

LATENT VARIABLE MODELS FOR DISCRETE DATA

- DISCRETE DATA ONLY: IE TEXT ANALYSIS BUT OTHER STUFF TOO. **TOKENS**: VARIABLE LENGTH SEQUENCES OF CATEGORICAL DATA $y \in \{1..V\}$ ONE-HOT ENCODING

BAG-OF-WORDS: IGNORES WORD ORDER, FIXED LENGTH VECTOR OF COUNTS. $N \times V$, N IS NO DOCUMENTS, IMAGE BUT SAME. CAN HAVE MULTIPLE ROW X DOCUMENT

GOAL: BUILDING JOINT MODELS FOR $P(y_i)$, $P(N_i)$ USING LATENTS

→ CHANNELS / FEATURES / PERS

MIXTURE MODELS

ASSOCIATE SINGLE DISCRETE HIDDEN VAR PER DOCUMENT. PRIOR: $q_i \sim \text{CAT}(\pi) \rightarrow$ TOPIC. LIKELIHOOD: $P(y_{1:L} | q_i = k) = \prod_{l=1}^L \text{CAT}(y_{il} | b_{kl})$ b_{kl} IS TOPIC WORD DISTRIBUTION

- CAN HAVE DIFFERENT TOPIC MATRIX FOR EACH OUTPUT VARIABLE → 'UNSUPERVISED NAIVE BAYES'
 - UNKNOWN LENGTH → USE POISSON
 - LENGTH KNOWN → MULTINOMIAL

EXPONENTIAL FAMILY PCA: UNSUPERVISED ANALOG OF GLM

USES VECTOR OF CONTINUOUS, REAL-VALUES HIDDEN VARS.

- LATENT-SEMANTIC ANALYSIS / INDEXING**: GAUSSIAN PRIOR, GAUSSIAN LIKELIHOOD → APPLY PCA TO TERM-BY-DOCUMENT COUNT MATRIX

- CATEGORICAL PCA**: USE MULTINOMIAL / MULTINOMIAL + SOFTMAX RESPONSE

- COUNTS**: IF WE HAVE THEM MULTINOMIAL OR POISSON

- TACKLE TO FIT → DEGENERATE EM, VARIATIONAL EM, MCMC

$$P(y_{1:L,L}) = \int \left[\prod_{l=1}^L P(y_{il} | z_l, w) \right] N(z_l | \mu, \Sigma) dz_l \leftarrow \text{VISIBLE DISTRIBUTION}$$

LDA, mPCA

- WE WANT TO USE DUAL PARAMETERS OF EXPONENTIAL FAMILY. LOG ODDS → PROBABILITY VECTOR

- COUNT VECTOR WITH KNOWN SUM → mPCA, MULTINOMIAL. MARGINAL IS $P(n_i | L) = \int M(n_i | L, \theta) \text{Dir}(\theta | \alpha) d\theta$

- VARIABLE LENGTH SEQUENCE (KNOWN LENGTH) → LATENT DIRICHLET ALLOCATION **LDA**. MARGINAL: $P(y_{1:L,L}) = \prod_{l=1}^L \text{CAT}(y_{il} | \theta_{l,i})$? PROBABILISTIC WA

GAP MODEL

- DOES NOT CONSIDER COUNT VECTOR SUM TO BE OBSERVED. LATENT VARIABLES FORCED POSITIVE.

- PRIOR: $P(z^i) = \prod_{n=1}^N \text{GA}(z_{in} | \alpha_n, \beta_n)$ LIKELIHOOD: $P(n_i | z^i) = \prod_{j=1}^V \text{POI}(m_{ij} | b_{ij}^T z^i)$

- $\alpha_n = \beta_n = 0 \rightarrow$ NON-NEGATIVE MATRIX FACTORIZATION: NOT PROPERLY PROBABILISTIC

