

AN EFFECTIVE ANALYSIS OF SPATIAL DATA MINING METHODS USING RANGE QUERIES

Gangireddy Ravikumar¹ and Mallireddy Sivareddy²

¹Student, M.Tech (CSE), Green Fields, K.L.University, Andhra Pradesh, India.
¹ravikumar.gv571@gmail.com

²Student, M.Tech (CSE), Green Fields, K.L.University, Andhra Pradesh, India.
²sivadoa47@gmail.com

Abstract: This paper reviews the data mining methods that are combined with Geographic Information Systems (GIS) for carrying out spatial analysis of geographic data. We will first look at data mining functions as applied to such data and then highlight their specificity compared with their application to classical data. We will go on to describe the research that is currently going on in this area, pointing out that there are two approaches: the first comes from learning on spatial databases, while the second is based on spatial statistics. We will conclude by discussing the main differences between these two approaches and the elements they have in common.

Index Terms — Spatial Data Mining, Spatial Databases, Rules Induction, Spatial Statistics, Spatial Neighborhood.

INTRODUCTION

The growing production of maps is generating huge volumes of data that exceed people's capacity to analyze them. It thus seems appropriate to apply knowledge discovery methods like data mining to spatial data. This recent technology is an extension of the data mining applied to alphanumeric data on spatial data. The main difference is that spatial analysis must take into account spatial relations between objects.

The applications covered by spatial data mining are decisional ones, such as geomarketing, environmental studies, risk analysis, and so on. For example, in geomarketing, a store can establish its trade area, i.e. the spatial extent of its customers, and then analyze the profile of those customers on the basis of both their properties and the properties related to the area where they live.

In our Analysis, spatial data mining is applied to traffic risk analysis [2]. The risk estimation is based on the information on the previous injury accidents, combined to thematic data relating to the road network, population, buildings, and so on. The project aims at identifying regions with a high level of risk and analyzing and explaining those risks with respect to the geographic neighborhood. Spatial data mining technology specifically allows for those neighborhood relationships.

Nowadays, data analysis in geography is essentially based on traditional statistics and multidimensional data analysis and does not take account of spatial data [3]. Yet the main specificity of geographic data is that observations located near to one another in space tend to share similar (or correlated) attribute values. This constitutes the fundamental of a distinct scientific area called "spatial statistics" which, unlike traditional statistics, supposes inter-dependence of nearby observations. An abundant bibliography exists in this area, including well-known geostatistics, recent developments in Exploratory Spatial Data Analysis (ESDA) by Anselin and Geographical Analysis Machine (GAM) by

Openshaw. For a summary, refer to Part 1.c of [4]. Multi-dimensional analytical methods have been extended to support contiguity [5, 6]. We maintain that spatial statistics is a part of spatial data mining, since it provides data-driven analyses. Some of those methods are now implemented in operational GIS or analysis tools.

In the field of databases, two main teams have contributed to developing data mining for spatial data analysis. The first one, DB Research Lab (Simon Fraser University, Vancouver), developed GeoMiner [7], which is an extension of DBMiner. The second one (Munich University) devised a structure-of-neighborhood graph [8], on which some algorithms are based. They have also worked on a clustering method based on a hierarchical partitioning (extension of DBSCAN with a R*Tree), classification (extension of ID3 and DBLearn), association rules (based upon an efficient spatial join), characterization and spatial trends. STING (University of California) uses a hierarchical grid to perform optimization on the clustering algorithm [9]. We might also mention work on Data warehouse dedicated to spatial data (University of Laval) [10].

This paper will describe data mining methods for Geographic Information Systems and highlight their value in performing spatial data analysis. It will survey both statistical approaches and those involving inference from databases.

It is structured as follows. In section 2 we define spatial data mining and subdivide it into generic tasks. Then in section 3 we classify spatial data mining methods, whether drawn from the realm of databases, statistics or artificial intelligence, in terms of these different tasks. We go on to compare the statistical analysis approach with the spatial database approach, with the aim of emphasizing their similarities and complementarity. Lastly, we conclude and discuss research issues.

