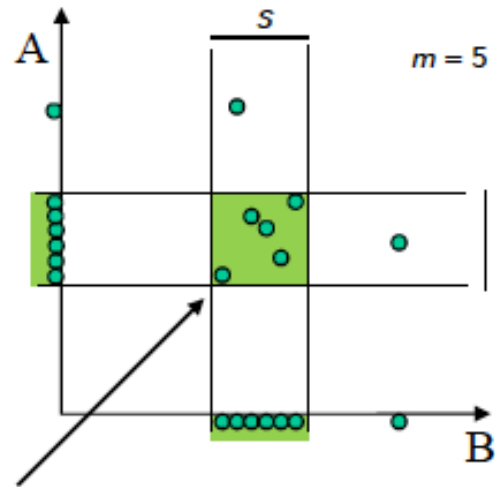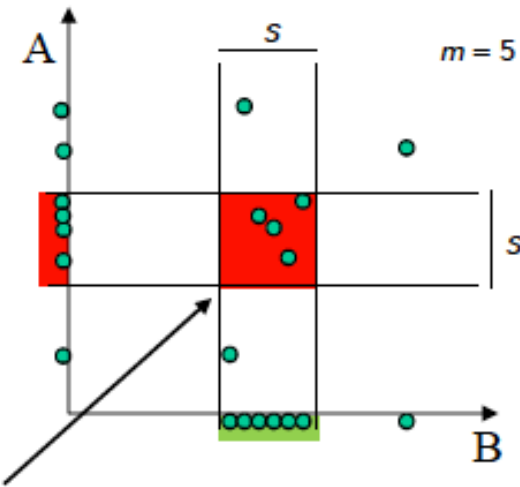# KDD II – Exercise 4

12.06.2017

## Downward-closure property: example

- Simple cluster criterion (density of grid cells):
  - If a cell $C$ of side length $s$ contains more than $m$ points, it represents a cluster

- Monotonicity:
  - if $C$ contains more than $m$ points in subspace $S$ then $C$ also contains more than $m$ points in any subspace $T \subset S$
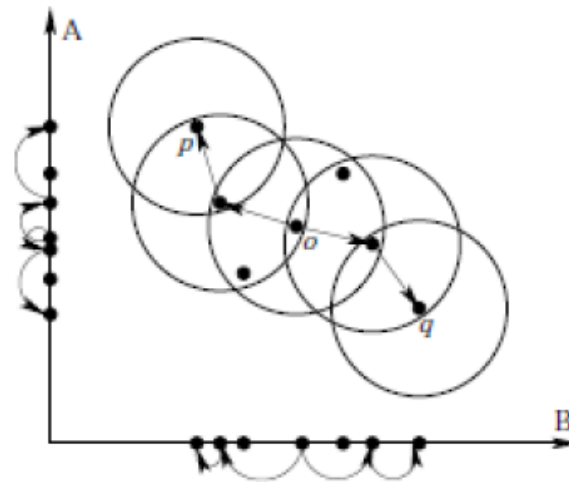  - Example: monotonicity (left) and reverse implication (right)



Cell $C$ contains more than $m=5$ points in subspace „AB"
=> Also in subspaces „A" $\subset$ „AB" and „B" $\subset$ „AB"

Cell $C$ contains less than $m=5$ points in subspace „A"
=> Also in subspace „AB"

If C is a density connected set in subspace S then C is a density connected set in any subspace $T \subset S$.
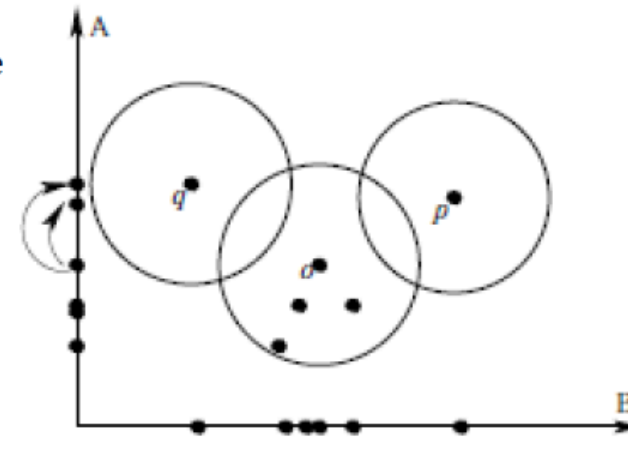
- But, if C is a cluster in S, it need not to be a cluster in $T \subset S$ – maximality might be violated

- All clusters in a higher-dimensional subspace will be subsets of the clusters detected in this first clustering.