- FIXED L → REDUCES TO mPCA

LATENT DIRICHLET ALLOCATION

- EVERY WORD IS ASSIGNED ITS OWN TOPIC. FROM DOCUMENT-SPECIFIC DISTRIBUTION **AD MIXTURE MIXTURE / MIXED MEMBERSHIP MODEL**. DOCUMENT BELONGS TO ? DISTRIBUTION OVER TOPICS

$$\begin{cases} \pi_i | \alpha \sim \text{Dir}(\alpha, 1_n) \\ q_{il} | \pi_i \sim \text{CAT}(\pi_i) \\ b_{kl} | \gamma \sim \text{Dir}(\gamma, 1_V) \\ y_{il} | q_{il} = k_{il} \sim \text{CAT}(b_{kl}) \end{cases}$$

MARGINALS: $P(y_{il} = v | \pi_i) = \sum_k \pi_{ik} b_{kv}$, π_i IS TOPIC DISTRIBUTION FOR DOCUMENT.

GEOMETRICALLY IS DIMENSIONALITY REDUCTION. WE PROJECT A POINT IN V-SPACE (WORDS) TO K-SPACE (TOPICS), $V > K$

TOPIC VECTORS 'LIVE' IN WORD SPACE. ALLOWS TO DISAMBIGUATE TOPIC BY LOOKING AT OTHER WORDS

GIVEN WORD L IN DOC I WE CAN INFER LATENT TOPIC $P(q_{il} = u | y_{il}, \theta)$

- USEFUL FOR TOPIC DISCOVERY. USES TAGS TO MATCH LDA TOPICS FOR IDENTIFIABILITY

- UNIGRAMS ONLY** → NOT A VERY GOOD LANGUAGE MODEL

- PERPLEXITY**: PERFORMANCE MEASURE FOR LANG MODELS; MEASURES 'BRANCHING FACTOR' OF PROSPECTIVE DISTRIBUTION

- PERPLEXITY ($P(q)$) = $2^{H(P(q))}$, IT IS CROSS-ENTROPY. FOR UNIGRAM MODELS $H = -\frac{1}{N} \sum_{n=1}^N \sum_{q=1}^V \log q(y_{il})$. ALSO MEAN OF INVERSE PROSPECTIVE PROBS.
- FOR LDA: $P(V)$ GOTTEN BY PLUGGING IN θ , POSTERIOR MEAN, AND INTEGRATE q OUT WITH MCMC OR VARIATIONAL INFERENCE

FITTING LDA

- GIBBS SAMPLING**, COLLAPSED VARIATION: WORDS IN DOC IS ASSIGNED ON HOW OFTEN WORD IS IN TOPIC AND HOW OFTEN TOPIC IS IN DOCUMENT.

OBTAIN FULL CONDITIONALS PER TOPICS $P(q_{il} = u | q_{-i,l}, y, \alpha, \gamma)$. RANDOMLY ASSIGN TOPIC TO WORDS. THEN FOR GIVEN WORD DOCUMENT COUNTS, BASED ON OUR ASSIGNED TOPIC. ~~REPEAT~~ DRAW NEW TOPIC, UPDATE COUNTS, REPEAT.

- BATCH VARIATIONAL INFERENCE**: USE VARIATIONAL EM. **SEQUENCE VERSION**: USES VARIATIONAL MEAN FIELD. FULLY FACTORED APPROX

$$Q(\pi_i, q_i) = \text{Dir}(\pi_i | \bar{\pi}_i) \prod_l \text{CAT}(q_{il} | \bar{q}_{il}), O((E, L) V K)$$

- ONLINE VARIATIONAL INFERENCE**: *

NORMALLY DO E STEP. COMPUTE PARAMS FOR θ AS IF SINGLE DATA

WERE FULL SET; PARTIAL UPDATES AGAIN FOR θ WITH q_u WEIGHT

ON NEW AM $(1 - p_u)$ ON OLD. **SEMI-ONLINE, DO MINI-BATCHES**

→ FASTEST

- COUNT VERSION**: ONLY $O(N V K)$. WORK ON mPCA PROBLEM BECAUSE LESS STORAGE

$$Q(\pi_i, q_i) = \text{Dir}(\pi_i | \bar{\pi}_i) \prod_l M(q_{il} | n_{il}, c_{il})$$

- V.B. VERSION**: NO INSTEAD THAN EM, WE ALSO INFER PARAMETERS, ENCOURAGES SPARSITY.