DEFINITION OF SPATIAL DATA MINING

Spatial data mining (SDM) consists of extracting knowledge, spatial relationships and any other properties which are not explicitly stored in the database. SDM is used to find implicit regularities, relations between spatial data and/or non-spatial data. The specificity of SDM lies in its interaction in space. In effect, a geographical database constitutes a spatio-temporal continuum in which properties concerning a particular place are generally linked and explained in terms of the properties of its neighborhood. We can thus see the great importance of spatial relationships in the analysis process. Temporal aspects for spatial data are also a central point but are rarely taken into account.

Data mining methods [11] are not suited to spatial data because they do not support location data nor the implicit relationships between objects. Hence, it is necessary to develop new methods including spatial relationships and spatial data handling. Calculating these spatial relationships is time consuming, and a huge volume of data is generated by encoding geometric location. Global performances will suffer from this complexity.

Using GIS, the user can query spatial data and perform simple analytical tasks using programs or queries. However, GIS are not designed to perform complex data analysis or knowledge discovery. They do not provide generic methods for carrying out analysis and inferring rules. Nevertheless, it seems necessary to integrate these existing methods and to extend them by incorporating spatial data mining methods. GIS methods are crucial for data access, spatial joins and graphical map display. Conventional data mining can only generate knowledge about alphanumeric properties.

Spatial Data Mining:

Spatial data mining is the application of data mining techniques to spatial data. Data mining in general is the search for hidden patterns that may exist in large databases. Spatial data mining is the discovery of interesting relationship and characteristics that may exist implicitly in spatial databases. Because of the huge amounts (usually, terabytes) of spatial data that may be obtained from satellite images, medical equipments, video cameras, etc. It is costly and often unrealistic for users to examine spatial data in detail. Spatial data mining aims to automate such a knowledge discovery process. Thus it plays an important role in

- Extracting interesting spatial patterns and features.
- Capturing intrinsic relationships between spatial and non spatial data.
- Presenting data regularity concisely and at higher conceptual levels and
- Helping to reorganize spatial databases to accommodate data semantics, as well as to achieve better performance.

Spatial database stores a large amount of space related data, such as maps, preprocessed remote sensing or medical imaging data and VLSI chip layout data. Spatial databases have many features distinguishing them from relational databases. They carry topological and/or distance information, usually organized by sophisticated, multi dimensional spatial indexing structures that are accessed by

spatial data access methods and often require spatial reasoning, geometric computation, and spatial knowledge representation techniques.

Spatial Data Mining Structure:

The spatial data mining can be used to understand spatial data, discover the relation between space and the non space data, set up the spatial knowledge base, excel the query, reorganize spatial database and obtain concise total characteristic etc.. The system structure of the spatial data mining can be divided into three layer structures mostly, such as the Figure 1 show [1]. The customer interface layer is mainly used for input and output, the miner layer is mainly used to manage data, select algorithm and storage the mined knowledge, the data source layer, which mainly includes the spatial database (camalig) and other related data and knowledge bases, is original data of the spatial data mining.

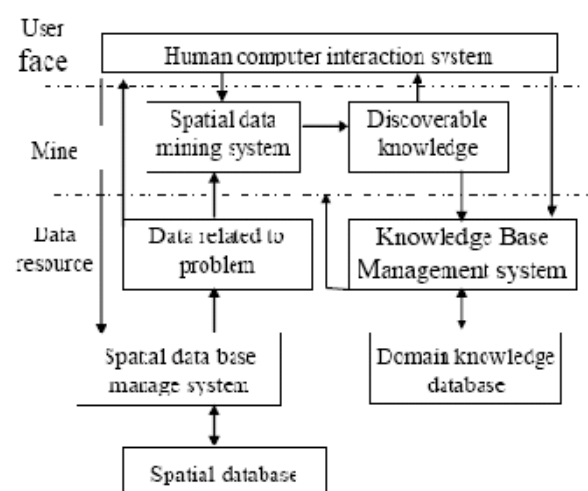


Figure.1 The systematic structure of spatial data mining

Primitives of Spatial Data Mining:

- Rules:** There are several kinds of rules can be discovered from databases in general. For example characteristic rules, discriminate rules, association rules, or deviation and evaluation rules can be mined [1]. A Spatial characteristic rule is a general description of the spatial data.

For example, a rule describing the general price range of houses in various geographic regions in a city is a spatial characteristic rule. A discriminate rule is general description of the features discriminating or contrasting a class of spatial data from other class(es) like the comparison of price ranges of houses in different geographical regions. A spatial association rule is a rule which describes the implication of one a set of features by another set of features in spatial databases. For example, a rule associating the price range of the houses with nearby spatial features, like beaches, is a spatial association rule.

