# Spatial Data Mining for Customer Segmentation
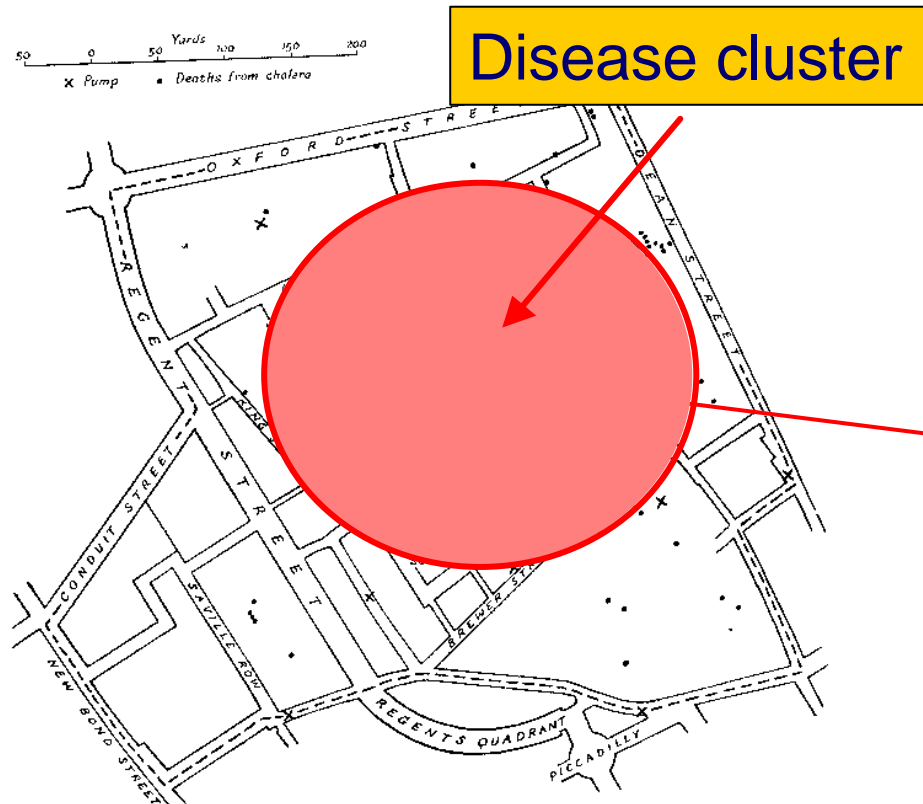
AIS

**Fraunhofer** Institut
Autonome Intelligente
Systeme

## Data Mining in Practice Seminar, Dortmund, 2003

**Dr. Michael May**
**Fraunhofer Institut Autonome Intelligente Systeme**

# Introduction: a classic example for spatial analysis



Disease cluster

Dr. John Snow
Deaths of cholera
epidemia
London, September 1854

Infected water pump?

A good representation is
the key to solving a problem

# Good representation because...

**Represents spatial relation of objects of the same type**

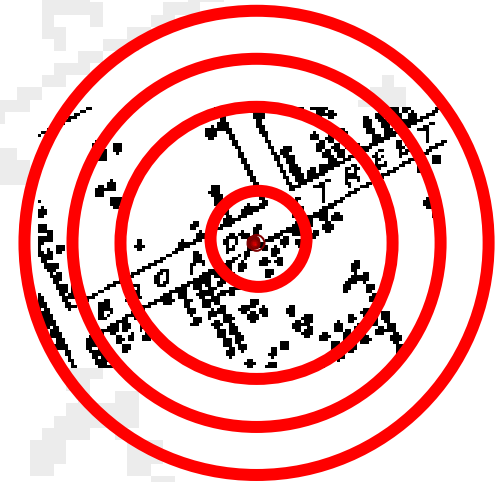**Represents spatial relation of objects to *other* objects**

**Shows only relevant aspects and hides irrelevant**

*It is not only important where a cluster is but also, what else is there (e.g. a water-pump)!*
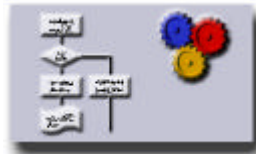
# Goals of Spatial Data Mining

- **Identifying spatial patterns**

- **Identifying spatial objects that are potential generators of patterns**

- **Identifying information relevant for explaining the spatial pattern (and hiding irrelevant information)**

- **Presenting the information in a way that is intuitive and supports further analysis**

