

CLUSTERING

SIMILARITY BASED CLUSTERING INPUT IS $N \times N$ DISSIMILARITY / DISTANCE MATRIX. $D_{ij} \geq 0, D_{ii} = 0, D_{ij} \geq 0$

FEATURE BASED CLUSTERING INPUT IS $N \times D$ FEATURE MATRIX

• DISTANCES $\Delta(x_i, x_j) = \sum \Delta_j(x_{ij}, x_{ij})$

- SQUARED DISTANCE - CORR COEFF

- L1/CITY BLOCK DISTANCE - HAMMING DISTANCE (CATEGORICALS)

BENCHMARKS IT'S UNSUPERVISED, HARD 2 EVAL. RELY ON EXTRA DATA

• Purity: EMPIRICAL PSIA. OVER CLASS LABELS FOR CLUSTER $P = \sum \frac{N_i}{N} \cdot p_i$ TRIVIALIZED IF EACH OBJECT OWN CLUSTER

• RAND INDEX: $R = \frac{TP + TN}{TP + FP + TN + FN}$ REQUIRES REFERENCE SET. COMPARES POINT CLUSTERINGS

• MUTUAL INFORMATION $I(U;V) = \sum_{i,j} \sum_{u,v} P_{UV}(i,j) \log \frac{P_{UV}(i,j)}{P_U(i)P_V(j)}$ $P_{UV} = \frac{|U \cap V|}{N}$ $0 \leq I(U;V) \leq \min\{H(U), H(V)\}$

- NORMALIZED $NMI = \frac{I(U;V)}{(H(U) + H(V))/2}$

DIRICHLET PROCESS MIXTURE MODELS

• BEST WAY TO CHOOSE $K \Rightarrow$ NOT HAVING TO CHOOSE K INFINITE MIXTURE MODELS. NON PARAMETRIC, DIRICHLET BASED PRIOR

• IS COMPUTATIONALLY SMART • DP GIVES A DISTRIBUTION, SAMPLING FROM SAMPLE GIVES US DATAPOINT.

• DIRICHLET PROCESS DISTRIBUTION OVER PROBABILITY MEASURES. $G \sim DP(\alpha, H)$ α CONCENTRATION PARAMETER, H BASE MEASURE
MV DISTRIB IS DIRICHLET. BASE K MARGINALS ARE BETA. DP DEFINES CONJUGATE PRIOR FOR ARBITRARILY MEASURABLE SPACES

- STICK BREAKING: IF INFINITE SEQUENCE OF MIXING WEIGHTS $\Rightarrow P_n = \text{BETA}(1, \alpha)$, $\pi_n = P_n(1 - \sum_{l=1}^{n-1} \pi_l)$. $\pi \sim \text{GEM}(\alpha)$
GENERATES DISCRETE NUMBER OF EVENTS; INCREASING WITH α . $G(\theta) = \sum \pi_n \delta_{\theta_n}(\theta)$

NEED SAMPLING \rightarrow MORE REPETITIONS

DATA SAMPLED FROM \hat{G} WILL CLUSTER MORE

$G \sim DP(\alpha, H)$

• α IS DISPERSION PARAMETER

- POLYA URN / CHINESE RESTAURANT PROCESS: $\theta_1 \sim G$ ARE N OBSERVATIONS FROM $G \sim DP(\alpha, H)$. THEY HAVE K DISTINCT VALUES θ_n .

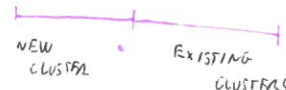
\rightarrow NEXT OBSERVATION PREDICTION $P(z_{n+1} = z | z_1:n, \alpha) = \frac{1}{\alpha + n} \left(\alpha I(z = u^*) + \sum N_n I(z = u) \right)$

• POLYA IS SAME WITH θ_n INSTEAD OF INDICATORS, ALSO ASSIGN PARAMS TO GRPS

• DISTRIBUTION OVER PARTITION OF INTEGERS

HOW TO FIT? FULL MODEL

$\pi \sim \text{GEM}(\alpha, H)$
 $z_i \sim \pi$
 $\theta_u \sim H(\lambda)$
 $x_i \sim F(\theta_{z_i})$
 G DRAWS UNOBSERVED PARAMS θ_u FROM BASE H . EACH HAS WEIGHT π_u .
 x_i GENERATED BY SAMPLING EACH OWN θ_i FROM G .
MORE DATA \rightarrow MORE LIKELY θ_i CLOSE TO θ_u ALREADY OUT.



