t. y

model

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

(??)

Dual optimizates find $\boxed{\lambda^*}$ = argmax $\Psi(\lambda)$ // maximize
$\qquad\qquad\qquad\qquad\qquad\;\lambda$

... $\lambda^*$ → parametric form

## Relation to maximum likelihood; Training data.

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

By definition, $\Rightarrow \Psi(\lambda) = L_{\tilde{P}}(P_\lambda)$ // see earlier section.

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

$\boxed{\text{Table 1}}$ $\boxed{\text{summary}}$

$\boxed{\text{Compute the Params}}$ : ① $f_i(x,y) \geq 0$

Algo ! Input iterative Scaling.

input : $f_1 \cdots f_n$ → empirical $\tilde{P}(x,y)$

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

1. $\lambda_i = 0 \; (i = 1 \cdots n)$
2. $i \in \{1 \cdots n\}$

See the algorithm

P12

$x_i \in \mathbb{R}^n$ ; $i: 1 \ldots l$ observations. $\rightarrow$ or. $y_i \in \{1, -1\}$

$f(\underline{x}, \alpha) \rightarrow$ Approximation

Expectation of test error.     parameter

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

$$R_{emp}(\alpha) = \frac{1}{2l} \sum_{i=1}^{l} |y_i - f(\underline{x}_i, \alpha)| \quad \text{// Approximation.}$$

Connection between them     VC confidence.↓

Vapnik, 1995

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

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

set to minimize this bound

$l \rightarrow$ Example

$\left\{ \begin{array}{l} h \uparrow \\ h \downarrow \text{ we want} \\ l \uparrow \end{array} \right.$

$h \rightarrow$ VC dimension.

[VC confidence lower $\Rightarrow$ the better]   (may overfit)
Approximation

The VC dimension: increases function capacity ↑ confidence boundary (higher is Bad)

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

true / Approximate

$x_i = 10^{-i}$

$y_i = \ldots$ Assign anything.

$\alpha = \pi \left( 1 + \sum_{i=1}^{l} \frac{(1 - y_i) 10^i}{2} \right)$

VC dimension $\neq \infty$

shattering depends on choice of points.

chose points that can be shattered.

## Seperable Case

$\boxed{11}$ projection to $HP = \boxed{-b}$

satisfying hyperplane: $\quad \underline{w}.x + b = 0$
(HP)

bias
(All points
projected as
$-b$ amount)
(nice)

& Linear $\nearrow$

$\Downarrow$

normal to HP.

Seperable case $\begin{cases} x_i.w + b \geq 1 \; ; \; \text{overshoot in projection} \; y_i = +1 \\ x_i.w + b \leq 1 \; ; \; y_i = -1 \end{cases}$

$\rightarrow \quad y_i (x_i.\underline{w} + b) - 1 \geq 0$ ‖

multiply & add alls.

introducing lagrangian $\alpha_i; i = 1 \cdots \ell; \quad \alpha_i \geq 0$

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

primal problem.

$\begin{cases} w = \sum_i \alpha_i y_i x_i \quad // \text{solution.} \\ \sum \alpha_i y_i = 0 \end{cases}$

Now the Dual problem.

use $k(x_i, x_j)$
in nonlinear case

putting values $\hookrightarrow$ $$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \left( \boxed{x_i . x_j} \right) \quad // \text{Dual formulation}$$

kernel idea.

## KKT condition:

$$\frac{\partial}{\partial w_\nu} L_P = w_\nu - \sum_i \alpha_i y_i x_{i\nu} = 0 \quad ; \nu = 1 \cdots d$$

$$\frac{\partial}{\partial b} L_P = \cancel{y_i (x_i + \underline{w} + b) - 1} \geq \sum_i \alpha_i y_i = 0$$

$$y_i (\underline{w}_i . x_i + b) - 1 \geq 0 \qquad i = 1 \cdots \ell$$

$$\alpha_i \geq 0 \qquad \forall i$$

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

(if seperable)

Non linear SVM:

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

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

Need $\underline{w}$ of $\lfloor \mathcal{H} \rfloor$ dimensional ??

How to use the kernel ??

→ we just Need dot product.

(using the train data )
No of support vector

test phase → $f(x) = \sum_{i=1}^{N_s} \alpha_i y_i \phi(s_i) \cdot \phi(x) = \sum_{i=1}^{N_s} \alpha_i y_i \underbrace{k(s_i, x)}_{\substack{\text{only} \\ \text{need this.}}} + b$

Mercers Conditiona.-  $(\mathcal{H}, \phi) (d \rightarrow \mathcal{H})$

$$k(x, y) = \sum_i \phi(x)_i \phi(y)_i \quad \underline{\text{mapping exists}}$$

if $\forall$ $g(x)$ 2 that satisfy.