# Approach to Spatial Knowledge Discovery

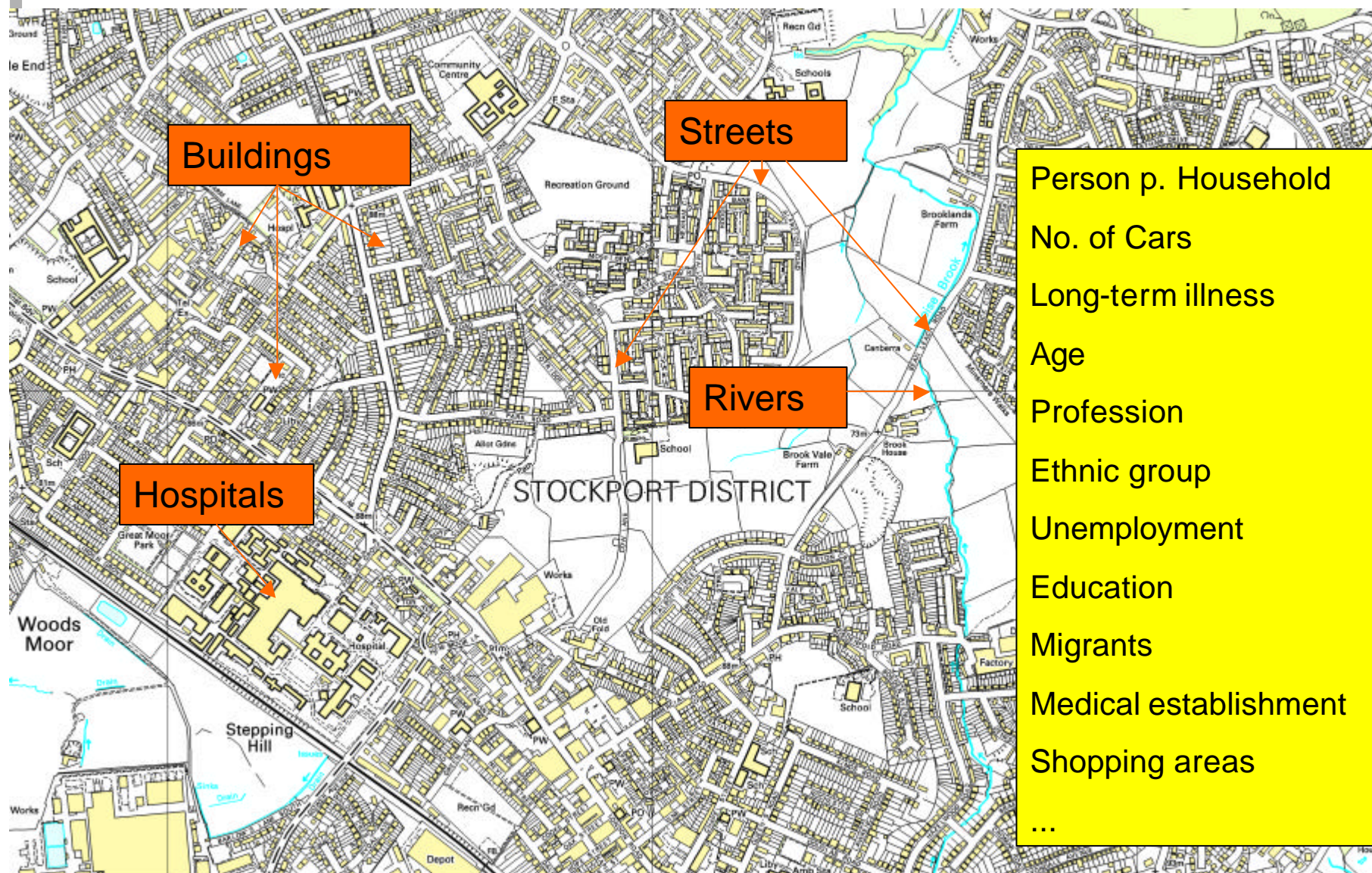Data Mining  $\sqrt{\dfrac{n}{p_0 \cdot (1-p_0)}}(p-p_0)$