$$Q(\pi_i, q_i, \beta) = \text{Dir}(\pi_i | \bar{\pi}_i) \prod_l M(q_{il} | n_{il}, \tilde{c}_{il}) \prod_u \text{Dir}(\tilde{b}_{ul} | \frac{\text{OF } \beta}{b_{ul}})$$

- PICKING K**

ANNEALED IMPORTANCE SAMPLING, CROSS-VALIDATION, VLM LOWER BOUNDS, NONPARAMETRIC METHODS

EXTENSIONS TO LDA

• CORRELATED TOPIC MODEL

FOR TI

VANILLA LDA DOESN'T DO IT BECAUSE PRIOR IS DIAGNOSTIC. REPLACE WITH LOGISTIC NORMAL, AS IN CATEGORICAL PCA. FITTING IS TRICKY BECAUSE PRIOR NO LONGER CONVEX TO LOG LIKELIHOOD. USE MEAN-FIELD OR VARIATIONAL MULTIVARIATE REG.

VISUALIZATION BY JOINING $\sum_{i=1}^K -1$ AND PRUNE LOW-STRENGTH EDGES. GET SPARSE GRAPHICAL MODEL.

• DYNAMIC TOPIC MODEL

MODELS TOPIC DISTRIBUTIONS EVOLVING OVER TIME. USE DYNAMIC LOGISTIC NORMAL MODEL. ASSUME TOPICS EVOLVE WITH GAUSSIAN RANDOM WALK THEN MAP TO PROBABILITIES VIA SOFTMAX. L WORD DISTRIBUTION OF

• LDA-HMM

COMBINES LDA AND HMM. HMM MODELS/GENERATES STOPWORDS. HAS SPECIAL STATE TO GENERATE L WORDS. LDA DOES SEMANTIC WORDS. L WORD DISTRIBUTION OF

• SUPERVISED LDA

- GENERATIVE: IE FOR SENTIMENT ANALYSIS. WE HAVE WORDS, AND ASSOCIATED LABEL. THIS LABEL GENERATED FROM TOPICS

$$P(c_i | \bar{q}_i) = \text{BER}(\text{SIGM}(W^T \bar{q}_i)) \cdot q_i \text{ IS EMPIRICAL TOPIC. FIT WITH MONTECARLO EM}$$

- DISCRIMINATIVE: THE TOPIC PRIOR IS NOW INPUT DEPENDENT $P(q_{il} | T_i, q_i = c, \theta) = \text{CAT}(A_c | T_i)$. IMAGE TAGGING

RETURN TAGS GIVEN INPUT. L COORDINATES. SIMPLEST: USE MIXTURE OF EXPERTS WITH MULTIPLE OUTPUTS

REDUCE DIAGNOSTIC WITH $\pi_i = S(Wx_i)$ DETERMINISTIC. \leftarrow MULTINOMIAL REGRESSION LDA

• OTHERS: MAKE W RANDOM VARIABLE; RANDOM FX MIXTURE OF EXPERTS, π_i DETERMINISTIC. π_i OBSERVED \rightarrow LITTLE INFLUENCE; TAG SCORES. \leftarrow MULTINOMIAL REGRESSION LDA, UNBIASED / PARTIALLY UNBIASED LDA

- DISCRIMINATIVE: EXPLAN CATEGORICAL PCA WITH INPUTS. LINEAR REGRESSION FOR IN/OUT MAPPING. \leftarrow MARKOV DYNAMICS

CATPCA FOR HID/OUT MAPPING. LIKE AN ANN WITH 1 PROBABILISTIC HIDDEN LAYER BUT WITH BOTTLENECK (HIDDEN).

LVM FOR GRAPH-STRUCTURED DATA

• STOCHASTIC BLOCK MODEL! INFER 'BLOCK' STRUCTURE OF GRAPH FROM ADJACENCY MATRIX. EVEN IF 'PITYPE' GRAPH HAS NO EDGES. PROPERTIES SUCH AS 'ALL n IN B CONNECT TO SAME NODE, ETC'. MODEL BLOCK AS LATENTS. • IF WE HAVE FEATURES FOR NODES, WE CAN MAKE DISCRIMINATIVE RELATIONAL EXTENSION OF MIXTURE OF EXPERTS

• MIXED MEMBERSHIP SBM

ALLOWS NODES TO BELONG TO MORE THAN ONE CLUSTER. AGAIN TO SOFT/FUZZY CLUSTERING. SOCIAL NETWORK ANALYSIS.