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

then

$\int k(x, y) g(x) g(y) \, dx \, dy \geq 0$ ⟹ Positive Semidefinite
$\underline{PSD}$

Open Question : How to formulate $\phi$ ?

since vc dimension is $\lfloor \mathcal{H} \rfloor + 1$ // in this case

so $\lfloor \mathcal{H} \rfloor \uparrow$ bad generalize.

Radial basis kernel :

$$k(\underline{x}, \underline{y}) = e^{-\|\underline{x} - \underline{y}\|^2 / 2\sigma^2} \quad // \text{ may be Infinite } d \\ \text{vc Dimension.}$$

↳ two layer sigmoid NN.

( generalize
n Vs 1 classifier )

1. Layer ⓪ $N_s$ set weights each $d$ dimensional

2. Layer Ⅱ $N_s$ weights ($\alpha_i$)

finally Sigmoid .

F19

Nonseparable Case:

may cause
⌐→error

$$\underline{x_i} \cdot \underline{w} + b \geqslant +1 - \varepsilon_i \quad ; y_i = +1$$

$$\underline{x_i} \cdot \underline{w} + b \leqslant -1 + \varepsilon_i \quad ; y_i = -1$$

$$\varepsilon_i \geqslant 0 \quad \forall i$$

⇒ if any $\varepsilon_i > 1$ error occures.

$$\sum_i \varepsilon_i = \text{upper bound of traing error.}$$

Dual problem⇒ $$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j \, y_i \, y_j \, \underline{x_i} \cdot \underline{x_j}$$

Subject to: $\boxed{0 < \alpha_i \leqslant C}$ → user parameter.

(higher penalty
to error)

$$\sum \alpha_i y_i = 0$$

Solution is $$w = \sum_{i=1}^{N_s} \alpha_i y_i \, \underline{x_i}$$

P20

skip gram model   maximize $\frac{1}{T} \sum_{t=1}^{T} \sum_{\substack{-c \leq j \leq c \\ j \neq 0}} \log P\left(w_{t+j} | w_t\right)$

sequence of words $\{w_1, w_2 \cdots w_T\}$

So, maximizing $P\left(w_{t+j} | w_t\right)$ ∝

Defined as   $P(w_0 | w_I) = \dfrac{\exp\left(v'_{w_0}{}^T v_{w_I}\right)}{\sum_{w=1}^{W} \exp\left(v'_{w}{}^T v_{w_I}\right)}$

$\left\{\begin{array}{l} c \uparrow \\ \text{training time} \uparrow \\ \text{Accuracy} \uparrow \end{array}\right.$

$\left[\begin{array}{l} v'_{w_0} = \text{output vector reps.} \\ v_{w_0} = \text{input vector reps} \end{array}\right]$   $\underset{w=1}{\hookrightarrow}$ All the words ?? [huge computation]

$v_{w_0} \to$ vector representation of $w_0$ (via Networks)

$W \to$ huge size !!

┌─────────────────────┐
│ Each word has │
│ two reps │
│ → input $v_w$ │
│ → output $v'_w$ │
└─────────────────────┘

**Hierarchical Softmax:** Need $\log_2 W$ nodes.

$P(w | w_I) = \prod_{j=1}^{L(w)-1} \sigma\left(\left[ n(w, j+1) = ch(n(w, j)) \right] v'_{n(w,j)}{}^T v_{w_I}\right)$

$\underset{\text{1 if true, else 0}}{}$   [computation ∝ L(w)]

→ care about [input/output] representation.

**Negative Sampling** ↙ should be high (interesting)

NEG → objective. $\log \sigma\left(v'_{w_0}{}^T v_{w_I}\right) + \sum_{i=1}^{K} \underset{w_i \sim P_n(w)}{\mathbb{E}} \uparrow \left[\log \sigma\left(-v'_{w_i}{}^T v_{w_I}\right)\right]$

$\underset{\text{row of matrix}}{\Downarrow}$   → different choices

positive   negative word   should be low & negative

how (reg)

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

//modified NCE

subsampling of freq words.

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

→ word frequency.

// Discarding prob.

∽ $f(w_i) \uparrow \quad P(w_i) \uparrow$

∽ $f(w_i) \downarrow \quad P(w_i) \downarrow$

( High freq words discarded more )

Balance between rare & frequent words.

