

## Erdős and Renyi (ER)

- randomly choose n out of  $N(N-1)/2$  possible vertices
- Parameter p: each of the  $N(N-1)/2$  possible connection is chosen with prob. p. Expected number of edges is  $N(N-1)/2 * p$
- Degree Distribution: Poisson
- diameter:  $\ln(N)$
- clustering coef:  $1/N$
- homogeneous networks influenced by variance of Poission variance

## Barabasi/Albert (BA)

- starting with a small number, vertices are added with m connections each
- Prob of edge depends on degree of vertex:  $k_i / \sum_j k_j$
- Degree Dist.: Power law
- diameter:  $\ln(N) / \ln\ln(N)$
- clustering coef:  $1/N^{0.75}$
- heterogeneous graphs since power law is long-tail distribution with large var
- scale free, hubs/cliques, small world

## PCA

- find the linear projection of the data which preserves as much variation in the data as possible
- Invariant against orthogormal operation
- PCA, ICA, LDA are all deterministic, work on vectors and are linear parametric projections

## ICA

- find linear projection such that output features are as independent as possible
- ICA regularizes the problem by aiming for such sources s which are maximum independent
- ICA cannot recover sources which are all Gaussian since linear mixtures of Gaussians remain Gaussian, cannot recover exact scaling of the sources and their order

## LDA

- given labels, find linear projection such that variance of points with same label is minimized, variance of points in between different labels is maximized
- Maximum dim. of the LDA projected space is C-1, C being number of classes

- BA more suited for web, since it reflects common phenomena like presence of hubs, the scale-free property, and small-world networks
- Tldr:
  - ER generates random graphs where edges are formed independently with a fixed prob. p between any pair of nodes
  - BA creates creates scale-free networks using preferential attachment, where new nodes are more likely to connect to already well-connected nodes (hubs)

Poisson dist: probability of x events in a given time if average rate lambda is known

$$P(x) = \exp(-\lambda) \cdot \lambda^x / x!$$

Expectation and variance lambda, both finite

Power law:  $E(x^m) = c / (\gamma - m - 1)$  long tails, probably expectation finite for  $\gamma > 2$  infinite var

$$P(x) \sim c \cdot x^{-\gamma}$$

variance finite for  $\gamma > 3$

for every vertex v pagerank  $r(v)$  should be high if pointed to by many sites with high pagerank and selective links

$$r(v) \sim \sum_{w \in pa(v)} \frac{r(w)}{|ch(w)|}$$

$$r := \epsilon/n + (1 - \epsilon) \cdot B \cdot r$$

$$r := (\epsilon/n \cdot \mathbf{1} + (1 - \epsilon) \cdot B)r$$

Introduce random jumps vector  $r$ , with probability epsilon = 0.15

- GPT = generative pretrained transformer
- generative: produces text via sampling from next word probability, which is modelled by means of Markov Model (with many parameters)
- transformer: specific neural architecture, which models the probabilities; acts on time window with attention heads
- pretrained: parameters are trained huge corpora from the web via next word prediction + human-feedback reinforcement learning

Models the liklihood of a random web surfer landing on a page, combining contributions from linked pages ( $B \cdot r$ ) and random jumps ( $\text{epsilon} / N$ ) to ensure stability and handle disconnected components

lower bound: given by the first summand ( $\text{epsilon} / N$ ) plus normalization  
upper bound is N number of nodes if all mass is centered on the page

Pages without incoming links (parents) rely only on the random jump term, so they have a low PageRank. Pages with many parents have a higher PageRank because they receive contributions from multiple terms, depending on their parents' PageRank.