• RELATIONAL TOPIC MODEL

EXTENDS SUPERVISED LDA. WHEN NODES HAVE ATTRIBUTES. IE PAPERS, CITATIONS. WE WANT LINKS GIVEN TEXT OR VICEVERSA. LATENT: FULLTEXT.

$$P(R_{ij} = 1 | \bar{q}_i, \bar{q}_j, \theta) = \text{SIGM}(W^T (\bar{q}_i \otimes \bar{q}_j) + w_0), \text{ FORCES LATENT SPACE TO BE PREDICTIVE OF GRAPH STRUCTURE AND WORDS, UNBIASED LDA}$$

LVM FOR RELATIONAL DATA

MULTIPLE TYPES OF OBJECTS (VARS), MULTIPLE TYPE OF RELATIONS. RELATION IS TYPED. REPRESENTED WITH BINARY MATRICES. STATISTICAL RELATIONAL LEARNING.

• INFINITE RELATIONAL MODEL! EXTENDS SBM, ASSOCIATES LATENT VAR $q_i \in \{1..K\}$ WITH EACH ENTITY i OF TYPE T . PROBABILITY OF RELATION BETWEEN TWO ENTITIES BY LOOKING UP PROD OF RELATION BETWEEN TYPES. MULTIDIM.

• CAN BE USED TO INFER/LEARN ONTOLOGIES FROM DATA. • CAN BE USED TO LEARN CLUSTERS BASED ON RELATIONS AND FEATURES AT SAME TIME

• COLLABORATIVE FILTERING

NETFLIX DATA EXAMPLE \rightarrow PROBABILISTIC MATRIX FACTORIZATION. REPLACES DISCRETE LATENT VARIABLES WITH UNCONSTRAINED CONTINUOUS ONES.

USERS AND MOVIES EMBEDDED IN SAME LOW-DIM CONTINUOUS SPACES. PMF IS CLOSE TO SVD BUT MISSING DATA \rightarrow NONCONVEX OBJECTIVE

MINIMIZE NLL WITH GRADIENT DESCENT METHODS. REGULARIZE WITH GAUSSIAN PRIORS. • MOVIE/USER SPECIFIC EFFECTS MOREOVER WITH BIAS TERM, CAN MAKE IT ADAPT OVER TIME. EXPLOIT SIDE INFORMATION

LATENT DIRICHLET ALLOCATION - ADDENDUM

• FOR EACH DOCUMENT

— DRAW TOPIC DISTRIBUTION $\theta \sim \text{Dir}(\alpha)$ → 'DRAWS ME A MULTINOMIAL'

— FOR EACH WORD IN DOCUMENT

• DRAW SPECIFIC TOPIC $z \sim \text{Multinomial}(\theta)$

• DRAW WORD $w_{dn} \sim \beta_{z,dn}$ FROM β DISTRIBUTION OF WORDS IN CHOSEN TOPIC

• DOCUMENT MODEL AS CONTINUOUS MIXTURE

$$P(w|\alpha, \beta) = \int P(\theta|\alpha) \left(\prod P(w_n|\theta, \beta) \right) d\theta$$

• LDA SPACE IS ALL POSSIBLE DISTRIBUTIONS OVER WORDS

— AXES ARE WORDS

• 'KEY INFERENCE':

POSTERIOR OF HIDDEN GIVEN DOCUMENT

$$P(\theta, z | w, \alpha, \beta) = \frac{P(\theta, z, w | \alpha, \beta)}{P(w | \alpha, \beta)}$$

• HIERARCHICAL DIRICHLET PROCESS

GENERALIZATION OF LDA. USES CRP TO AUTOTUNE NO OF TOPICS FROM DATA! CLUSTER → TOPIC