- Thematic Maps:** Thematic map is map primarily design to show a theme, a single spatial distribution or a pattern, using a specific map type. These maps show the distribution of features over limited geography areas [1]. Each map defines a partitioning of the area into a set of closed and disjoint regions; each includes all the points with the same feature value. Thematic maps present the spatial distribution of a single or a

few attributes. This differs from general or reference maps where the main objective is to present the position of the object in relation to other spatial objects. Thematic maps may be used for discovering different rules. For example, we may want to look at temperature thematic map while analyzing the general weather pattern of a geographic region. There are two ways to represent thematic maps: Raster, and Vector.

In the raster image form thematic maps have pixels associated with the attribute values. For example, a map may have the altitude of the spatial objects coded as the intensity of the pixel (or the color). In the vector representation, a spatial object is represented by its geometry, most commonly being the boundary representation along with the thematic attributes. For example, a park may be represented by the boundary points and corresponding elevation values.

SPATIAL DATA MINING TASKS

As shown in the table below, spatial data mining tasks are generally an extension of data mining tasks in which spatial data and criteria are combined. These tasks aim to: (i) summarize data, (ii) find classification rules, (iii) make clusters of similar objects, (iv) find associations and dependencies to characterize data, and (v) detect deviations after looking for general trends. They are carried out using different methods, some of which are derived from statistics and others from the field of machine learning.

SDM Tasks	Statistics	Machine Learning
Summarization	Global autocorrelation Density analysis Smooth and contrast analysis Factorial analysis	Generalization Characteristic rules
Class identification	Spatial classification	Decision trees
Clustering	Point pattern analysis	Geometric clustering
Dependencies	Local autocorrelation Correspondence analysis	Association rules
Trends and deviations	Kriging	Trend rules

Table 1: Comparison between statistical and machine learning approaches to SDM

The rest of this section is devoted to describing data mining tasks that are dedicated to GIS.

Spatial data summarization:

The main goal is to describe data in a global way, which can be done in several ways. One involves extending statistical methods such as variance or factorial analysis to spatial structures. Another entails applying the generalization method to spatial data.

Statistical analysis of contiguous objects:

Global autocorrelation: The most common way of summarizing a dataset is to apply elementary statistics, such as the calculation of average, variance, etc., and graphic tools like histograms and pie charts. New methods have been developed for measuring neighborhood dependency at a global level, such as local variance and local covariance, spatial auto-correlation by Geary, and Moran indices [12, 13]. These methods are based on the notion of a contiguity matrix that represents the spatial relationships between objects. It should be noted that this contiguity can

correspond to different spatial relationships, such as adjacency, a distance gap, and so on.

Density analysis: This method forms part of Exploratory Spatial Data Analysis (ESDA) which, contrary to the autocorrelation measure, does not require any knowledge about data. The idea is to estimate the density by computing the intensity of each small circle window on the space and then to visualize the point pattern. It could be described as a graphical method.

Smooth, contrast and factorial analysis: In density analysis, non-spatial properties are ignored. Geographic data analysis is usually concerned with both alphanumeric properties (called attributes) and spatial data. This requires two things: integrating spatial data with attributes in the analysis process, and using multidimensional data to analyze multiple attributes. To integrate the spatial neighborhood into attributes, two techniques exist that modify attribute values using the contiguity matrix. The first technique performs a smoothing by replacing each attribute value by the average value of its neighbors. This highlights the general characteristics of the data. The other contrasts data by subtracting this average from each value. Each attribute (called variable) in statistics can then be analyzed using conventional methods. However, when multiple attributes (above tree) have to be analyzed together, multidimensional data analysis methods (i.e. factorial analysis) become necessary [6]. Their principle is to reduce the number of variables by looking for the factorial axes where there is maximum spreading of data values. By projecting and visualizing the initial dataset on those axes, the correlation or dependencies between properties can be deduced. In statistics and especially in the above methods, the analyzed objects were originally considered to be independent. The need to look at spatial organization spawned several research studies [6, 14]. The extension of factorial analysis methods to contiguous objects entails applying common Principal Component Analysis or Correspondence Analysis methods once the original table is transformed using smoothing or contrasting techniques.