Bigram skip:

(Phrase selection)

bigram

$$score(w_i, w_j) = \frac{count(w_i, w_j) - \delta}{count(w_i) \times count(w_j)}$$

unigram    unigram
$w_i$          $w_j$

P21

The probabilistic model:

$$P(\{s_t, y_t\}) = P(s_1) \, P(x_1 | s_1) \prod_{t=2}^{T} P(s_t | s_{t-1}) \, P(y_t | s_t)$$

observed, D

↓

Hidden state, k

Separable.
Conditional Independence.

Let, k states,

$$P(s_t | s_{t-1}) \Rightarrow k \times k \text{ matrix (state transition matrix)}$$

$$P(y_t | s_t) \Rightarrow k \times D \text{ observation matrix}$$

↳ modeled by Gmm/Neural Network.

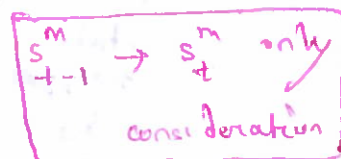what if: $s_t = s_t^{(1)}, s_t^{(2)} \cdots, s_t^{(m)}$ } factorial HMM !!

$k^{(m)}$ possible values. ←each

$k^{(m)} = k$ // simplicity

Then, $k^M$ states [ $k^M \times k^M$ ] state transition matrix !!

↳ impossible to work with

↳ Requires constraint on state tx. mat.

factorial HMM → underlying state tx is constrained

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

Decoupled.

$s_{t-1}^m \to s_t^m$ only
consideration

p21

observation:

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

↗ D×D = covariance Matrix

$$\mu_t = \sum_{m=1}^{M} \underline{w^{(m)}} s_t^{(m)} \quad \text{// for all M-th contribution} \rightarrow \text{probabilistic}$$

→ D×k [columns are contribution of each states] ⇒ D×1 finally.

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

⇒ Depends on $s_t^{(m)}$ value 𝒘

## Learning & Inference :

### Expectation maximization : param learning

$$Q(\phi^{new} | \phi) = E\left\{ \log P(\{s_t, Y_t\} | \phi^{new}) \,\Big|\, \phi, \{Y_t\} \right\} \quad — ⑤$$

↓ new params.     ↓ current params

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

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

$\downarrow$
$P(s_1^{(m)})$

### E step :

Compute Q : Expand S by using earlier equations.

→ Can be expressed as Expectation of

$$E\{ \cdot | \phi, Y_t \} \Rightarrow \langle s_t^{(m)} \rangle \; ; \; \langle s_t^{(m)} . s_t^{(n)} \rangle \; ; \; \langle s_{t-1}^{(m)} s_t^{(m)} \rangle$$

state occupation $\gamma_t$  K×1 vec.

(two states jointly) ??

state transition $\sum_t$ K×K mat

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

solved by : weighted linear regression.

P21

Gibbs Sampling: Inference:

$$s_t^{(m)} \sim P\left(s_t^{(m)} \mid \underbrace{\{s_t^n : n \neq m\}}_{\text{Neigh}}, \underbrace{s_{t-1}^{(m)}}_{\text{past}}, \underbrace{s_{t+1}^{(m)}}_{\text{future}}, \underrightarrow{y_t}_{\text{observation}}\right)$$

$$\propto \underbrace{P\left(s_t^{(m)} \mid s_{t-1}^{(m)}\right)}_{\text{state transition}} \underbrace{P\left(s_{t+1}^{(m)} \mid s_t^{(m)}\right)}_{\text{state transition}} \underbrace{P\left(y_t \mid s_t^{(1)}, \cdots s_t^{(m)}, \cdots s_t^{(m)}\right)}_{}$$

$\boxed{\text{graphical model design itself}}$

↓ markovian

completely factorized Variational Inference:

$$\underbrace{\log P(\{y_t\})}_{\text{①}} = \log \sum_{\{s_t\}} P(\{s_t, y_t\})$$

$$= \log \sum_{\{s_t\}} Q(\{s_t\}) \frac{P(\{s_t, y_t\})}{Q(\{s_t\})}$$

$$\geq \underbrace{\sum_{\{s_t\}} Q(\{s_t\}) \log \left[\frac{P(\{s_t, y_t\})}{Q(\{s_t\})}\right]}_{\text{②}}$$

The difference between ② & ① is $|① - ②|$  //simple math.

$$KL(Q \| P) = \sum_{\{s_t\}} Q(\{s_t\}) \log\left[\frac{Q(\{s_t\})}{P(\{s_t \mid y_t\})}\right]$$

→ change parameter of $Q(\{s_t\})$ to minimize:

$p=1$

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

$\underbrace{}_{\text{time step}}$  $\underbrace{}_{\substack{\text{possible} \\ \text{steps at each } t.}}$

→ vector itself.