• $N \rightarrow \infty$; $K \rightarrow \alpha \log(N)$ BECAUSE OF PRIOR

• MORE FLEXIBLE PRIOR \rightarrow PITMAN-YOR PROCESS

HOW TO FIT?

• COLLAPSED GIBBS SAMPLING: MOSTLY. BUT HAS CASE FOR $z_i = K^*$ NEW CLUSTER. • IS EFFICIENT BECAUSE CREATES EXTRA REDUNDANT CLUSTERS EARLY ON

• OTHER METHODS: STAR/BEAM SEARCH, PARTICLE FILTERING.

• HYPERPARAMS \rightarrow PUT $G(\alpha, b)$ AS PRIOR FOR α TO CONTROL NO OF CLUSTERS

AFFINITY PROPAGATION

IDEA: EACH POINT MUST CHOOSE ANOTHER DATAPOINT AS ITS EXEMPLAR/CENTROID. SOME WILL CHOOSE THEMSELVES. CLUSTERING VIA MSG PASSING

→ MINIMIZE $S(c) = \sum s(i, c) + \sum \delta u(i)$

SIMILARITY TO
CENTAURIO

PENALTY

PENALTY: -∞ IF U
DIDN'T CHOOSE ITSELF.
ELSE 0.

- OBJECTIVE FCN REPRESENTED AS FACTOR GRAPH. N NODES W/ N POSSIBLE VALUES. \rightarrow USE MAX-PRODUCT BELIEF PROPAGATION FOR FINDING MAXIMUM.
- $C_i \rightarrow \delta_i$ VAR \rightarrow FACTOR RESPONSIBILITY.
- $\delta_i \rightarrow C_i$ FACTOR \rightarrow VAR AVAILABILITY.
- NOT STABLE.
- $O(N^2)$ BUT $O(E)$ NO EDGES IN SPARSE MATRICES
- CONTROL NO OF CLUSTERS VIA DIAGONAL TERMS $S(i,i)$ HOW MUCH EACH DP WANTS TO BE EXPENSIVE.

SPECTRAL CLUSTERING

USES GRAPH CUTS. WEIGHTED UNDIRECTED GRAPH W FROM SIMILARITY MATRIX $\text{MINIMIZE } \text{CUT}(A_1, \dots, A_n) = \frac{1}{2} \sum_{i=1}^n W(A_i, \bar{A}_i)$, $\bar{A}_i = \text{COMPLEMENT OF } A_i$

→ **NORMALIZED CVT** $N_{CVT} = \frac{1 \sum_{i=1}^N CVT(A_i, \bar{A}_i)}{\sum_{i=1}^N vol(A_i)}$; $vol(A_i) = \sum_j d_{ij}$, $d_i = \sum_j w_{ij}$ **WEIGHTED DEGREE OF NODE.**

→ FORMULATION OF SEARCHING FOR BINARY VECTORS SUCH. → RELAX BINARY MEMBERSHIP CONSTRAINT, REAL VALUES. → EIGENVALUE PROBLEM

GRAPH LAPLACIAN: W IS SYMMETRIC WEIGHT MATRIX. $D = \text{diag}(D_i)$ NODE DEGREES. $L = D - W$ • EACH ROW SUMS TO 0 $\rightarrow 1$ IS EIGENVECTOR WITH $\lambda_1 = 0$

THEOREM: EIGENVECTORS OF L WITH EIGENVALUES 0 IS SPANNED BY

INDICATED VECTORS, $1_A \dots 1_n$ WHERE n ARE CONNECTED COMPONENTS OF GRAPH

SMALLEST

ALGO IDEA: COMPUTE n EIGENVECTORS u_i OF L . U EIGENVAL COLUMN MATRIX. y_i ARE ROWS. ASSIGN POINT i TO CLUSTER k IFF ROW y_i WAS ASSIGNED TO k .