**+**

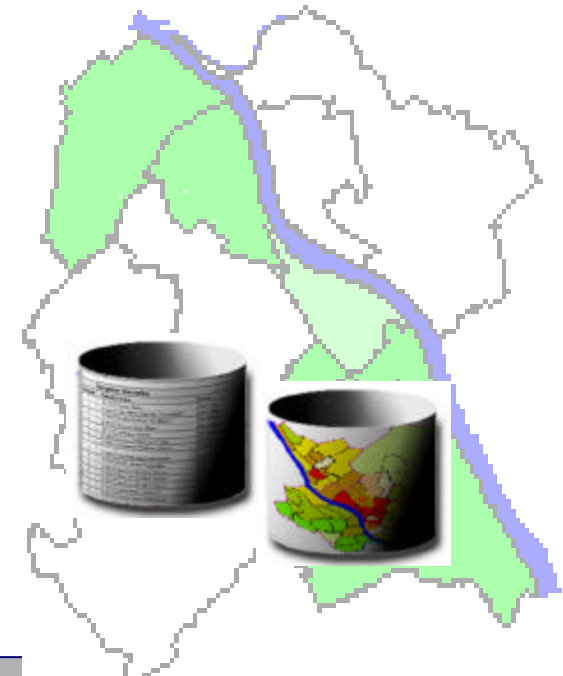Geographic Information Systems 

= SPIN!
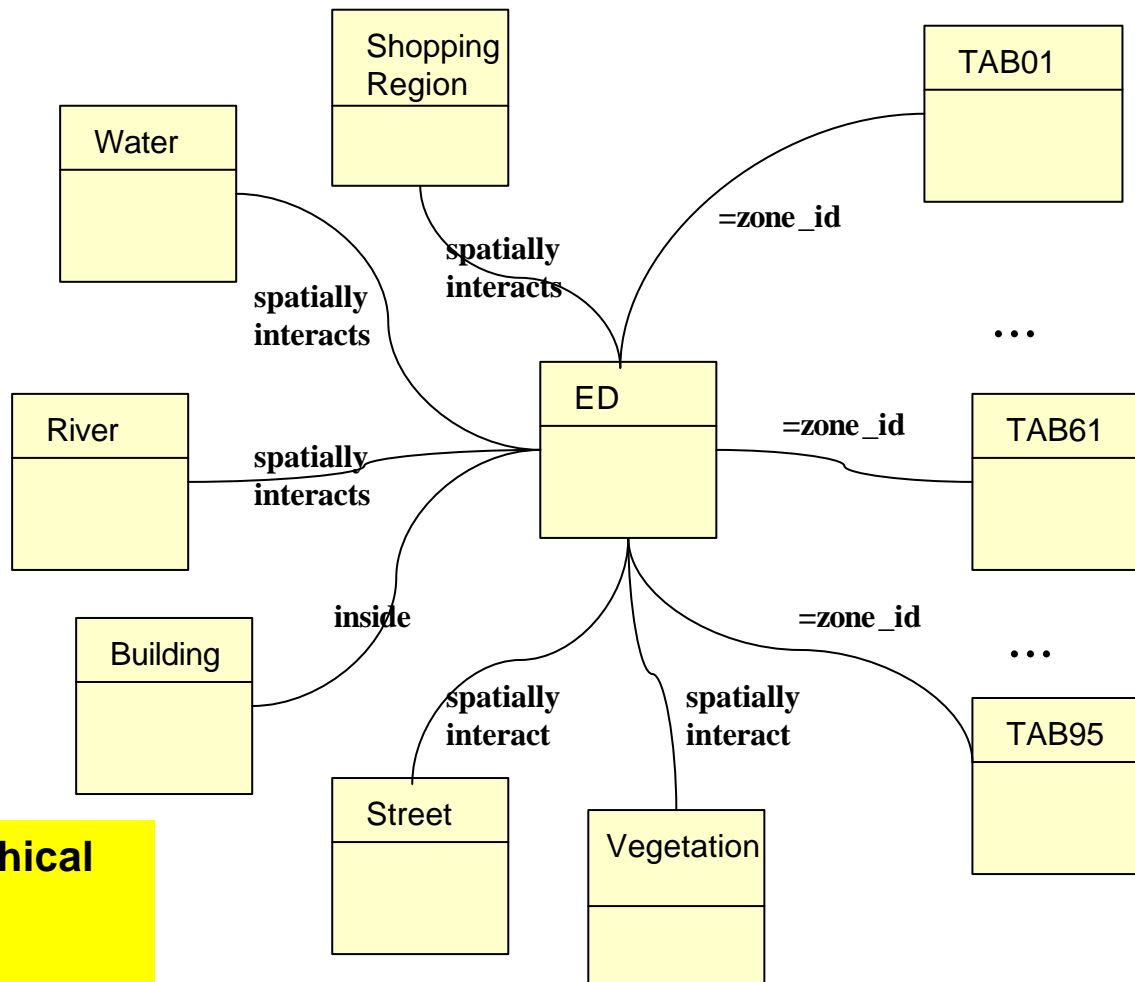
# UK, Greater Manchester, Stockport