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

vector element

(m)

$\rightarrow$ $m$-th markovian chain

$\rightarrow$ state $k$, at time $t$

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

$\underbrace{}_{\text{multiply.}}$

$\downarrow$ All the params.

with $\sum_{k=1}^{K} s_{t,k}^{(m)} = 1$

$\underbrace{}_{\text{only one is } = 1 \text{ else } 0}$

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

$\theta_t^{(m)} \rightarrow$ state occupation prob. with multinomial var $s_t^{(m)}$

under distribution $Q$

→ vector → softmax element wise.

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

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

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

## 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

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

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

$$Q\left(S_t^{(m)} \mid S_{t-1}^{(m)}, \theta\right) = \prod_{k=1}^{K} \left( h_{t,k}^{(m)} \prod_{j=1}^{K} \left(P_{kj}^{(m)}\right)^{S_{t-1,j}^{(m)}} \right)^{S_{t,k}^{(m)}}$$

one hot vector.

$$\theta = \left\{ \pi^{(m)}, P^{(m)}, h_t^{(m)} \right\}$$

$k \times 1 \to$ prob of observation $P(Y_t \mid S_t)$

for each $k$ setting $S_t^{(m)}$

$$Q\left(S_{1,j}^{(m)} = 1 \mid \theta\right) = h_{1,j}^{(m)} \; P\left(S_{1,j}^{(m)} = 1 \mid \phi\right)$$

$\Rightarrow$ having an observation at $t = 1$, under $S_{1,j}^{(m)} = 1$

has prob of $h_{1,j}^{(m)}$

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

$$h_t^{(m) \, new} = \exp\left\{ W^{(m)'} C^{-1} \tilde{Y}_t^{(m)} - \frac{1}{2} \Delta^{(m)} \right\}$$

$$\text{reci} \longrightarrow = Y_t - \sum_{l \neq m}^{M} W^{(l)} \langle S_t^{(l)} \rangle$$

connected to ELBO bound

$$F(Q, \phi) = \mathbb{E}_Q\left\{ \log P(Y, S \mid \phi) \right\} - \mathbb{E}_Q \log\left\{ Q(S) \right\} \leq \log P(Y_s)$$

$(-)$ value. $\to$ reverse it (itererting)

No more info than $|ELBO|$

P22

Problem formulation:

$\{A_1, A_2 \cdots A_m\} \longrightarrow$ m smart phone.

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

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

Query $Q = \{q_1 \cdots q_f\}$; $f \ll L$ 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:

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

$LCSS_{\delta,\epsilon}$ → time matching window / spatial matching window

Both are application specific.

[ iterative Algorithm ]

(x, y at time 1)

$Head(A) = \left((a_{x:1}, a_{y:1}), \cdots, (a_{x:L-1}, a_{y:L-1})\right)$

P>2

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 & Env High[i]
Low envelop.    Env Low[i]

$$\text{Env High}[i] = \max(Q[j] + \epsilon); \quad |i-j| \leq \delta$$

$$\text{Env Low}[i] = \min(Q[j] - \epsilon); \quad |i-j| \leq \delta$$

unique Solution

$\begin{cases} G \to \text{recovers the training Data Distribution} \\ D \to \frac{1}{2} \text{ everywhere} \end{cases}$

## Adversarial Nets:

$gen \times \sim p_g$

Noise vector $p_z(z) \to$ map $G(z, \theta_g)$

$\Downarrow$

Differentiable function (MLP)

multi layer perceptron

$D(x, \theta_d) \to$ Differentiable MLP.

↳ Discreemination → x from data / G ??

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

→ Iterative Numerical approach:

$\begin{cases} k \text{ step of } D \to \text{to keep near optimal (inner loop)} \\ 1 \text{ step of } G \to \text{changes slowly enough.} \end{cases}$ {SML/PCD way ??}

## Theory : Algorithm 1 → crack of jack.

understanding sequence: when D is optimal ⟹ ?? ✓

     (keep it) 2^near

How/when G is optimal | D is optimal ↵

what happens when D is optimal ?? ✓

what happens when both are optimal

for no of train

for _ k in range (4)

sample $z^1, \ldots z^{(m)} \to$ $P_g(z)$ ∅)

sample $x^1 \cdots x^m \to$ $P_{data}$ ☺)

update $\theta_d$, D by ascending ~~grad~~ stochastic grad.