ε: circles indicate

$MinPts = 4$

(a) p and q are density-connected via o

(b) p and q are not density-connected

p and q density connected in {A,B}.
Thus, they are also density connected in {A} and {B}

p and q not density connected in {B}.
Thus, they are not density connected in{A,B}, although they are density connected in {A}.

- Algorithm

  - All subspaces that contain any density-connected set are computed using the bottom-up approach (similar to CLIQUE/APRIORI)

  - Density-connected clusters are computed using a specialized DBSCAN run in the resulting subspace to generate the subspace clusters

- Discussion

  - Input: $\varepsilon$ and *MinPts* specifying the density threshold

  - Output: all clusters in all subspaces, clusters may overlap

  - Uses a fixed density threshold for all subspaces

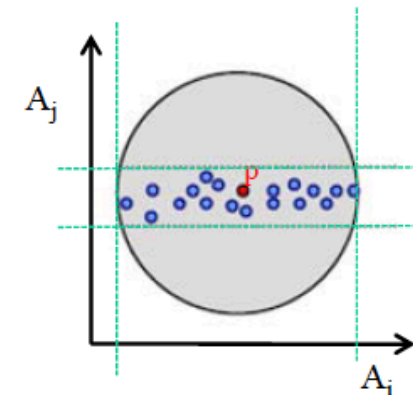  - Advanced but costly cluster model

- Instance-based top-down approach: we learn the subspace for each instance
- Extends DBSCAN to high dimensional spaces by incorporating the notion of dimension preferences in the distance function

- For each point p, it defines its subspace preference vector:

$$\bar{\mathbf{w}}_p = (w_1, w_2, ... w_d) \qquad w_i = \begin{cases} 1 & if & \mathrm{VAR}_i > \delta \\ \kappa & if & \mathrm{VAR}_i \leq \delta \end{cases}$$

- $\mathrm{VAR}_i$ is the variance along dimension j ($A_i$) in $N_\varepsilon(p)$:

$$\mathrm{VAR}_{A_i}(\mathcal{N}_\varepsilon(p)) = \frac{\sum_{q \in \mathcal{N}_\varepsilon(p)} (dist(\pi_{A_i}(p), \pi_{A_i}(q)))^2}{|\mathcal{N}_\varepsilon(p)|}$$

*δ, κ (κ>>1) are input parameters*

- Preference weighted distance function:

$$dist_p(p,q) = \sqrt{\sum_{i=1}^{d} \frac{1}{w_i} \cdot (\pi_{A_i}(p) - \pi_{A_i}(q))^2}$$
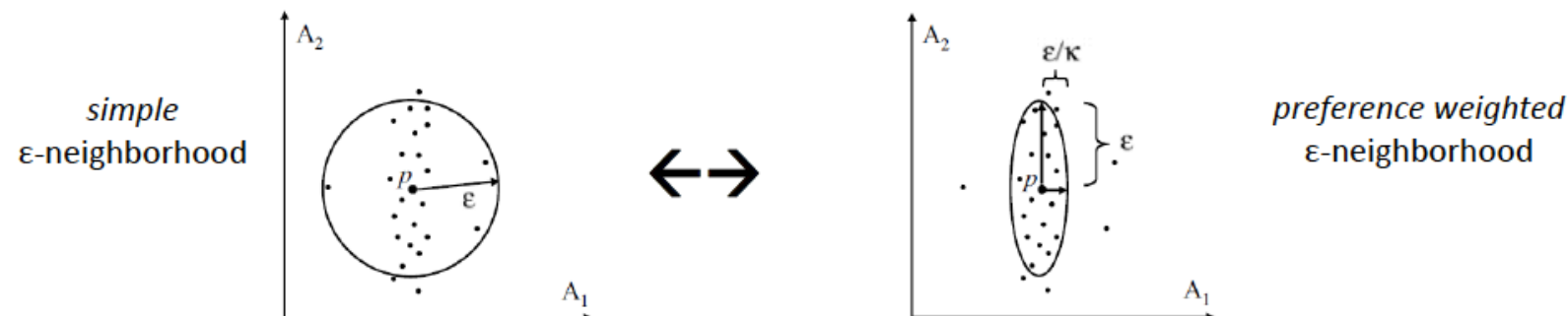
<span style="color:red">$w_i$</span>