↳ CLUSTER y . WITH k -MEANS IS $O(N^3)$
TO RECOVER CONNECTED COMPONENTS

IN PRACTICE: NORMALIZE L TO ACCOUNT FOR NODES MORE CONNECTED THAN OTHERS. DIFFERENT WAYS. $L_{norm} = D^{-1/2} L = I - D^{-1/2} W D^{-1/2}$

RELATION TO RANDOM WALKS: CONNECTED, NON-BIPARTITE GRAPH HAS UNIQUE STATIONARY DISTRIBUTION $\pi = (\pi_1, \dots, \pi_n)$ $\pi_i = d_i / \text{vol}(V)$ $P = D^{-1}W$

$N_{cut}(A, \bar{A}) = P(\bar{A}|A) + P(A|\bar{A})$ — WE ARE LOOKING FOR CUT RARELY TRANSITIONING $A \rightleftharpoons \bar{A}$

P is random walk
TRANSITION MATRIX

ALSO RELATED TO USUAL PCA

HIERARCHICAL CLUSTERING

- AGGLOMERATIVE / DIVISIVE.
- INPUT = DISIMILARITY MATRIX
- HEURISTIC METHODS, NO OBJECTIVE FUNCTION \rightarrow HARD TO EVALUATE

AGGLOMERATIVE: STARTS WITH N GROUPS, EACH STEP MERGES 2 MOST SIMILAR GROUPS. TOTAL $O(N^3)$ BUT $O(N^2 \log N)$ WITH TREES. COMMONLY WE RUN WARDEN

- **SINGLE LINK:** NEAREST NEIGHBOR. DISTANCE IS CLOSEST DISTANCE BETWEEN MEMBERS. MINIMUM SPANNING TREE OF DATA/POSSIBLE WEIGHTS. CAN BE $\ln(N)$
- **COMPLETE LINK:** FURTHEST NEIGHBOR. IN A LARGEST n n n n . TENDS TO PRODUCE SMALL CLUSTERS
- **AVG LINK:** AVERAGES DISTANCES BETWEEN ALL PAIRS OF OBJECTS. INTERMEDIATE BEHAVIOR. PREFERRED

DIVISIVE: STARTS WITH SINGLE CLUSTER. SPLITS IT IN 2, TOP-DOWN. HEURISTICS TO PICK OPTIMAL SPLIT. CAN BE FASTER THAN AGGREGATIVE $O(n)$. SEES ALL DATA.

- PICK n WITH LARGEST DIAMETER \rightarrow SPLIT IT WITH $n\text{-MEANS}/\text{MEDOIDS} = 2 \rightarrow$ **BISECTING $n\text{-MEANS}$**

- MIN SPANNING TREE \rightarrow BREAK LONGEST DISSIMILARITY LINK

- **DISSIMILARITY ANALYSIS** PICK MOST DISSIMILAR OBJECT AND PUT IN OTHER CLUSTER. UNTIL WE PICK 1st MAXIMIZING AVG G DISSIMILARITY / MINIMIZING H DISSIMILARITY

L FROM G

L_H

- TRICKY TO PICK OPTIMAL k , EYEBALL GAPS ON DENDROGRAM

BAYESIAN HIERARCHICAL CLUSTERING

SIMILAR TO AGGLOMERATIVE, ANALYSIS HYPOTHESIS TESTING TO PICK CLUSTERS TO MERGE • INPUT IS DATA MATRIX.

- $P(D_{ij} | T_{ij}) = P(D_{ij} | M_{ij} = 1) P(M_{ij} = 1) + P(D_{ij} | M_{ij} = 0) P(M_{ij} = 0)$ FOR EACH CLUSTER

- PICK PAIR WITH HIGHEST $R_{ij} = \frac{P(D_{ij} | M_{ij} = 1) P(M_{ij} = 1)}{P(D_{ij} | T_{ij})}$

- CONNECTED TO DIRECTLY PROCESS MM \rightarrow BHC GIVES US WORK NUM FOR DPMM MAXIMAL LINES/POS
- IS GREEDY SEARCH FOR BEST FUSE POSITION AT EACH STEP.