# Representation of spatial data
# in Oracle Spatial

**A set of relations $R_1,...,R_n$ such that each relation $R_i$ has a geometry attribute $G_i$ or an identifier $A_i$ such that $R_i$ can be linked (joined) to a relation $R_k$ having a geometry attribute $G_k$**

- Geometry attributes $G_i$ consist of ordered sets of x,y-pairs defining points, lines, or polygons

- Different types of spatial objects are organized in different relations $R_i$ (geographic layers), e.g. streets, rivers, enumeration districts, buidlings, and

- each layer can have its own set of attributes $A_1,..., A_n$ and at most one geometry attribute G

# Stockport Database Schema
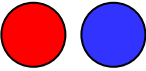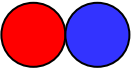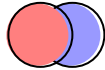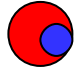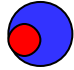


**Attribute data**

95 tables with census data, ~8000 attributes

**Spatial Hierarchy**
- County
- District
- Wards
- Enumeration district

**Geographical Layers**

85 tables

# Spatial Predicates in Oracle Spatial

**Topological relation** (Egenhofer 1991)

| | |
|---|---|
| A disjoint B, B disjoint A | |
| A meets B, B meets A | |
| A overlaps B, B overlaps A | |
| A equals B, B equals A | |
| A covers B, B covered by A | |
| A covered-by B, B covers A | |
| A contains B, B inside A | |
| A inside B, B contains A | |

**Distance relation**: Minimum distance between 2 points

# Typical Data Mining representation

'spreadsheet data'

exactly 1 table

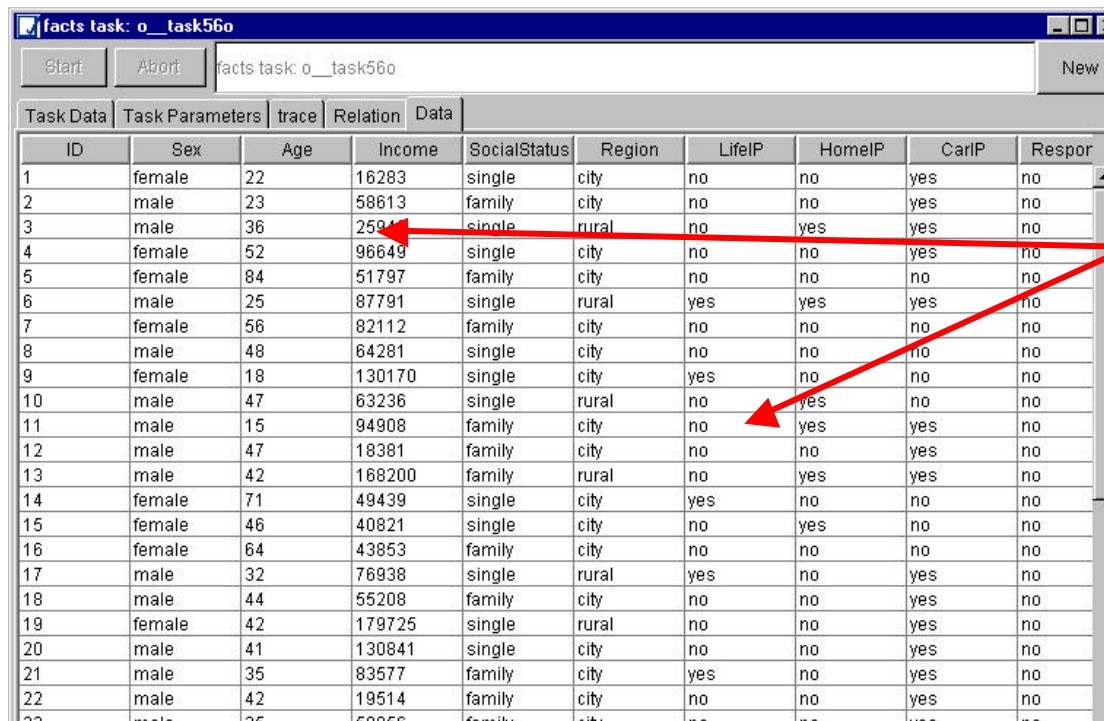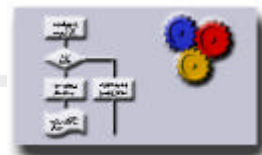| facts task: o__task56o | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Start | Abort | facts task: o__task56o | | | | | | | New |