Important dimensions weighted more heavily!

$$dist_{pref}(p,q) = \max\{dist_p(p,q), dist_q(q,p)\}$$

- Preference weighted ε-neighborhood:

$$\mathcal{N}_\varepsilon^{\bar{w}_p}(p) = \{x \in \mathcal{D} \mid dist_{pref}(p,x) \leq \varepsilon\}$$



simple ε-neighborhood ↔ preference weighted ε-neighborhood

- Preference weighted core points:

$$\mathrm{CORE}^{\mathrm{pref}}_{\mathrm{den}}(p) \Leftrightarrow \mathrm{PDIM}(\mathcal{N}_\varepsilon(p)) \leq \lambda \wedge |\mathcal{N}^{\bar{\mathrm{w}}}_\varepsilon(p)| \geq \mu$$

p

  p is core point      Subspace preference dimensionality      Preference weighted neigborhood
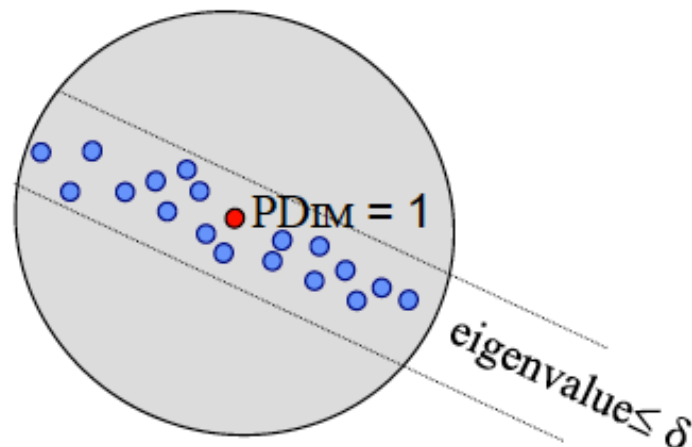
- Direct density reachability, reachability and connectivity are defined based on preference weighted core points

- A *subspace preference cluster* is a maximal density connected set of points associated with a certain subspace preference vector.

4C = Computing Correlation Connected Clusters

Idea: Integrate PCA into density-based clustering.

Approach:
- Check the core point property of a point p in the complete feature space
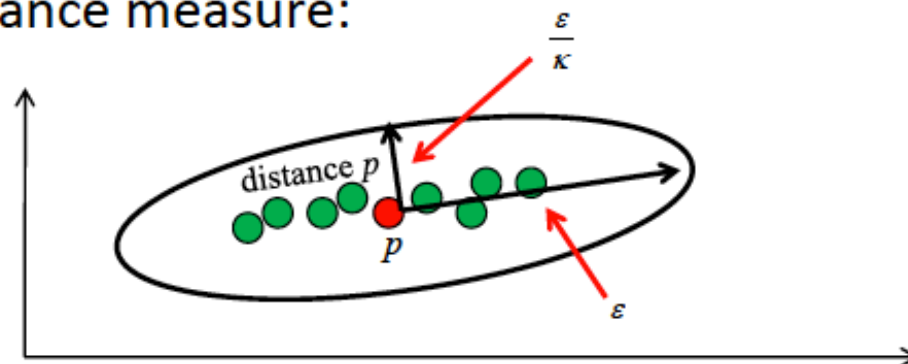- Perform PCA on the local neighborhood S of p to find subspace correlations