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^{m} \left[ \log D(u^{(i)}) + \log \left( 1 - D\left( G\left( z^{(i)} \right) \right) \right) \right] \quad // maximize$$

*m more loop*

*keep nearn optimal*

→ end. for

sample m noise $\{ z^{(1)} \ldots z^{(m)} \}$

Update the gen. $\theta_g$ descending gradient

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^{m} \log \left( 1 - D\left( G\left( z^{(i)} \right) \right) \right) \quad // \cdot minimize.$$

end _ for.

global optimality: $P_g = P_{data}$.

step 1: For G fixed,

$$D_G^*(u) = \frac{P_{data}(u)}{P_{data}(u) + P_g(u)} \quad // optimal \text{ Discriminator (maximize)}$$

training criterion:

$$V(G, D) = \int_u P_{data}(x) \log(D(x)) dx + \int_z P_z(z) \log \left( 1 - D(G(z)) \right) dz$$

$$= \int_x P_{data}(u) \log(D(u)) du + \int_x P_g(u) \log \left( 1 - D(u) \right) du$$

for any function of $y \to a \log y + b \log (1-y)$ achieves its

maximum in $[0, 1]$ at $\frac{a}{a+b} = y)$

$$C(G) = \max_D V(G, D)$$

$$= \mathbb{E}_{x \sim P_{data}}\left[\log D_G^*(x)\right] + \mathbb{E}_{z \sim P_z}\left[\log\left(1 - D_G^*(G(z))\right)\right]$$

$$= \mathbb{E}_{x \sim P_{data}}\left[\log D_G^*(x)\right] + \mathbb{E}_{x \sim P_g}\left[\log\left(1 - D_G^*(x)\right)\right] \quad // \text{ from our earlier argument } \left(y = \frac{a}{a+b}\right) \text{ maximize}$$

$$= \mathbb{E}_{x \sim P_{data}}\left[\log \frac{P_{data}(x)}{P_{data}(x) + P_g(x)}\right] + \mathbb{E}_{x \sim P_g}\left[\log \frac{P_g(x)}{P_{data}(x) + P_g(x)}\right]$$

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

<u>The global minima</u> for $C(G)$ is if $\boxed{P_g(x) = P_{data}(x)}$   // Now D fails to comprehend

in that case,

$$\mathbb{E}_{x \sim P_{data}}\left[\log \frac{P_{data}(x)}{P_{data}(x) + P_g(x)}\right] + \mathbb{E}_{x \sim P_g}\left[\log\left[\frac{P_g(x)}{P_{data}(x) + P_g(x)}\right]\right] =$$

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

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

Re organizing the equation ① we get   [$P_g$ may not be optimal]

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

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

$$= -\log 4 + 2 \cdot JSD\left(P_{data} \| P_g\right) \quad // \text{ jenson shanon divergence}$$

for optimal $G^* \Rightarrow C(G) = -\log 4$ !!  // if $P_g$'s optimal → if only

[ minimum for optimal D ]  ( $P_g = P_{data}$ )

(IV)

Convergence of Algo: let discriminator reaches optimal $\textcircled{1}$ for

$$\boxed{D(x) = P_{data} / (P_g + P_{data})}$$ $\textcircled{2}$

and $P_g$ is updated to improve.

$\textcircled{1}$

$\boxed{then \left( G \Longleftrightarrow (P_g) \right) \rightarrow (P_{data})}$

$\Rightarrow \quad \underset{x \sim P_{data}}{E} \left[ log \; \overset{k}{D}_G(x) \right] + \underset{x \sim P_g}{E} \left[ log \left( 1 - \overset{\ell}{D}_G(x) \right) \right]$

$\textcircled{ii}$ Now $P_g$ converge to data.

$\|$

MLP via function $G(z; \theta_g)$

optimize this instead of $P_g$

(multiple critical point) $\Rightarrow$ multiple local optimal

theory $\Rightarrow$ (not yet)

---

Dis-Ad: G must not be trained too much !! $\left( \begin{array}{c} D \; must \; be \\ sync \; with \; G \end{array} \right)$

mode collapse !!   convergence ??

Ad: only backprogpaga $\checkmark$

Large distribution learning.