Task Data | Task Parameters | trace | Relation | **Data**

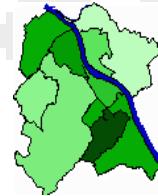| ID | Sex | Age | Income | SocialStatus | Region | LifelP | HomelP | CarlP | Respor |
|---|---|---|---|---|---|---|---|---|---|
| 1 | female | 22 | 16283 | single | city | no | no | yes | no |
| 2 | male | 23 | 58613 | family | city | no | no | yes | no |
| 3 | male | 36 | 259.. | single | rural | no | yes | yes | no |
| 4 | female | 52 | 96649 | single | city | no | no | yes | no |
| 5 | female | 84 | 51797 | family | city | no | no | no | no |
| 6 | male | 25 | 87791 | single | rural | yes | yes | yes | no |
| 7 | female | 56 | 82112 | family | city | no | no | no | no |
| 8 | male | 48 | 64281 | single | city | no | no | no | no |
| 9 | female | 18 | 130170 | single | city | yes | no | no | no |
| 10 | male | 47 | 63236 | single | rural | no | yes | no | no |
| 11 | male | 15 | 94908 | family | city | no | yes | yes | no |
| 12 | male | 47 | 18381 | family | city | no | no | yes | no |
| 13 | male | 42 | 168200 | family | rural | no | yes | yes | no |
| 14 | female | 71 | 49439 | single | city | yes | no | no | no |
| 15 | female | 46 | 40821 | single | city | no | yes | no | no |
| 16 | female | 64 | 43853 | family | city | no | no | no | no |
| 17 | male | 32 | 76938 | single | rural | yes | no | yes | no |
| 18 | male | 44 | 55208 | family | city | no | no | yes | no |
| 19 | female | 42 | 179725 | single | rural | no | no | yes | no |
| 20 | male | 41 | 130841 | single | city | no | no | yes | no |
| 21 | male | 35 | 83577 | family | city | yes | no | yes | no |
| 22 | male | 42 | 19514 | family | city | no | no | yes | no |
| 23 | male | 35 | 58856 | family | city | no | no | yes | no |

atomic values

Data Mining for spatial data: strong discrepancy between usual and adequate problem representation