- THIS MATES US HAVE TO $TN = f(MK=1)$ • HYPERMANNS α, λ OF DPMN THROUGH 'BACKPOW' THROUGH THE TREE

A BHC WICKS ASS!

BICLUSTERING

CLUSTERS ROWS AND COLUMNS. WHEN WE WANT TO CLUSTER ON FEATURES. BIOINFORMATICS. COLLABORATIVE FILTERING.

ASSOCIATE ROW/ COLS TO LATENT INDICATORS AND ASSUME DATA ARE IID W/ \$M\$ SAMPLES AND FEATURES. $P(X|R, C, \theta) = \prod_i \prod_j P(x_{ij} | r_i, c_j, \theta) = P(x_{ij} | \theta_{r_i, c_j})$

\$\theta_{r,c}\$ PARAMS FOR ROW AND COL CLUSTER. • WE CAN USE DIRICHLET PROCESS VS FINITE NUMBER OF \$K\$!

MULTI-VIEW CLUSTERING

WE WANT TO MODEL DIFFERENT CLUSTERS ON THE BASIS OF DIFFERENT FEATURES. PARTITION COLUMNS INTO \$V\$ VIEWS WHERE WE 'PICK' FEATURES

DIRICHLET PROCESS FOR \$P(C)\$ SO \$V\$ VIEWS AUTOMATICALLY. THEN FOR EACH VIEW WE PARTITION ROWS \$R_{iv} = \{1 \dots K(v)\}\$ CLUSTER WHERE \$R\$ IS IN ROW \$V\$

DIRICHLET PROCESS FOR ROWS TOO. ASSUME ALL ROWS/COL IN DATA ARE IID.

• BINARY DATA \$\rightarrow\$ BETA-BERNOULLI MODEL • HYPERPARAMS: \$M\$-IT WORKS WELL • ROBUST TO IRRELEVANT FEATURES \$\rightarrow\$ THEY ARE FILTERED OUT.

DBSCAN

DENSITY-BASED SPATIAL CLUSTERING OF APPLICATION WITH NOISE • IDEA: DENSITY IN NEIGHBORHOOD HAS TO EXCEED THRESHOLD

• DIRECT - DENSITY REACHABLE. CORE/BORDER POINT. $PEN(a), |N(a)| \geq \epsilon$. • DENSITY REACHABLE: CHAIN OF DBR POINTS FROM \$P_1 \dots P_n\$ • DENSITY CONNECTED: \$P, Q\$ DBC IFF \$\exists 0\$ OR \$P, Q\$

• CLUSTER: MAXIMAL SET OF DB POINTS • NOISE: OUTLIER PTS

ALGORITHM: SAME \$\epsilon\$, MINPTS FOR ALL CLUSTERS. MANUALLY SET OR HEURISTIC: ESTIMATION VIA \$KNN\$. - GET \$N(P)\$, CHECK, EXPAND CLUSTER, REPEAT.

- ACCOMMODATIVE IN SUCCESSIVE ITERATIONS - MAY BE VASTLY IMPROVED - RESOLVES NOISE IN SUCCESSIVE ITERATIONS - WORST CASE! \$O(N^2)\$ - BEST: \$O(N \log m)\$

- SUFFERS FROM \$O(N^2)\$

BIRCH

BALANCED ITERATIVE REDUCING AND CLUSTERING USING HIERARCHIES. IS ONLINE. CAN WORK IN \$O(N)\$. CLUSTERING FEATURE: ALSO MAINTAINS NO OF POINTS, SUM, SQUARES SUM OF CLUSTERS

• CF-TREE: A TREE STRUCTURE MAINTAINING THE CFs. PARENT-CHILD RELATIONSHIPS. BINARY TREE. INSERTION/SPLIT FOR ADDITION. LIMITED STRUCTURE. REFINING STOPS TO CONSOLIDATE STRUCTURE. ALGORITHM: SCAN DATA, PRUNE / REMOVE OUTLIERS, CONSOLIDATE / GLOBAL CLUSTERING, REFINED. AGGLOMERATIVE HIERARCHICAL CLUSTERING

• DIFFERENT SIMILANCE METRICS. \$\hookrightarrow\$ AGGLOMERATE \$\hookrightarrow\$ CLUSTERING ALSO \$\hookrightarrow\$ RECOMPUTE CENTROIDS

• IS MORE OF A FRAMEWORK FOR EFFICIENT CLUSTERING UNDER MEMORY CONSTRAINT THAN AN ALGORITHM.