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

ⓐ <u>Paper 30</u>

Joint Dist    $P\left(x, \overset{\text{hidden}}{h} \mid n\right) = P(x \mid h)\, P(h \mid n)$

      obs     param

$P(h \mid x, n) \propto P(h, x \mid n)$

predictive,    $P(x_{new} \mid x) = \int P(x_{new} \mid h)\, P(h \mid x, n)\, dh$.

     mixture (dir)

$P\left(\mu_{1:k}, \theta, z_{1, N} \mid x_{1:N}\right)$

       ↓         ↓      ⇓

    <u>mean</u>       class    observ

<u>Generative probabi model</u>:

                global hid van

mixture prop $\left( \theta \sim \text{Dirichlet}(\underline{\alpha}) \right.$

            $\left. \mu_k \sim N(0, \underline{\underline{\sigma_0^2}}) \right.$

   ⓐ mixture assignment   $z_n \mid \theta \sim \text{Discrete}(\theta)$ ?

   ⓑ   Data point   $x_n \mid z_n, \mu \sim N(\mu_{z_n}, 1)$.

easy peasy:- sample $\theta \to$ corresponding $\mu_{1 \cdots k} \to$ for each

data point estimate $z_n \in 1 \cdots k \to$ govern the distribution.

Paper 30

Joint distribution:

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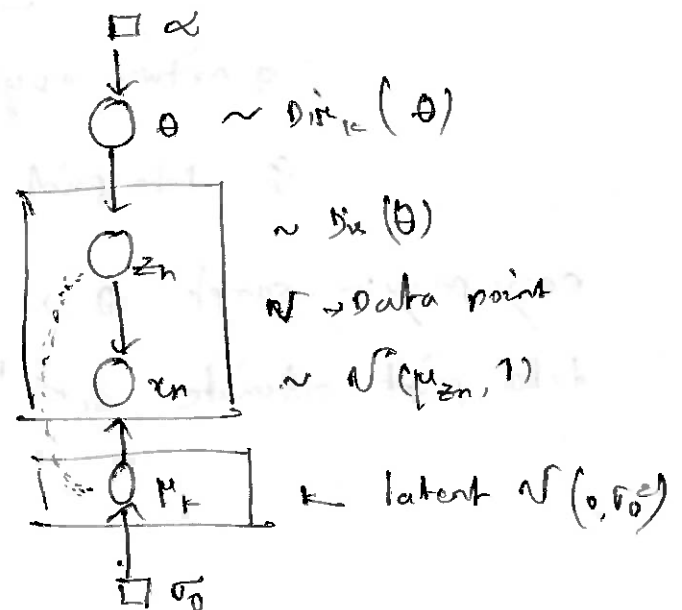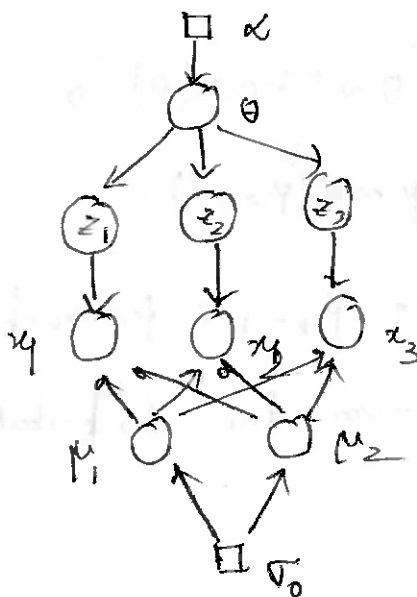
Posterior

$$P(\theta, \mu, z \mid x, \sigma_0^2, \alpha) = \frac{P(\theta, \mu, z, x \mid \sigma_0^2, \alpha)}{\oint P(x \mid \sigma_0^2, \alpha)}$$

Predictive Distribution.

$$P(x_{new} \mid x, \sigma_0^2, \alpha) = \int \left( \sum_{z_{new}} P(z_{new} \mid \theta) P(x_{new} \mid z_{new}, \mu, \sigma_0^2) P(\theta, \mu \mid x, \sigma_0^2, \alpha) \right. \\ \left. d\theta \, d\mu \right.$$

The graphical model:



$\theta \sim Dir_k(\theta)$

$\sim B_k(\theta)$

$N \sim$ Data point

$\sim N(\mu_{z_n}, 1)$

$k$ latent $N(0, \sigma_0^2)$

Example model:

① linear factor model: PCA, factor model. (graph 3)

⑪ mixed membership model:

⑭ Matrix factorization model:

④ Time Series
   - ~~Time Series models~~ Hidden Markov model
   - Kalman filter

Posterior Inference  With mean field: Variational method.