Generalization:

This method consists of raising the abstract level of non-spatial attributes and reducing the detail of geometric description by merging adjacent objects. It is derived from the concept of attribute-oriented induction as described in [7]. Here, a concept hierarchy can be spatial (like the hierarchy of administrative boundaries) or non-spatial (thematic) [15]. An example of thematic hierarchy in agriculture can be represented as follows: "cultivation type (food (cereals (maize, wheat, rice), vegetable, fruit, other)". That kind of hierarchy can be directly introduced by an expert in the field or generated by an inference process related to the attribute. A spatial hierarchy may preexist, like the administrative boundaries one, or it may be based on an artificial geometric splitting like a quad-tree [16], or it may result from a spatial clustering (see below). There are two kinds of generalization: non-spatial dominant generalization, where we first use a thematic hierarchy and then merge adjacent objects; and spatial dominant generalization, which is based on a spatial hierarchy to begin with, followed by the aggregation or generalization of non-spatial values for each generalized spatial value. The complexity of the

corresponding algorithms is $O(N \log N)$, where N is the number of actual objects. This approach could be treated as a first step towards a method of inferring rules, such as association rules or comparison rules.

Characteristic rules:

The characterization of a selected part of the database has been defined in [17] as the description of properties that are typical for the part in question but not for the whole database. In the case of a spatial database, it takes account not only of the properties of objects, but also of the properties of their neighborhood up to a given level.

Consider a subset S of objects to analyze. This method uses the following parameters: 1) significance (relative frequency to the database in S); 2) confidence (ratio of objects in S which satisfy the significance threshold in the neighborhood); and 3) the maximum extension \max -neighbors to the neighbors. This method throws up the properties $p_i = (\text{attribute}, \text{value})$, the relative frequency factors $\text{freq-fac } i$ (higher than the significance parameter) and the number n_i of neighbors on which the frequency of the property is extended. The characterization can be expressed by the following rule:

$$S \Rightarrow p_1 (n_1, \text{freq-fac } 1) \wedge \dots \wedge p_k (n_k, \text{freq-fac } k).$$

Class identification:

This task, also called supervised classification, provides a logical description that yields the best partitioning of the database. Classification rules constitute a decision tree where each node contains a criterion on an attribute. The difference in spatial databases is that this criterion could be a spatial predicate and, because spatial objects are dependent on neighborhood, a rule involving the non-spatial properties of an object should be extended to neighborhood properties. In spatial statistics, classification has essentially served to analyze remotely-sensed data, and aims to identify each pixel with a particular category. Homogeneous pixels are then aggregated in order to form a geographic entity [4]. In the spatial database approach [18], classification is seen as an arrangement of objects using both their properties (non-spatial values) and their neighbors' properties, not only for direct neighbors but also for the neighbors of neighbors and so on, up to degree N . Let us take as an example the classification of areas by their economic power. Classification rules are described as follows:

High population \wedge neighbor = road \wedge neighbor of neighbor = airport \Rightarrow high economic power (95%).

In GeoMiner, a classification criterion can also be related to a spatial attribute, in which case it reflects its inclusion in a wider zone. These zones could be determined by the algorithm, whether by clustering or by merging adjacent objects, or it could arise from a predefined spatial hierarchy. A new algorithm [19] extends this classification method in GeoMiner to spatial predicates. For example, to determine high level wholesale profits, a decision factor can be the proximity to densely populated districts.

Clustering:

This task is an automatic or unsupervised classification that yields a partition of a given dataset depending on a similarity function.

Database approach:

Paradoxically, clustering methods for spatial databases do not appear to be very revolutionary compared with those applied to relational databases (automatic classification). The clustering is performed using a similarity function which was already classed as a semantic distance. Hence, in spatial databases it appears natural to use the Euclidean distance in order to group neighboring objects. Research studies have focused on the optimization of algorithms. Geometric clustering generates new classes, such as the location of houses in terms of residential areas. This stage is often performed before other data mining tasks, such as association detection between groups or other geographic entities, or characterization of a group.

GeoMiner combines geometric clustering applied to a point set distribution with generalization based on non-spatial attributes. For example, we may want to characterize groups of major cities in the United States and see how they are grouped. Cluster results will be represented by new areas, which correspond to the convex hull of a group of towns. A few points could stay outside clusters and represent noise. A description of each group may be generated for each attribute specified.