Normalized cut decreases as the number of clusters increases, which can lead to trivial solutions

$$\text{cut} \quad \text{cut}(C_1, \dots, C_k) := \frac{1}{2} \sum_{i=1}^k W(C_i, V \setminus C_i)$$

$$\text{ratio cut} \quad \text{ratio cut}(C_1, \dots, C_k) := \frac{1}{2} \sum_{i=1}^k W(C_i, V \setminus C_i) / |C_i|$$

$$\text{normalized cut} \quad \text{ncut}(C_1, \dots, C_k) := \frac{1}{2} \sum_{i=1}^k W(C_i, V \setminus C_i) / \text{vol}(C_i)$$

$$\text{where } \text{vol}(A) = \sum_{i \in A} d_i \text{ and } d_i = \sum_j w_{ij} \text{ degree of } i$$

CUT should not be used for clustering, as it favors small imbalanced clusters, NCUT more reasonable ensures that clusters are balanced and meaningful

One cluster:

- Cut yields 0, QE yields sum of squared dist. of points to the centre
- N clusters: CUT yields sum of weights, QE yields 0 when choosing prototypes as data points itself

Explain objectives of ratio cut (CUT) and quant. error (QE), NCUT (Normalized Cut)

- NCUT: addresses issue by dividing the edge weights by the total degree of each cluster, balancing the cluster. It ensures that both clusters are meaningful and of reasonable size
- ratio cut: given weighted graph, decompose weighted graph into clusters such that the sum of weights which are pointing outside a cluster normalized by the size of the cluster is minimized
- quantization error: given vectors, find prototypes representing clusters and an assignment of data points to these prototypes such that the sum of (squared) distances of data points to their prototypes is minimum

Method	vector	dissimilarity	deterministic	parametric
LLE	X		X	
UMAP	X	X		
PCA	X		X	X

- Parametric
- PCA:  $x \rightarrow Wx$
  - SOM:  $x \rightarrow$  position of closest w in a lattice
  - GTM:  $x \rightarrow$  most likely position in latent space
  - autoencoder:  $x \rightarrow f(x) \rightarrow f^{-1}f(x)$
- Non-parametric
- MDS, SNE, t-SNE, MVU, LLE, Isomap, Sammon, CCA, Laplacian Eigenmap
- Nearest neighbour desc.
- efficient algo to approx. nearest neighbours
  - iteratively improve initial random or approx. neighbor graph by focusing on pairs of points that are likely to yield better neighbours

## t-SNE:

- goal: find low-dimensional vectors such that the original dist. of pairwise points is preserved as much as possible
- match pairwise dist in original space and pairwise dist in projection space
- take gaussian in orginal space and student-t in projection space
- pros: moderate original distance can be modelled by much larger distance in projection, accuracy of large distances in projection space is not so important
- optimize  $KL(p_{ij} || q_{ij})$  with gradient descent

## Laplacian Eigenmap:

- deterministic
- local graph
- weighting: one for every link or heat kernel
- compute degree matrix
- compute graph laplacian
- embed via generalized eigenvectors

## MDS:

- find vector to preserve distances
- Constructs a configuration of points in the desired dim.
- Minimizes the difference between the original distances and the distances in the low-dim space
- Assumes linear relationships between distances

## Isomap

- input: vectors or pairwise distances
- parameters: k or epsilon for neighbourhood computation
- output: low-dimensional projection of points
- objective: preserve geodesic distances on data manifold
- Algo: Compute neighbourhood graph, compute shortest distances on neighbourhood graph, project preserving distances using Multidimensional Scaling
- effort: dominated by shortest paths in sparse graphs hence  $N^2 \log N$ , rest is  $N^2$
- is non-linear
- goal: find real vectors  $x_i$  such that the original manifold distances are preserved as much as possible
- works if intrinsic data manifold is two-dim and flat (swiss roll) fails if sampling is not sufficient to get proper neighbourhood graph
- non-parametric hence does not directly give you out-of-sample extension.