Latent var. mod          Variational family:

Conditional Conjugate model:

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

local latent ↑ (over $z$)

↓ global latent  ↓ obs

↑ Always fixed  mixture prop

↓ changes so local

Posterior:

$$P(\beta \mid x) = \frac{P(\beta, z, x)}{\int P(\beta, z, x) d\beta, dz} \rightarrow \text{how? Problem!}$$

Dependant ↗          natural param. ↗

$$P(x|\eta) = h(x) \exp\left(\eta^T t(x) - a(\eta)\right) \quad \sim \text{exponential family.}$$

some suff. sufficient stat (great) → log normalizer.

base measure.

## mean field variational model: Approximate Posterior.

variational Objective function

$$v^* = \underset{v}{\arg\min} \, KL\left(q(\beta, z|v) \| P(\beta, z| x)\right)$$

(good read)

$$L(v) = E\left[\log P(\beta, z, x|\eta)\right] - E\left[\log q(\beta, z|v)\right]$$

## mean field variational family:

$$q(\beta, z|v) = q(\beta|\lambda) \prod_{n=1}^{N} q(z_n|\phi_n)$$

$$v = \{\lambda, \phi_{1-N}\} \quad \} \quad \text{easy.}$$

## coordinate Ascent Varia. Inference:

$$\lambda^* = E_q\left[\eta_g(z, x)\right] \quad \text{global}$$

$$\phi_n^* = E_q\left[\eta_\ell(\beta, x_n)\right]$$

↳ local

$$q(\mu, z) = \prod_{i=1}^{k} q(\mu_k|\lambda_k) \prod_{n=1}^{N} q(z_n|\phi_n)$$

<u>Model Criticism:</u>   Exploration & prediction.

$$\Downarrow$$

inference about
hidden vars.

$$\downarrow \text{distribution}$$

$$P(x_{new}|x) = \underbrace{\int P(\beta|x)}_{\parallel}\left(\underbrace{\int P(z_{new}|\beta)\, P(x_{new}|z_{new},\beta)\, dz}_{}\right)d\beta$$

not available

<u>Predictive Sample Reuse:</u>   $n$ removed   $P(\beta, z, x)$

$$l_n = \log P(x_n | x_{[n]})^{\uparrow}$$   $$P(\beta, z | x_{[n]})$$

$$\left[= \log \int \left(\int P(x_n|z_n)\, q(z_n)\, dz\right) q(\beta)\, d\beta \right]^{q_{[n]}(\beta,z)}_{[n]}$$

full likelihood $\overset{N}{\underset{n=1}{\Sigma}} l_n$

<u>Posterior Predictive Check:</u> $\rightarrow$ Test statistics

$$PPC = P(T(x^{rep}) > T(y|x))$$

$$\Downarrow$$

Data drawn from hypothetical future. obs.

$$PPC = P(T(x^{rep},\beta) > T(x,\beta)|x)$$

$$P(\beta, x^{rep}|x) = P(\beta|x)\, P(x^{rep}|\beta)$$

$$PPC = \int P(\beta|x) \int P(x^{rep}|\beta)\, \mathbb{1}_{[T(x^{rep},\beta) > T(x,\beta)]}\, dx^{rep}\, d\beta$$

P 10

$$T(x^+, \beta^+) > T(\mu; \sigma^+)$$

$$T(x, \mu) = \frac{1}{N} \sum_{n=1}^{N} \log(x_n | \beta) \quad // \text{ may be } T)$$

PPC is adaptive.

Basics ①

Batch Normalization in CNN: input $[N, C, H, W] \Rightarrow [i, c, j, k]$

Batch size
Height
width
filter channel

$\Rightarrow \quad \mu_{B,C} = \frac{1}{NHW} \left( \sum_{i=1}^{N} \right) \sum_{j=1}^{H} \sum_{k=1}^{W} x_{i,c,j,k}$

for all items → same filter ©

for each channel we get 1 value
So total C value.

$\Rightarrow \quad \sigma_{B,C}^2 = \frac{1}{NHW} \sum_{i=1}^{N} \sum_{j=1}^{H} \sum_{k=1}^{W} \frac{(x_{i,c,j,k} - \mu_{B,c})^2}{\cancel{1}}$

$\Rightarrow \quad$ final output $\quad \hat{x}_{i,c,j,k} = \frac{x_{i,c,j,k} - \mu_{B,c}}{\sqrt{\sigma_{B,c}^2 + \epsilon}}$

$\Rightarrow \quad$ Do calculation for each $c \in C$ channels.

More formally,

$$\hat{x}_{i,c,j,k} = \gamma \left( \frac{x_{i,c,j,k} - \mu_{BC}}{\sqrt{\sigma_{B,c}^2 + \epsilon}} \right) + \beta.$$

Solves : ① Internal Covariate Shift. (each zero mean, var-1)
→ features distribution differs internally.
→ inside the neural Network (layer - layer)

② Robust Network creation → less prone to perturbation

③ Learning faster.

