Data Mining Anomaly Detection			
L	ecture Notes for Chap	ter 10	
Introduction to Data Mining			
Tan, Steinbach, Kumar			
© Tan,Steinbach, Kumar	Introduction to Data Mining	4/18/2004	1

Anomaly/Outlier Detection • What are anomalies/outliers? The set of data points that are considerably different than the remainder of the data • Variants of Anomaly/Outlier Detection Problems Given a database D, find all the data points $\boldsymbol{x} \in \mathsf{D}$ with anomaly scores greater than some threshold t Given a database D, find all the data points $\boldsymbol{x} \in \mathsf{D}$ having the topn largest anomaly scores f(x) Given a database D, containing mostly normal (but unlabeled) data points, and a test point **x**, compute the anomaly score of **x** with respect to D Applications: - Credit card fraud detection, telecommunication fraud detection, network intrusion detection, fault detection C Tan, Steinbach, Kumar Introduction to Data Mining 4/18/2004 (#)

















Outliers in Lower Dimensional Projection

- In high-dimensional space, data is sparse and notion of proximity becomes meaningless
 - Every point is an almost equally good outlier from the perspective of proximity-based definitions
- Lower-dimensional projection methods

Introduction to Data Mining

 A point is an outlier if in some lower dimensional projection, it is present in a local region of abnormally low density

© Tan,Steinbach, Kumar

4/18/2004

(#)

Density-based: LOF approach

- For each point, compute the density of its local neighborhood
- Compute local outlier factor (LOF) of a sample *p* as the average of the ratios of the density of sample *p* and the density of its nearest neighbors
- Outliers are points with largest LOF value



In the NN approach, p_2 is not considered as outlier, while LOF approach finds both p_1 and p_2 as outliers

4/18/2004

(#)