## Locally linear embedding (LLE)

- given: real vectors
- goal: find low-dimensional vectors such that the original local data structure is preserved as much as possible
- step 1: define local neighbourhood
- step 2: find linear relation in neighbourhood
- local structure should be preserved
- step 3: find mapped points such that the linear relationships are similar
- is a deterministic dimensionality reduction method

## UMAP

- tries to preserve local topological structure
- compute knn graph, weight connections by a term, cost given by cross-entropy

## Measure Reduction Dim. Methods

**Entropy:** Measures randomness or disorder of a variable. Max. entropy means the data is evenly distributed. For a fixed variance, entropy is highest if the distribution is Gaussian

**Relevance:** Evaluate structure of complex data, determine rather data follows gaussian dist.

**Negentropy:** Measures deviation of a dist. from a Gaussian dist. It is calculated as the difference between the entropy of a Gaussian dist. and the entropy of the actual data

- Used to measure non-Gaussianity

**Mutual Information:** Measure degree of dependence between multiple variables. If two variables are independent, their mutual information is 0.  $I(t_1, \dots, t_n) = \sum_i H(t_i) - H(t_1, \dots, t_n)$

MI and negentropy are inversely corr. to any source, only if data is whitened

## Prototype-based Models

	LVQ1	LVQ2.1	GLVQ	RSLVQ
Cost Func.	no	Yes but unbounded	Yes: correlates to margin optimization	Yes: optimize posterior probability for GMM

GLVQ	
	$\sum \Phi \left( \frac{d(\vec{x}_i - \vec{w}_+) - d(\vec{x}_i - \vec{w}_-)}{d(\vec{x}_i - \vec{w}_+) + d(\vec{x}_i - \vec{w}_-)} \right)$
Max. Margin with Loss function	
LVQ is supervised learning	
Shape:	
<ul style="list-style-type: none"> <li>GMLVQ: elipsoidal shape, with arbitrary rotated axes</li> <li>GLVQ: round shape, standard unit sphere</li> <li>GRLVQ: elipsoidal shape, with axis-aligned symmetry-achsces</li> <li>LVQ: round shape</li> </ul>	
LVQ, GLVQ, GMLVQ are convex	
GMLVQ: $\Lambda$ should be positive definite, can be achieved through reparametrization: $\Lambda = \Omega^* \Omega^T$	
GLVQ uses euclidian <sup>2</sup> and GMLVQ uses adaptive squared euclidian	

Metric parameters	
Relevance learning (GRLVQ)	Adjust feature weighting only
Metric learning (GMLVQ)	Adjust quadratic form via adaptive linear transformation of data
Local (IGLVQ or LGRLVQ)	Adjust individual parameters per receptive field
Low rank (LIRALVQ or use eigenvectors of matrix in GMLVQ)	Use a low rank matrix for the data transformation
Robust soft LVQ (RSLVQ)	
<p>data are given as Gaussian mixture over the classes:  <math>p(\vec{x} W) = \sum_c \sum_{c(\vec{w}_j)=c} p(\vec{x} \vec{w}_j) P(\vec{w}_j)</math>  <math>p(\vec{x} \vec{w}_j) = (\sigma/2\pi)^{n/2} \exp(-\sigma/2\ \vec{x} - \vec{w}_j\ ^2)</math>  <math>P(\vec{w}_j)</math> uniform distribution</p> <p>probability of data labeled by <math>c</math>:  <math>p(\vec{x}_i, y_i W) = \sum_{c(\vec{w}_j)=y_i} p(\vec{x}_i \vec{w}_j) P(\vec{w}_j)</math></p> <p>optimize the ratio Log-Likelihood-Ratio  <math>\sum_i \log(p(\vec{x}_i, y_i W)/p(\vec{x}_i W))</math></p> <p>benefit: the term is bounded → no divergence</p>	

init $\alpha_j$ repeat compute $d_{i,j} = [D \cdot \alpha_j]_i - 0.5 \cdot \alpha_j^T D \alpha_j$ determine receptive field $R_j = \{x_i \mid d_{i,j} \leq d_{i,k} \forall k\}$ determine coefficients $[\alpha_i]_j = 1_{x_i \in R_j} /  R_j $	LVQ1: yields good results in typical model situations
	LVQ2.1: without window does not, diverges, strange behaviour when using Early Stopping