(n)

## Instance Normalization (IN) / Layer Normalization

$$IN(x_{i,j,k}) = \gamma \left( \frac{x_{i,c,j,k} - \mu_{i,c}}{\sqrt{\sigma_{i,c}^2 + \epsilon}} \right) + \beta$$

↳ can be conditioned → conditional IN

where,

$$\mu_{i,c} = \frac{1}{HW} \sum_{j=1}^{H} \sum_{k=1}^{W} x_{i,c,j,k} \rightarrow \text{summed out.}$$

$$\sigma_{i,c}^2 = \frac{1}{HW} \sum_{j=1}^{H} \sum_{k=1}^{W} \left( x_{i,c,j,k} - \mu_{i,c} \right)^2$$

Instance: foreach n & c

## Adaptive Instance Normalization (AdaIN)

$$AdaIN(x_{i,c,\cdot}, y) = \sigma(y) \left( \frac{x_{i,c} - \mu_{i,c}}{\sqrt{\sigma_{i,c}^2 + \epsilon}} \right) + \mu(y)$$

pvar ↱  Adaptive term.

contional

↓

→ 0 if y is constant → retriive if y is constant

→ High if y varies a lot.

→ ↓ shold be 0 mean.

## Alternative Layer Normalization.

$$LN(x_{i,c,j,k}) = \gamma \frac{x_{i,c,j,k} - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$$

→ selecting subset leads to group Normaliza.

$$\mu_i = \frac{1}{H} \frac{1}{CW} \sum_{c=1}^{H} \sum_{j=1}^{W} \sum_{k=1}^{W} x_{i,c,j,k}$$

summed out.

Basic

## multinomial distribution:

total n trial

Each trial : Possible k outcomes $\{E_1, E_2 \cdots E_k\}$
with prob $\{P_1, \cdots P_k\}$ sterpetively.

Let's assume $E_1 \overset{happens}{-} n_1$, $E_2$ happens $n_2$ $\cdots$ $E_k \to n_k$ times.

So, $n_1 + n_2 \cdots + n_k = n$ // as n trial.

So, $P = \dfrac{n!}{n_1! \, n_2! \cdots n_k!} P_1^{n_1} P_2^{n_2} \cdots P_k^{n_k}$

$= \dfrac{n!}{n_1! \, n_2! \cdots n_k!} \prod_{i=1}^{k} P_i^{n_i}$

$= n! \prod_{i=1}^{k} \dfrac{P_i^{n_i}}{n_i!}$

→ straight extensions.

## Binomial distribution: Bernoulli trials.

n times fliping $x$ positive $n-x$ negatives.

$P(X = x| n, p) = {}^n C_x \, P^x (1-P)^{n-x}$

↓ (Distribution over Distribution)

## Beta distribution (Prior to Drichlet Distribution)

what if the p is a distribution itself ??

Beta is conjugate prior for binomial

$P \in [0, 1]$

Now, $P(P|\alpha, \beta) = \dfrac{1}{B(\alpha, \beta)} P^{\alpha-1} (1-P)^{\beta-1}$

$B(\alpha, \beta) = \dfrac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \, \Gamma(\beta)} = \binom{\alpha-1}{\alpha+\beta-2}$ //continuous estimation.

conjugate priors:

for some likelihood function, if we choose certain priors, the posterior ends up being the same function → then conjugate priors.

Dirichlet Distribution: Extended from multinomial probs.

Conjugate prior
tof.

what about the probs of multinomials $P_1 - - P_F$ ??

① $\Sigma P_i = 1$ // already known, as each events prob needs to be 1.

condition

$$P(P = \{P_i\} \mid \alpha_i) = \frac{\Pi_i \Gamma(\alpha_i)}{\Pi \Gamma(\Sigma_i \alpha_i)} \Pi_i P_i^{\alpha_i - 1}$$

↳ to multinomial.

→ Generalization of Beta.

→ Distribution over multinomials.

→ Conjugate prior of multinomial

↳ given the data the $\{P_i\}$ will also be Dirichlet distribution.

✳ Beta/Gamma functions are different:

$$\text{Beta}(x, y) = B(x, y) = \frac{\Gamma(x) \Gamma(y)}{\Gamma(x+y)} \text{ // binomial coeffs.}$$

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx \; ; \; \mathfrak{R}(z) > 0$$

$$\Gamma(z) = (z-1)! \text{ if } z \text{ is a natural number} \geq 0 \text{ positive integer.}$$