T-SNE

USED FOR DIM/REDUCTION & VISUALIZATION. PRESERVES LOCAL AND GLOBAL STRUCTURE IN LOW-DIM SPACE. T-DISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING

VANILLA SNE: SIMILARITY OF \$x_i, x_j\$ IS CONDITIONAL PROB I WOULD PICK \$j\$ AS NEIGHBOR UNDER GAUSSIAN CENTROIDS ON \$x_i\$.

\$P_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma^2)}\$ - \$Q_{j|i} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq i} \exp(-\|y_i - y_k\|^2)}\$ IDEA: MINIMIZE MISMATCH BETWEEN \$P_{j|i}\$ AND \$Q_{j|i}\$ OVER ALL DATAPOINTS VIA GRADIENT DESCENT

COST FCN: \$C = \sum_i KL(P_i || Q_i) = \sum_i \sum_j P_{j|i} \cdot \log \frac{P_{j|i}}{Q_{j|i}}\$ FIM: \$\sigma\$ TO MINIMIZE \$P_i\$ WITH FIXED, SPECIFIED PERPLEXITY \$P(i) = 2^{H(P)}\$, ENTROPY

GRADIENT: \$\frac{\partial C}{\partial y_i} = 2 \sum_j (P_{j|i} - Q_{j|i} + P_{i|i} + Q_{i|i}) (y_i - y_j)\$ INIT: RANDOM SAMPLES AROUND ORIGIN. GRADDESC W/ MOMENTUM. NUMERICAL W/ ARTIFICIAL NOISE REDUCED OVER TIME.

SYMMETRIC SNE: ASSUME \$P_{i|j} = P_{j|i}, Q_{i|j} = Q_{j|i} \forall i, j\$. \$\rightarrow C = KL(P || Q) = \sum_i \sum_j P_{i|j} \cdot \log \frac{P_{i|j}}{Q_{i|j}} \rightarrow \frac{\partial C}{\partial y_i} = 4 \sum_j (P_{i|j} - Q_{i|j}) (y_i - y_j)\$

ISSUE: CROWDING PROBLEM: AREA FOR DISTANT POINTS NOT ENOUGH COMPARED TO THAT FOR CLOSE POINTS \$\rightarrow\$ IN LOW DIM \$\rightarrow\$ USE A STUDENT'S \$T\$ IN LOW DIM

T-SNE: \$Q_{i|j} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|y_i - y_k\|^2)^{-1}}\$ • \$T\$ IS 1-DOF CAUCHY. NUMERATOR \$\approx\$ INVERSE SQUARE LAW. • INVARIANT OF CHANGE OF SCALE. • LARGE CLUSTERS INTERACT LIKE INDIVIDUAL POINTS. LONG-RANGE FORCES

\$\frac{\partial C}{\partial y_i} = 4 \sum_j (P_{i|j} - Q_{i|j}) (y_i - y_j) (1 + \|y_i - y_j\|^2)^{-1}\$ • EFFICIENT TO COMPUTE

OPTIMIZATION: MOMENTUM IN FIRST ITERATIONS. \$L_2\$ PENALTY TO COST FCN; PROPORTIONAL TO SQUARES, MAP INITIALLY CLOSE. ON EXAGGERATE HIGH-DIM DISTANCES AT START. UNFILLS EMPTY SPACE IN LOW-DIM. BETTER VIZ W/ ETL. • INITIALLY REDUCE DATA WITH PCA TO \$D \le 30\$ TO AVOID TOO MUCH OVERFIT

• 70% POINTS: REDUCE COMPLEXITY VIA NEIGHBORHOOD GRAPH FOR ALL POINTS. RANDOM WALKS. TO PICK \$P_{i|j}\$ • UNKNOWN POINTS.