# SPIN! – The Elements

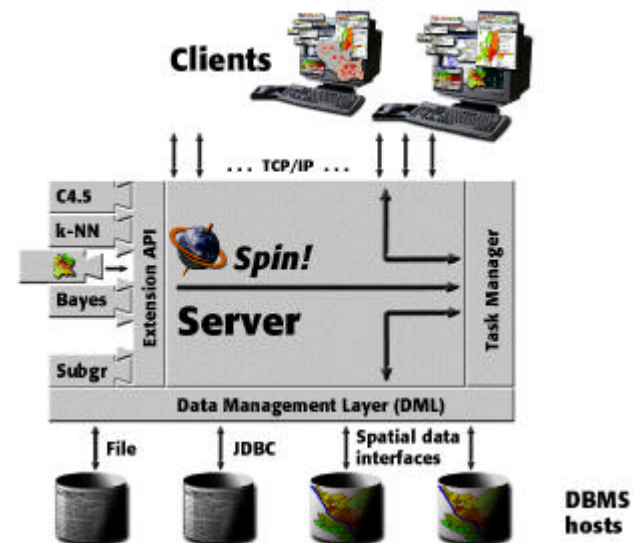$$\sqrt{\frac{n}{p_0 \cdot (1 - p_0)}}(p - p_0)$$
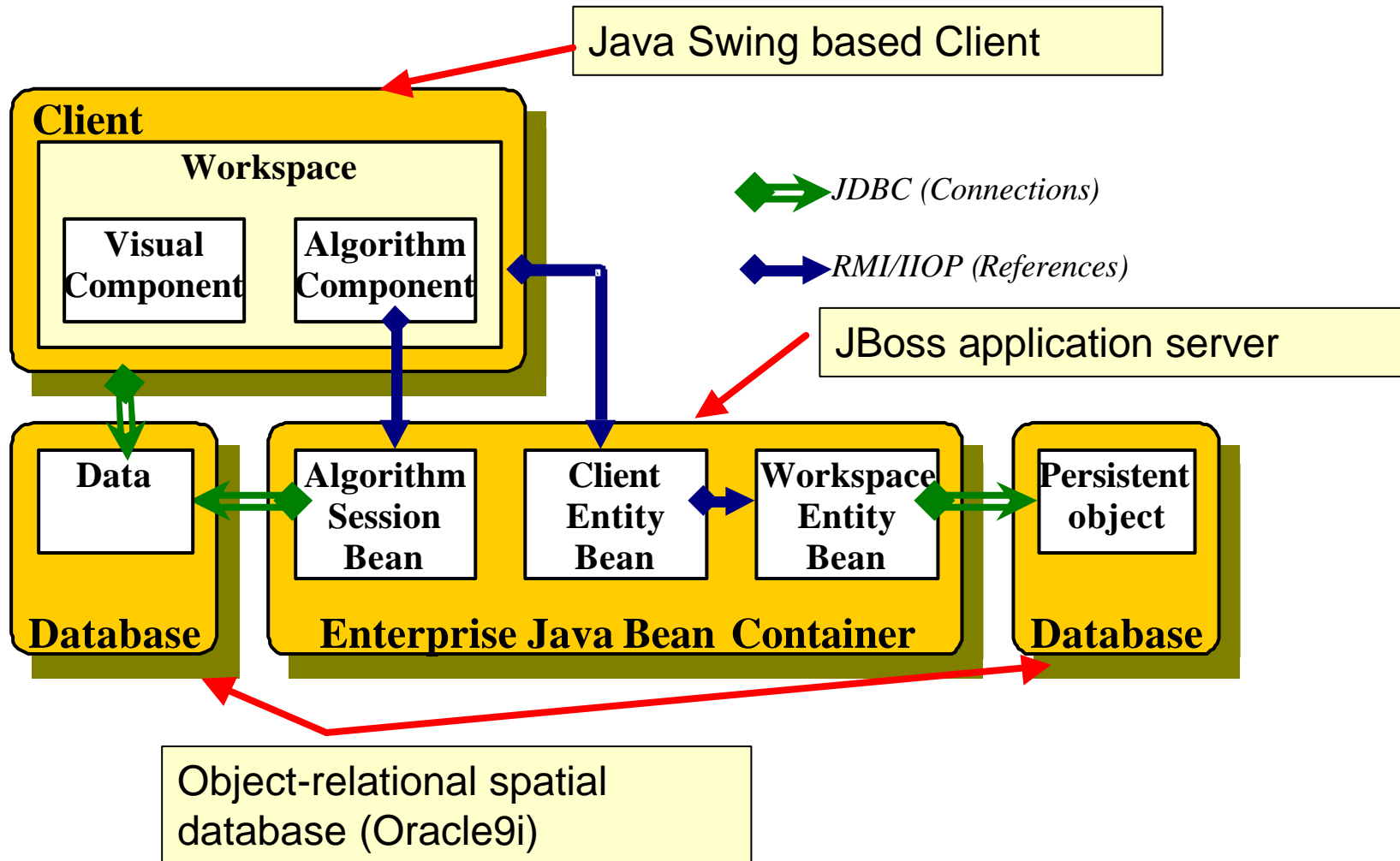
# 1. Spatial Data Mining Platform

# Providing an integrated data mining platform

- Data access to heterogeneous and distributed data sources (Oracle RDBMS, flat file, spatial data)
- Organizing and documenting analysis tasks
- Launching analysis tasks
- Visualizing results

Note: Same software basis as MiningMart!

# SPIN! Architecture: Enterprise Java Bean-based

Java Swing based Client

**Client**

**Workspace**

| Visual Component | Algorithm Component |
|---|---|

JDBC (Connections)

RMI/IIOP (References)

JBoss application server

| Data | Algorithm Session Bean | Client Entity Bean | Workspace Entity Bean | Persistent object |
|---|---|---|---|---|

**Database**

**Enterprise Java Bean Container**

**Database**

Object-relational spatial database (Oracle9i)