$d^+ := d(\vec{x}_i, \vec{w}^+)$ closest prototype with label = $y^i$ $d^- := d(\vec{x}_i, \vec{w}^-)$ closest prototype with label $\neq y^i$	GLVQ: Has linear boundary and the two receptive fields are convex
	GMLVQ: Similiar since choosing $\Lambda = \Omega^* \Omega^T$ . Linear transformation preserves the shape of a linear classification boundary and receptive fields are convex
	LGMVQ: Boundary has quadratic form, receptive fields can be seperated by nonlinear boundaries, might not be convex
	Example of not convex receptive fields, which both have two connected components 'two cigar' example: $w_1 = w_2 = (0,0)^T$ , $\Lambda_1 = (1,1)^T (1,1)$ , $\Lambda_2 = (-1,-1)^T (-1,-1)$

	k-means / NG	spectral	affinity propagation	relational NG		silhouette	Davis Bouldin	Calinski Harabasz	gap		method	non i.i.d.	fixed part of D
objective	quantization error	(normalized /ratio) graph cut	quantization error for exemplars	quantization error in kernel space	range	[-1,1]	>0	>0	usually >0	Nyström	matrix based techniques: • spectral clustering • relational NG	no	yes
data	vectors	(sparse) graph with similarities	(sparse) similarity matrix	dissimilarity matrix	which k	large value	small value	large value	smallest k with large value	patch	prototype based techniques: • relational NG • AP	yes	no
representation	prototypes	cluster indices	exemplars	prototypes as implicit linear combination	k=1	no	no	no	yes				
OOS	yes	Needs additional effort	yes	yes	data type	every distance	vectors	vectors	vectors				
parameters	cluster prototypes and assignments to cluster	assignments and assignments and number of clusters	exemplars and assignments and number of clusters										

Methods for large data sets or data streams:  
BIRCH (balanced iterative reducing and clustering using hierarchies), Path neural gas, Patch relational neural gas, Patch affinity propagation

k-means		Spectral Clustering	Affinity propagation	Neural Gas																	
goal: find prototypes so they min quantization error:		<ul style="list-style-type: none"> <li>SC requires non-negative pairwise similarities as input, missing values are ok and implicitly treated as zero.</li> <li>Vectorial data need to be transformed using e.g. neighborhood graphs. The number of clusters k is an input.</li> <li>Cannot handle vectorial data</li> </ul>		<ul style="list-style-type: none"> <li>Treats missing info as -inf</li> <li>does not need number of clusters</li> <li>uses self similarities</li> <li>New cluster is build whenever the costs induced by the self-similarity of the clusters exemplar are smaller than the costs when using a smaller number of clusters measured in terms of the quant. error</li> </ul>																	
$E = \frac{1}{2} \sum_{\vec{w}^i} \sum_{\vec{x}^j \in R(\vec{w}^i)} \ \vec{w}^i - \vec{x}^j\ ^2$																					
k-means++: init centroids in better way start with one center x take next center x with probability proportional to D(x) where D(x) = distance from already chosen centers repeat until k prototypes are chosen																					
perform k-means <ul style="list-style-type: none"> <li>unsupervised learning</li> <li>distance-based</li> </ul>				<ul style="list-style-type: none"> <li>Algorithm that adapts a set of prototype vectors to represent the strcutre of the input data by min. a distortion measure. Prototypes updated in a neighbourhood-based manner, with the influence of updates decreasing over time and distance from the winning prototypes</li> </ul>																	
First step: each prototype $w_j$ is implicitly represented via a N-dimensional vector $\alpha_j$ of coefficients representing $w = \sum_i [\alpha_j]_i x_i$ , whereby $[\alpha_j]_i \geq 0$ and $\sum_i [\alpha_j]_i = 1$ . This allows a computation of (squared) distances: $d(x_i, w_j) = [D \cdot \alpha_j]_i - 0.5 \cdot \alpha_j^T D \alpha_j$ . Relational k-means clustering results by entering this into the vectorial algorithm:		<table border="1"> <thead> <tr> <th>Method</th> <th>cost function</th> <th>parameters</th> <th>out-of-sample</th> </tr> </thead> <tbody> <tr> <td>spectral clustering</td> <td>normalized/ratio cut</td> <td>assignments to cluster</td> <td>no</td> </tr> <tr> <td>kmeans++</td> <td>quantization error</td> <td>cluster prototypes and assignments</td> <td>yes</td> </tr> <tr> <td>affinity propagation</td> <td>(exemplar based) quantization error</td> <td>exemplars and assignments and number of clusters</td> <td>yes</td> </tr> </tbody> </table>		Method	cost function	parameters	out-of-sample	spectral clustering	normalized/ratio cut	assignments to cluster	no	kmeans++	quantization error	cluster prototypes and assignments	yes	affinity propagation	(exemplar based) quantization error	exemplars and assignments and number of clusters	yes		
Method	cost function	parameters	out-of-sample																		
spectral clustering	normalized/ratio cut	assignments to cluster	no																		
kmeans++	quantization error	cluster prototypes and assignments	yes																		
affinity propagation	(exemplar based) quantization error	exemplars and assignments and number of clusters	yes																		

