







- All A

X

-200

Two major types of clustering

- **Partition-based clustering**: divide items into groups so that items within a group are similar (close) and items from different groups are less similar (more distant)
- · Examples: k-means, self-organizing map, hierarchical
- Property of the result: each item belongs to some group
- Density-based clustering: find groups of highly similar (close) items and separate from them items that are less similar (more distant) to others
- Examples: DBScan, OPTICS
- Properties of the results: some items belong to groups, other items remain ungrouped and are treated as "noise"



Use of the two types of clustering

- Partition-based:
- Typically applied to multiple thematic (non-spatial) attributes or to time series of thematic attributes
- Objective: divide objects into groups such that objects within a group have similar attribute values and differ from the objects in the other groups
- Density-based:
 - Typically applied to spatial and temporal attributes of spatial or spatiotemporal objects
 - Objective: find concentrations of objects in space or in space and time (i.e., groups of objects with close spatial locations and existence times)
 - concentrations of objects may have special meanings; e.g., spatio-temporal cluster of low speed events \Rightarrow traffic jam

Density-based clustering (DBC)

Goal: find dense groups of close or similar objects

For a given object *o*, the objects whose distances from *o* are within a chosen distance threshold (radius) R are called <u>neighbours</u> of the object *o*.

R

- An object is treated as a <u>core</u> object of a cluster if it has at least N neighbours.
- To make a cluster:
- some core object with all its neighbours is taken;
 for each core object already included in the cluster, all its neighbours are
- also added to the cluster (if not added yet).
- Some objects may remain out of any cluster (when they have not enough neighbours and do not belong to the neighbourhood of any core object). These objects are treated as "<u>noise</u>".

Density-based clustering

Parameters

- For DBC, the user needs to specify the neighbourhood radius (distance threshold) R and the minimum number of neighbours N.
 ⇒ The use of DBC requires an understandable definition of distance between objects, e.g., spatial distance or spatio-temporal distance.
 - It may be hard to choose R for a more abstract "distance" between combinations of values of multiple diverse attributes.
- · Results of DBC greatly depend on the parameter choice.
- Visualisation and interactive exploration help the analyst to find suitable values for R and N that lead to good results.







X

DBC by spatio-temporal distances

- + For any two objects, there is a distance in space $d_{\rm space}$ and a distance in time $d_{\rm time}.$
- To cluster the objects by their spatio-temporal proximity, the analyst may choose two neighbourhood radii R_{space} and R_{time} e.g., R_{space} = 300 m and R_{time} = 30 minutes.
- However, the clustering algorithm requires a single measure of distance and a single radius.
- \Rightarrow Spatial and temporal distances need to be combined together \bullet e.g., d = max(d_{space}/R_{space}, d_{time}/R_{time}) * R_{space}



















 \sim

X



Distance functions in DBC

- · Elementary distances: spatial, temporal, difference of values of a single thematic attribute
- · It may be necessary to group objects on the basis of two or more elementary distances, e.g., spatial and temporal \Rightarrow A <u>distance function</u> integrating the elementary distances is needed
- · General approach:
 - 1) Set a separate threshold for each elementary distance
- 2) Transform the absolute elementary distances to relative w.r.t. the respective thresholds
- 3) Combine the relative distances:
- take their maximum or compute the Euclidean or Manhattan distance
- Defining more complex distance functions is also possible
- May be needed for complex objects, such as trajectories

Investigation of parameter impact

- · The results of DBC greatly depend on the parameter settings (values of R and N)
- \Rightarrow It is necessary to run the clustering tool multiple times with different parameter settings
- · Choose clear, easily interpretable results
- · Results from different runs may complement each other and contribute to better understanding
- Interactive visual interfaces are used for investigating the results of different runs.

X

X

Visual investigation of DBC results

- · Analogously to PBC, clusters are given distinct colours, which are used for colouring marks in a map, space-time cube, various graphs and plots, segments in histograms, ...
- · Noise is usually shown in grey
- · The analyst should be able to interactively hide and unhide the noise
- · Problems in visualising density-based clusters:
- [®]DBC may produce more clusters than there are distinguishable colours The analyst should not rely too much on cluster colours but use them mainly for distinguishing clusters from nois

® Visual displays showing individual cluster members may be too cluttered \Rightarrow The clusters need to be represented in a summarized form

· E.g., by spatial or spatio-temporal convex hulls

\sim Reading • IEEE VAST 2011 paper (best paper award): G.Andrienko, N.Andrienko, C.Hurter, S.Rinzivillo, S.Wrobel From Movement Tracks through Events to Places: Extracting and Characterizing Significant Places from Mobility Data IEEE Visual Analytics Science and Technology (VAST 2011). Proceedings, IEEE Computer Society Press, 183-192 Extended version, covering also scalable clustering of events: G.Andrienko, N.Andrienko, C.Hurter, S.Rinzivillo, S.Wrobel Scalable Analysis of Movement Data for Extracting and Exploring Significant Places IEEE Transactions on Visualization and Computer Graphics, 2013, 19(7), 1078-1094 http://dx.doi.org/10.1109/TVCG.2012.311



X

Origin-destination movement data

Shortly: OD moves

- · Movement data specify spatial positions of moving objects at different times
- · Origin-destination data: specify only the positions and times of trip starts and ends; intermediate positions are not available.