# SPIN! User Interface



Point & Click-Tool for defining analysis tasks

Workspace Tree

Property editor

# 2. Visual Exploratory Analyis

# Interactive Exploratory Analysis

**Parallel Coordinate Plot**

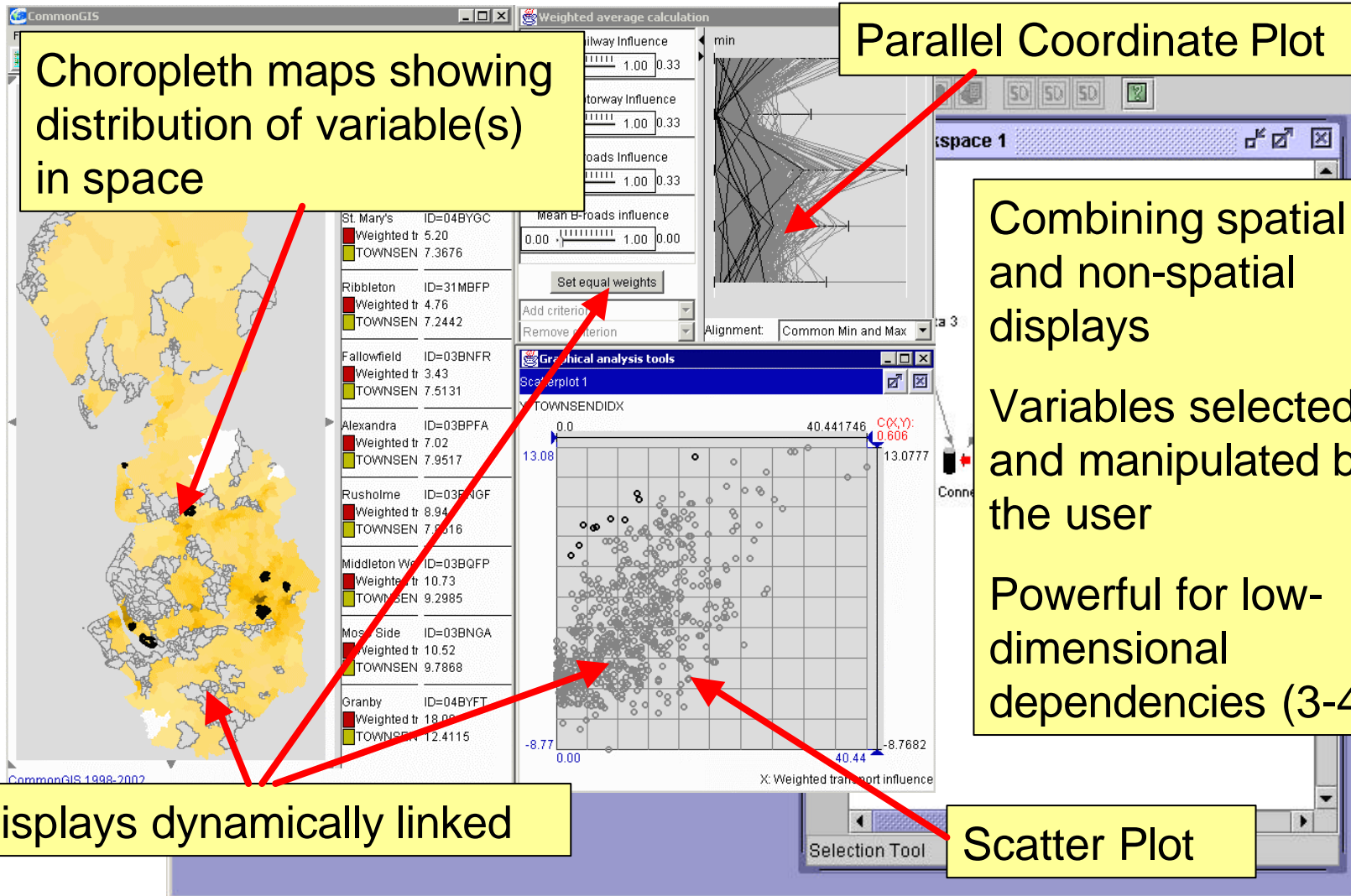**Choropleth maps showing distribution of variable(s) in space**

**Combining spatial and non-spatial displays**

**Variables selected and manipulated by the user**

**Powerful for low-dimensional dependencies (3-4)**

**Displays dynamically linked**

**Scatter Plot**

# 3. Searching for Explanatory

# Data Mining Tasks in SPIN!

- Looking for associations between subsets of spatial and non-spatial attributes

  ⇨ **Spatial Association Rules**

- A phenomenon of interest (e.g. death rate) is given but it is not clear which of a large number of spatial and non-spatial attributes is relevant for explaining it