PCA factorizes $M_p$ into $M_p = V \, E \, V^T$

V:     eigenvectors

E:     eigenvalues

- A parameter δ discerns large from small eigenvalues.
- CorDim(S)=#eigenvalues>δ
- In the eigenvalue matrix of p, large eigenvalues are replaced by 1, small eigenvalues by a value κ >>1 → adapted eigenvalue matrix $E'_p$

# 4C: Distance measure

- effect on distance measure:



$$\hat{e}_i = \begin{cases} 1 & if \quad \Omega(e_i) > \delta \\ \kappa & if \quad \Omega(e_i) \leq \delta \end{cases} \qquad where \; \Omega \; is \; the \; normalization \; of \; the \; eigenvalues \; onto \; [0,1]$$

- distance of $p$ and $q$ w.r.t. $p$:  $\sqrt{(p-q)\cdot V_p \cdot E'_p \cdot V_p^{\mathrm{T}} \cdot (p-q)^{\mathrm{T}}}$

- distance of $p$ and $q$ w.r.t. $q$:  $\sqrt{(q-p)\cdot V_q \cdot E'_q \cdot V_q^{\mathrm{T}} \cdot (q-p)^{\mathrm{T}}}$

- symmetry of distance measure by choosing the maximum:

- *p* and *q* are correlation-neighbors if

$$\max\left\{\begin{array}{c}\sqrt{(p-q)\cdot V_p \cdot E_p' \cdot V_p^{\mathrm{T}} \cdot (p-q)^{\mathrm{T}}}, \\ \sqrt{(q-p)\cdot V_q \cdot E_q' \cdot V_q^{\mathrm{T}} \cdot (q-p)^{\mathrm{T}}}\end{array}\right\} \le \varepsilon$$

```
algorithm 4C(𝒟, ε, μ, λ, δ)

    // assumption: each object in 𝒟 is marked as unclassified

    for each unclassified O ∈ 𝒟 do

STEP 1. test CORE_den^cor(O) predicate:

        compute 𝒩_ε(O);
        if |𝒩_ε(O)| ≥ μ then
            compute M_O;                          ──────────── Covariance matrix   M_O = V_P E_P V_P^T
            if CORDIM(𝒩_ε(O)) ≤ λ then
                compute M̂_O and 𝒩_ε^{M̂_O}(O);    ──────────── Correlation similarity matrix   M̂_O = V_P Ê_P V_P^T
                test |𝒩_ε^{M̂_O}(O)| ≥ μ;

STEP 2.1. if CORE_den^cor(O) expand a new cluster:

        generate new clusterID;
        insert all X ∈ 𝒩_ε^{M̂_O}(O) into queue Φ;
        while Φ ≠ ∅ do
            Q = first object in Φ;
            compute ℛ = {X ∈ 𝒟 | DIRREACH_den^cor(Q, X)};
            for each X ∈ ℛ do
                if X is unclassified or noise then
                    assign current clusterID to X
                if X is unclassified then
                    insert X into Φ;
            remove Q from Φ;

STEP 2.2. if not CORE_den^cor(O) O is noise:

        mark O as noise;

end.
```

- Basic idea of CASH (= Clustering in Arbitrary Subspaces based on the Hough transform)

  - Transform each object into a so-called *parameter space* representing all possible subspaces accommodating this object (i.e. all hyper-planes through this object)

  - This parameter space is a *continuum* of all these subspaces

  - The subspaces are represented by a considerably small number of parameters

  - This transform is a generalization of the Hough Transform (which is designed to detect linear structures in 2D images) for arbitrary dimensions

- Transform

  - For each $d$-dimensional point $p$ there is an infinite number of $(d\text{-}1)$-dimensional hyper-planes through $p$

  - Each of these hyper-planes $s$ is defined by $(p, \alpha_1, \ldots, \alpha_{d-1})$, where $\alpha_1, \ldots, \alpha_{d-1}$ is the normal vector $\boldsymbol{n}_s$ of the hyper-plane $s$

  - The function $f_p(\alpha_1, \ldots, \alpha_{d-1}) = \delta_s = <p, \boldsymbol{n}_s>$ maps $p$ and $\alpha_1, \ldots, \alpha_{d-1}$ onto the distance $\delta_s$ of the hyper-plane $s$ to the origin

  - The parameter space plots the graph of this *function*



data space

Parameter space