A running data example

Use of shared bikes in London

- · Barclays Cycle Hire
- · 569 docking stations to pick up and return bicycles
- The data are publicly available from Transport for London (http://www.tfl.gov.uk/info-for/open-data-users/our-feeds)
- Data describe the journeys made with the rented bikes:
- Journey ID
- Bike ID
- Start date & time • End date & time
- · Start docking station ID · End docking station ID
- The geographic coordinates of the docking stations are known and can be joined with the transaction data



X









Relation of OD moves to spatial events

• Spatial events are objects having positions in space and times of occurrence or existence (= positions in time).

- Space and time (and space \times time) can be treated as containers of events.
- Analysis is generally concerned with the spatio-temporal distribution of events or values of their thematic attributes.

× H

- An OD move includes 2 instant spatial events: start and end.
 - The spatio-temporal distributions of these events may be of interest in analysis.
 These events can be analysed with the same aims and the same methods as any other spatial events.



Example: instant point events extracted from the London bike trip data from Wednesday, 25/07/2012. The upper STC image shows the spatio-temporal distribution of the trip start events and the lower image represents the trip end events. There are noticeable differences on the east, the other parts of the territory need to be explored with the help of spatial filtering, to decrease the display clutter and overplotting.

-200



Relation of OD moves to spatial events

(continued)

- An OD move as a whole is a spatial event.
- It has existence time = the time of the trip.
- It has a position in space consisting of the start and end positions.
- Such discontinuous spatial positions also occur for other spatial objects, e.g., a university.
- The spatial position of a move can be represented by a directed line segment (vector), which is a spatial object.
- However, the visualisation and analysis methods designed for spatial events may not be applicable to this kind of events.
 - Most methods assume that the spatial positions of the events do not change during the time of event existence.
 - Moreover, most of them assume that the spatial positions of the events can be represented as points in space.
- Hence, OD moves require special approaches.





General approaches to dealing with large data and display clutter

- Aggregation: divide attribute domains into bins; obtain counts of objects and summaries of their thematic attributes for the bins; analyse the distributions of the aggregates.
- $\bullet\,$ In particular, aggregate spatio-temporal objects by areas and time steps.
- **Partition-based clustering**: divide objects into groups by similarity or closeness; analyse and compare group summaries (aggregates) and internal variations.
 - In particular, cluster spatio-temporal objects by their positions in space and time.
 - · PBC is efficient when a relatively small number of clusters is sufficient.
- Density-based clustering: find dense groups (concentrations) of similar or close objects; analyse their relations to the remainder.
 In particular, find spatio-temporal concentrations of spatio-temporal
 - objects; analyse where and when they occurred.



DBC: what distance function to use?

- DBC requires setting a distance threshold \Rightarrow there should be a meaningful notion of the distance between objects.
- The **spatial distance** between OD moves can be defined as the mean of the spatial distances between the origins and between the destinations:
- + $s_distance(m_1,m_2) = (s_distance(o_1,o_2) + s_distance(d_1,d_2)) / 2$
- Analogously, the **temporal distance** may be defined as
- + $t_{distance(m_1,m_2)} = (t_{distance(o_1,o_2)} + t_{distance(d_1,d_2)}) / 2$
- The **spatio-temporal distance** between OD moves can be defined taking the same approach as for instant spatial events:
- + Spatial distance threshold (radius) $\rm R_{S}$ + temporal distance threshold $\rm R_{T}$
- + $s_t_distance = \max(s_distance/R_s + t_distance/R_T)*R_s$

47



- A spatio-temporal cluster means collective movement of multiple objects between some origin and destination regions.
- Further questions: How frequent are the occurrences of collective movement? Based on their positions in space and time, how can they be interpreted?

DBC by the spatial distance (examples)

25/07/2012 (Wednesday); round trips excluded. Distance threshold R = 200 m; minimal number of neighbours N = 10. Result: 93 spatial clusters with sizes from 10 to 63 include 1,618 trips (4%); the noise consists of 38,119 trips (96%). Hence, not many very similar trips occurred throughout the day.

