  ⇨ **Spatial Subgroup Discovery**

- A quantitative variable of interest is given and we ask how much this variable changes when one of the relevant independent variables is changed

  ⇨ **Bayesian Local regression**

# Subgroup Discovery Search

- Subgroup discovery is a multi-relational approach that searches for probabilistically defined deviation patterns (Klösgen 1996, Wrobel 1997)
- Top-down search search from most general to most specific subgroups, exploiting partial ordering of subgroups ($S_1 \geq S_2$ $S_1$ *more general* than $S_{2)}$)
- Beam search expanding only the $n$ best ones at each level of search
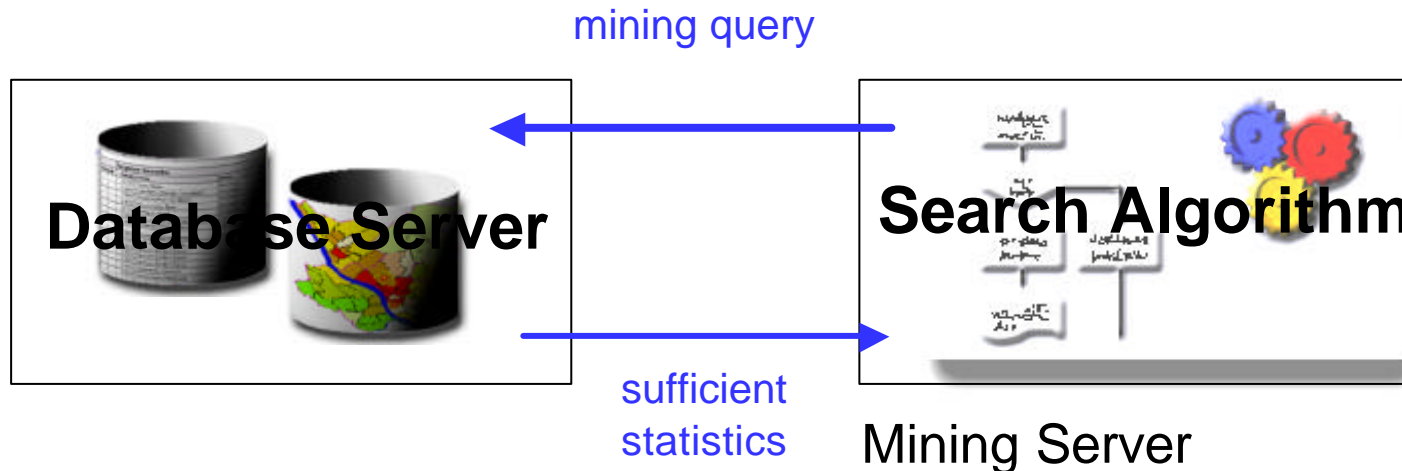- Evaluating hypothesis according to quality function:

T= target group
C= concept

$$\frac{p(T\,|\,C) - p(T)}{\sqrt{p(T)(1-p(T))}} \sqrt{n} \sqrt{\frac{N}{N-n}}$$

```
T = long-term illness=high

C = unemployment=high
```

# Division of labour between Oracle RDBMS and Search Manager

mining query



**Database Server**

**Search Algorithm**

sufficient statistics

Mining Server

- Database integration: efficiently organize mining queries

- Mining query delivers statistics (aggregations) sufficient for evaluating **many** hypotheses

- search in hypothesis space

- generation and evaluation of hypotheses (subgroup patterns)

# Data Mining visualization



High long-term illness in districts crossed by M60

p(T|C) vs. p(C)

Spatial Venn Diagram

Subgroup Overview
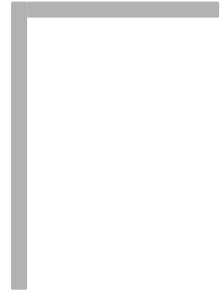
Subgroup

Linked Display

# Customer Analysis Rodgau, Germany

# System Demo:
# Customer Analysis
# using
# MiningMart and SPIN!

# Summary & Outlook

- SPIN! tightly integrates Data Mining analysis and GIS-based visualization

- Main features:
  - A spatial data mining platform
  - New spatial data mining algortihms for subgroup discovery, association rules, Baysian MCMC
  - New visualization methods

- Integration of Spatial Data allows to get results that could not be achieved otherwise

- MiningMart can usefully applied for some pre-processing tasks

- Future tasks: Integrating spatial preprocessing in MiningMart