Gap statistics		Silhouette
Method to determine number of clusters deviation of clustering coefficient from expected value for baseline		scaled difference of distances between clusters and within clusters. Should be large
$E(\log W_k) - \log W_k$ pick smallest k where this is significantly larger than 0		$S_i = (b_i - a_i) / \max(a_i, b_i)$ where $a_i = \text{average distance of point } i \text{ to points in its cluster}$ $b_i = \min_j \text{ average distance of point } i \text{ to points in cluster } j$ average over all points i, centered at 0, no reasonable value for k=1
$W_k = \sum_{i=1}^k d_{ii} / (2n_i)$ $n_i = \text{number of points in cluster } i$ $d_{ii} = \text{sum of pairwise distances in cluster } i$ $E$ refers to expectation over reference distribution		
Calinski Harabasz		Davis Bouldin
ratio between cluster variance and within cluster variance		ratio within cluster and between cluster distances, should be small
$\frac{\sum_{i=1}^k n_i (m_i - m)^2}{\sum_{i=1}^k n_i (m_i - m)^2}$ where $m_i = \text{mean of cluster } i$ , $m = \text{data mean}$ , $n_i$ size cluster $i$		$\sum_{i=1}^k \max_{j \neq i} D_{ij} / k$ where $D_{ij} = (\bar{d}_i + \bar{d}_j) / \bar{d}_{ij}$ $\bar{d}_i = \text{average distance of points in cluster } i \text{ to cluster center}$ $\bar{d}_{ij} = \text{distance of cluster centers } i \text{ and } j$ normalized by k, nonnegative values not reasonable value for k=1
external/internal clust. measures		External: Adjusted Rand index (ARI) measures the similarity between pred. clusters and ground truth by counting matching pairs of points. Not suited when no labeled data
Internal: Silhouette Score: See up. Not suited for clusters of varying density or non-spherical shapes		Internal: Silhouette Score: See up. Not suited for clusters of varying density or non-spherical shapes

Challenges when learning on data streams:		
• Label information, what are good learning signals?		
• Algorithmic challenge, how to efficiently update model for new data points		
• Model selecteion challenge, how to adjust model complex. and hyper para		
• Efficient memory models, how to store required information		
Method	Dealing with drift	Model
ISVM	no	Kernel method
Learn++	Implicit virtual drift	Ensemble
VFDT	Window	Decision tree
JIT	Drift detection	Generic concept (eg kNN)
LVGB	Drift detection	Ensemble
ARF	Drift warning and detection	Ensemble of trees
SAM	ITTE based adaptive window	kNN
SAME	Drift detection and ITTE based adaptive window	Ensemble of kNN
Name	Type	Method
DDM	Error-based	Statistics of error
ADWIN	Error-based	Significant distance of mean error
HDDDM	Distribution-based	Hellinger distance for bins
PCA-CD	Distribution-based	Difference of PCA projections
DAWIDD	Time-dependency	Change of time characteristics in space
Shape method	Time-dependent shape fit	Detects rapid drift