What we have learned so far

- · High variety of trips; not many similar trips.
- A large number of similar trips from Waterloo in the morning and to Waterloo in the evening.

 \sim

- These may be trips of commuters who travel to London by railway and then get to the final destinations by bike.
- Existence of opposite groups of trips.
- Some of these may be trips of people going to their work or study places in the morning and back in the evening.
- Existence of groups of similar trips occurring close in time.
 Most such groups occur in the morning or in the evening and go from or
- Most such groups occur in the morning or in the evening and go from or to areas around railway stations.























Time division possibilities

- Treat the time as a directed <u>line</u>.
- Divide the range from the earliest to the latest time value in the data into consecutive intervals.
- To aggregate: for each time value in the data, find the containing interval.
- Treat the time as a <u>cycle</u>.
- Choose a relevant cycle for the data: daily, weekly, annual, domain-specific, e.g., production cycle
- Divide the chosen <u>cycle</u> into intervals.
- To aggregate: for each time value in the data, determine its position in the cycle and find the containing interval of this position.
- I.e., the absolute time value is transformed to relative w.r.t. the time cycle.

Aggregation result: spatial time series

• Data with two referrers: a set of space compartments and a set of time intervals

.

- Can be viewed in two complementary ways:
 as a set of time series of attribute values associated with the spatial compartments
- as a set of *spatial situations* associated with the time steps
- A spatial situation is a distribution of attribute values over the space in some time step

X

× S

X

Analysis of spatial time series (a reminder)

- Analysis tasks address two aspects of the overall behaviour:Spatial distribution of the local temporal variations of the attribute
- values in different compartments • Temporal variation of the overall spatial distribution of the attribute
- values in different time moments
- Supporting visual analytics techniques (considered in the previous lecture) include two-way partition-based clustering
- Grouping of places (compartments) by similarity of the local time series
 followed by visual exploration of the distribution of the group members over the space
- Grouping of time steps by similarity of the spatial situations
- followed by visual exploration of the distribution of the group members over time

















Spatial link-referenced time series

- We transformed the complex data structure $S\times S\times T\to A$ into a simpler data structure $L\times T\to A$

X

- + L is the set of spatial links defined as pairs (s_o, s_d), s_o \in S, s_d \in S.
- Hence, we obtained "normal" space-referenced time series, in which attribute values refer to spatial objects (specifically, the spatial links) and time steps.
- \Rightarrow We can apply the same visual analytics techniques as for usual time series
 - \dots but we need to deal with the difficulties in the visual representation of the spatial links in a map
 - \otimes This is a visualisation problem that is not completely solved yet.



















X

Analysis of spatial events and OD moves

Complementarity of DBC and spatio-temporal aggregation

X

X

- · Task: analyse the spatio-temporal distribution of events/moves
 - · Space and time are treated as containers of the events/moves
- · Spatio-temporal aggregation:
 - Creates a data structure (spatial time series) that appropriately represents the distribution:
 - Space and time become referrers; attributes express containment of objects by spatio-temporal bins
- However, it may conceal important features in the ST distribution
 particularly, spatio-temporal concentrations of events and spatio-temporal coincidences of moves with close origins and/or destinations
- Possible reasons:
- Large bins: features become concealed by averaging
- Small bins: features become dispersed over multiple bins

⇒ST aggregation need to be complemented by ST density-based clustering, which reveals concentrations and coincidences.



Purposes of data transformations

$a\ reminder$

Aggregation

- Supports abstraction, gaining an overall view of characteristics and behaviour
- Reduces large data
- Simplifies complex data
- Extraction of events, etc.
- Selects a portion of data relevant to a task, enables focusing
- Allows dealing with complex data portion-wise
- Integration, disintegration, projection (taking one of possible aspects) $% \left({{{\mathbf{F}}_{i}}_{i}} \right)$
- Adapts data to analysis tasks



a reminder

- **Partition-based clustering**: divide items into groups so that items within a group are similar (close) and items from different groups are less similar (more distant)
 - · Examples: k-means, self-organizing map, hierarchical
- · Property of the result: each item belongs to some group
- **Density-based clustering**: find groups of highly similar (close) items and separate from them items that are less similar (more distant) to others
- Examples: DBScan, OPTICS
- Properties of the results: some items belong to groups, other items remain ungrouped and are treated as "noise"
- DBC and PBC can be combined: first use DBC to find and filter out the noise (i.e., outliers), then apply PBC to the remaining data.
 This may result in cleaner and clearer clusters.

Questions? Analysis of spatial events and OD moves