Many algorithms have been proposed for performing clustering, such as CLARANS [20], DBSCAN [8] or STING [9]. They usually focus on cost optimization. Recently, a method that is more specifically applicable to spatial data, GDBSCAN, was outlined in [21]. It applies to any spatial shape, not only to points data, and incorporates attributes data.

Statistic approach:

Clustering arises from point pattern analysis [22, 23] and was mainly applied to epidemiological research. This is implemented in Openshaw's well-known Geographical Analysis Machine (GAM) and could be tested by using the K-function [24]. The clusters could also be detected by the ratio of two density estimates: one of the studied subset and the other of the whole reference dataset.

Trend and Deviation Analysis:

In relational databases, this analysis is applied to temporal sequences. In spatial databases, we want to find and characterize spatial trends.

Database approach:

Using the process described in [18], which is based on the central places theory, the analysis is performed in four stages. The first one involves discovering centers by computing local maxima of particular attributes; in the second, the theoretical trend of these attributes is determined by moving away from the centers; the third stage determines the deviations in relation to these trends; and finally, we explain these trends by analyzing the properties of these zones. One example is the trend analysis of the unemployment rate in comparison with the distance to a metropolis like Munich. Another example is the trend analysis of the development of house construction.

Geostatistical approach:

Geostatistics is a tool used for spatial analysis and for the prediction of spatio-temporal phenomena. It was first used

for geological applications (the geo prefix comes from geology). Nowadays, geostatistics encompasses a class of techniques used to analyze and predict the unknown values of variables distributed in space and/or time. These values are supposed to be connected to the environment. The study of such a correlation is called structural analysis. The prediction of location values outside the sample is then performed by the "kriging" technique [25]. It is important to remember that geostatics is limited to point set analysis or polygonal subdivisions and deals with a unique variable or attributes. Under those conditions, it constitutes a good tool for spatial and spatio-temporal trend analysis.

CONCLUSION

Different methods of data mining in spatial databases have been outlined in this paper, which has shown that these methods have been developed by two very separate research communities: the Statistics community and the Database community. We have summarized and classified this research and compared the two approaches, emphasizing the particular utility of each method and the possible advantages of combining them. This work constitutes a first step towards a methodology incorporating the whole process of knowledge discovery in spatial databases and allowing the combination of the above data mining techniques.

Among the other issues in the area of spatial data mining, one approach is to consider the temporality of spatial data, while another is to see how linear or network shape (like roads) can have a particular influence on graphical methods. In any event, it remains essential to continue enhancing the performance of these techniques. One reason is the enormous volumes of data involved, another is the intensive use of spatial proximity relationships. In the case of graphical methods, these relationships could be optimized using spatial indexes. As regards the other methods that use neighborhood structures, instantiation of the structure is costly and should be pre-computed as far as possible.

ACKNOWLEDGEMENTS

We are greatly delighted to place my most profound appreciation to Dr. K. Satyanarayana Chancellor of K.L.University, Dr. K. Raja Sekhara Rao Principal and Dr. K. Subramanyam coordinator for M.Tech under their guidance and encouragement and kindness in giving us the opportunity to carry out the paper. Their pleasure nature, directions, concerns towards us and their readiness to share ideas enthused us and rejuvenated our efforts towards our goal. We also thank the anonymous references of this paper for their valuable comments.

