

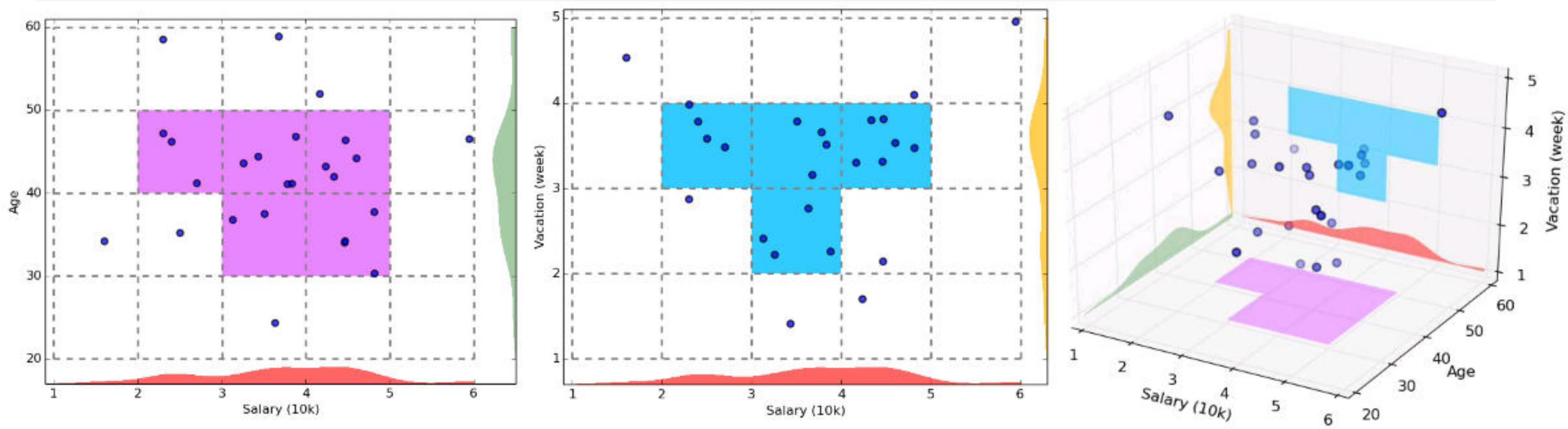


CLIQUE: Grid-Based Subspace Clustering

CLIQUE: Grid-Based Subspace Clustering

- CLIQUE (Clustering In QUES) (Agrawal, Gehrke, Gunopulos, Raghavan: SIGMOD'98)
- CLIQUE is a **density-based** and **grid-based** **subspace clustering** algorithm
 - **Grid-based**: It discretizes the data space through a grid and estimates the density by counting the number of points in a grid cell
 - **Density-based**: A cluster is a maximal set of connected dense units in a subspace
 - A unit is dense if the fraction of total data points contained in the unit exceeds the input model parameter
 - **Subspace clustering**: A subspace cluster is a set of neighboring dense cells in an arbitrary subspace. It also discovers some minimal descriptions of the clusters
- It automatically identifies subspaces of a high dimensional data space that allow better clustering than original space using the Apriori principle

CLIQUE: SubSpace Clustering with Aprori Pruning



- Start at 1-D space and discretize numerical intervals in each axis into grid
- Find dense regions (clusters) in each subspace and generate their minimal descriptions
- Use the dense regions to find promising candidates in 2-D space based on the Apriori principle
两个颜色的色块.
- Repeat the above in level-wise manner in higher dimensional subspaces

Major Steps of the CLIQUE Algorithm

- Identify subspaces that contain clusters
 - Partition the data space and find the number of points that lie inside each cell of the partition
 - Identify the subspaces that contain clusters using the Apriori principle
- Identify clusters
 - Determine dense units in all subspaces of interests
 - Determine connected dense units in all subspaces of interests
- Generate minimal descriptions for the clusters
 - Determine maximal regions that cover a cluster of connected dense units for each cluster
 - Determine minimal cover for each cluster

Additional Comments on *CLIQUE*

□ Strengths

- *Automatically* finds subspaces of the highest dimensionality as long as high density clusters exist in those subspaces
- *Insensitive* to the order of records in input and does not presume some canonical data distribution
- Scales *linearly* with the size of input and has good scalability as the number of dimensions in the data increases

□ Weaknesses

- As in all grid-based clustering approaches, the quality of the results crucially depends on the appropriate choice of the number and width of the partitions and grid cells

Recommended Readings

- ❑ M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases. KDD'96
- ❑ W. Wang, J. Yang, R. Muntz, STING: A Statistical Information Grid Approach to Spatial Data Mining, VLDB'97
- ❑ R. Agrawal, J. Gehrke, D. Gunopulos, and P. Raghavan. Automatic Subspace Clustering of High Dimensional Data for Data Mining Applications. SIGMOD'98
- ❑ A. Hinneburg and D. A. Keim. An Efficient Approach to Clustering in Large Multimedia Databases with Noise. KDD'98
- ❑ M. Ankerst, M. M. Breunig, H.-P. Kriegel, and J. Sander. Optics: Ordering Points to Identify the Clustering Structure. SIGMOD'99
- ❑ M. Ester. Density-Based Clustering. In (Chapter 5) Aggarwal and Reddy (eds.), Data Clustering: Algorithms and Applications . CRC Press. 2014
- ❑ W. Cheng, W. Wang, and S. Batista. Grid-based Clustering. In (Chapter 6) Aggarwal and Reddy (eds.), Data Clustering: Algorithms and Applications. CRC Press. 2014