REFERENCES

- [1]. M.Hemalatha.M; Naga Saranya.N. A Recent Survey on Knowledge Discovery in Spatial Data Mining, IJCI International Journal of Computer Science, Vol 8, Issue 3, No.2, may,2011.
- [2]. Zeitouni, K.: Etude de l'application du data mining à l'analyse spatiale du risque d'accidents routiers par l'exploration des bases de données en accidentologie, Final report of the contract PRISM -INRETS, December 1998, 33 p.
- [3]. Sanders, L.: L'analyse statistique des données en géographie, GIP Reclus, 1989
- [4]. Longley P. A., Goodchild M. F., Maguire D. J., Rhind D. W., Geographical Information Systems - Principles and Technical Issues, John Wiley & Sons, Inc., Second Edition, 1999.
- [5]. Lebart L. et al., "Statistique exploratoire multidimensionnelle" , Editions Dunod, Paris, 439 p., 1997.
- [6]. Lebart, L. (1984) Correspondence analysis of graph structure. Bulletin technique du CESIA, Paris:2, 1-2, pp 5-19.
- [7]. Lu, W., Han, J. and Ooi, B.: Discovery of General Knowledge in Large Spatial Databases, in Proc. of 1993 Far East Workshop on Geographic Information Systems (FEGIS'93), Singapore, June 1993, pp. 275-289
- [8]. Ester, M., Kriegel ,H.-P., Sander, J., Xu, X.: Density-Connected Sets and their Application for Trend Detection in Spatial Databases, Proc. 3rd Int. Conf. on Knowledge Discovery and Data Mining, Newport Beach, CA, 1997, pp.10-15
- [9]. Wang, W., Yang, J., and Muntz, R.: STING : A Statistical Information Grid Approach to Spatial Data Mining, Technical Report CSD-97006, Computer Science Department, University of California, Los Angeles, February 1997
- [10]. Bédard, Y., Lam, S., Proulx, M.J., Caron, P.Y., Létourneau, F.: Data Warehousing for Spatial Data: Research Issues, Proceedings of the International Symposium Geomatics in the Era of Radarsat (GER'97), Ottawa, May 1997, pp. 25-30
- [11]. Fayyad et al., "Advances in Knowledge Discovery and Data Mining", AAAI Press / MIT Press, 1996
- [12]. Geary R.C.: The contiguity ratio and statistical mapping, The incorporated Statistician, 5 (3), pp 115-145.
- [13]. Moran P.A.P., The interpretation of statistical maps, Journal of the Royal Statistical Society, B: 10, pp 234-251.,1948.
- [14]. Benali, H., Escofier, B.: Analyse factorielle lissée et analyse factorielle des différences locales, Revue Statistique Appliquée, 1990, XXXVIII (2), pp 55-76
- [15]. Han J., Cai Y. & Cerone N., "Knowledge Discovery in Databases; An Attribute-Oriented Approach." Proceedings of the 18th VLDB Conference. Vancouver, B.C., August 1992. pp. 547-559
- [16]. Samet H., "Design and Analysis of Spatial Data Structures: Hierarchical (quadtree and octree) data structures " , Addison-Wesley Edition, 1990
- [17]. Ester, M., Frommelt, A., Kriegel, H.-P., Sander J.: Algorithms for Characterization and Trend Detection in Spatial Databases, Proc. 4th Int. Conf. on Knowledge Discovery and Data Mining, New York, NY, 1998
- [18]. Ester, M., Kriegel, H.-P., Sander, J.: Spatial Data Mining: A Database Approach, Proc. 5th Symp. on Spatial Databases, Berlin, Germany, 1997

- [19]. Koperski, K., Han, J., and Stefanovi,c N.: An Efficient Two-Step Method for Classification of Spatial Data, In Proc. International Symposium on Spatial Data Handling (SDH'98) , pp. 45-54, Vancouver, Canada, July 1998
- [20]. Ng, R. and Han, J.: Efficient and Effective Clustering Method for Spatial Data Mining, in Proc. of 1994 Int'l Conf. on Very Large Data Bases (VLDB'94), Santiago, Chile, September 1994, pp. 144-155
- [21]. Knorr E. M., and Ng R. T.: Finding Aggregate Proximity Relationships and Commonalities in Spatial Data Mining, IEEE Transactions in Knowledge and Data Engineering, Vol 8(6), December 1996.
- [22]. Openshaw S., Charlton M., Wymer C., Craft A., 1987 : "A mark 1 geographical analysis machine for the automated analysis of point data sets", International Journal of Geographical Information Systems, Vol. 1, n° 4, pp. 335-358
- [23]. Fotheringham S., Zhan B., 1996 : "A comparison of three exploratory methods for cluster detection in spatial point patterns", Geographical Analysis, Vol. 28, n° 3, pp. 200-218
- [24]. Diggle P.J., 1993, Point process modeling in environmental epidemiology. In Barnett V., Turkman K. (eds) Statistics for the environment, Chichester, John Wiley & Sons, pp 89-110.
- [25]. Isobel C., "Practical geostatistics", Applied Science Publisher, Reprinted 1987. Also at URL: <<http://curie.ej.jrc.it/faq/introduction.html>>