



LUDWIG-
MAXIMILIANS-
UNIVERSITY
MUNICH


DEPARTMENT
INSTITUTE FOR
INFORMATICS


DATABASE
SYSTEMS
GROUP

Can Shared-Neighbor Distances Defeat the Curse of Dimensionality?

SSDBM 2010

Michael E. Houle¹, Hans-Peter Kriegel², Peer Kröger², Erich Schubert², Arthur Zimek²

¹National Institute of Informatics
Tokyo, Japan
meh@nii.ac.jp

²Ludwig-Maximilians-Universität München
Munich, Germany
<http://www.dbs.ifi.lmu.de>
{kriegel,kroegerp,schube,zimek}@dbs.ifi.lmu.de



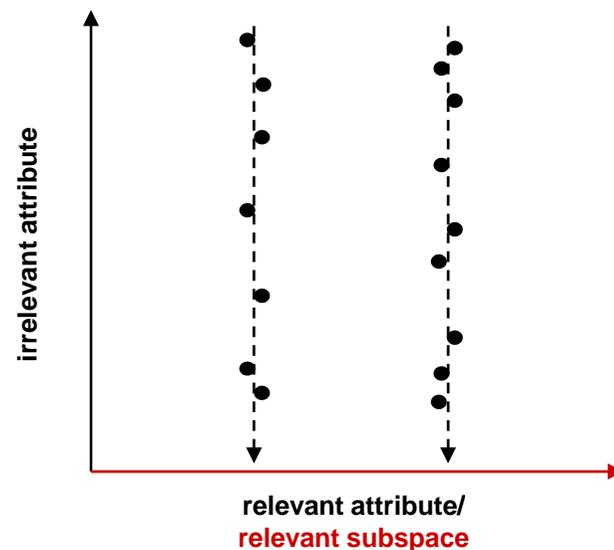
1. The Curse of Dimensionality
2. Shared-Neighbor Distances
3. Experimental Set-Up
4. Observations
5. Conclusions

- Beyer et al. (1999): distances to near and to far neighbors become more and more similar with increasing data dimensionality (loss of *relative contrast* or *concentration effect* of distances):

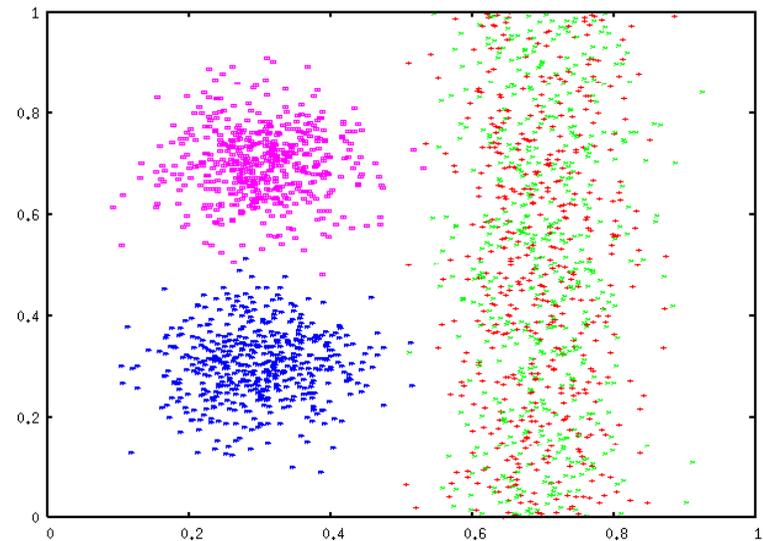
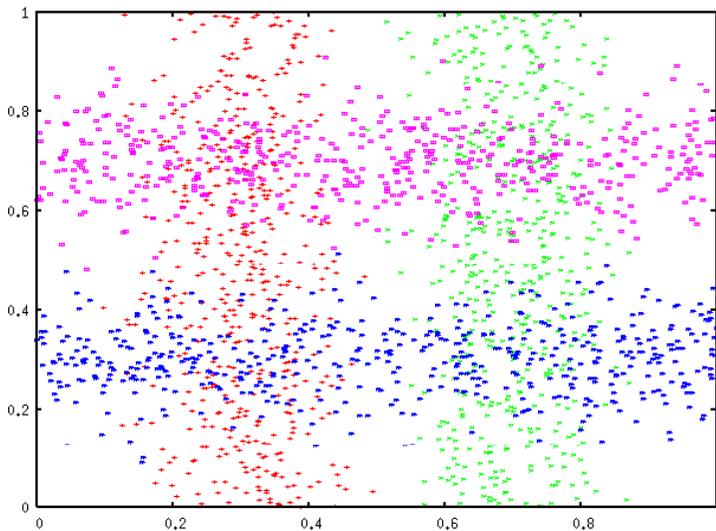
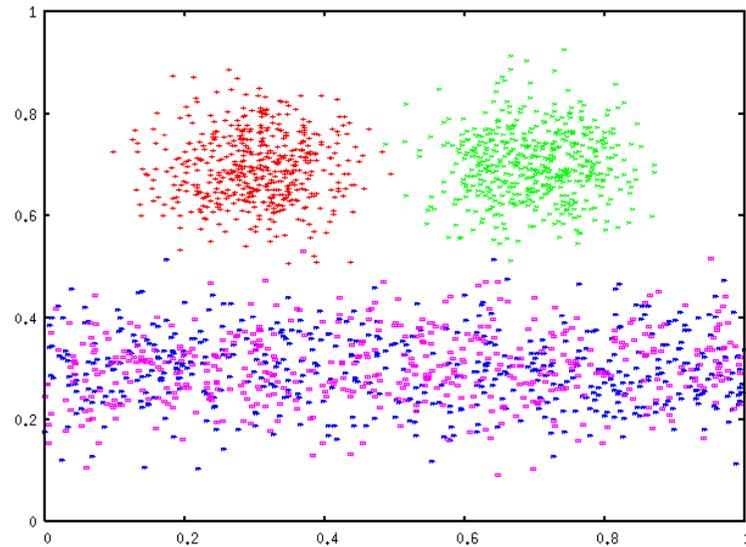
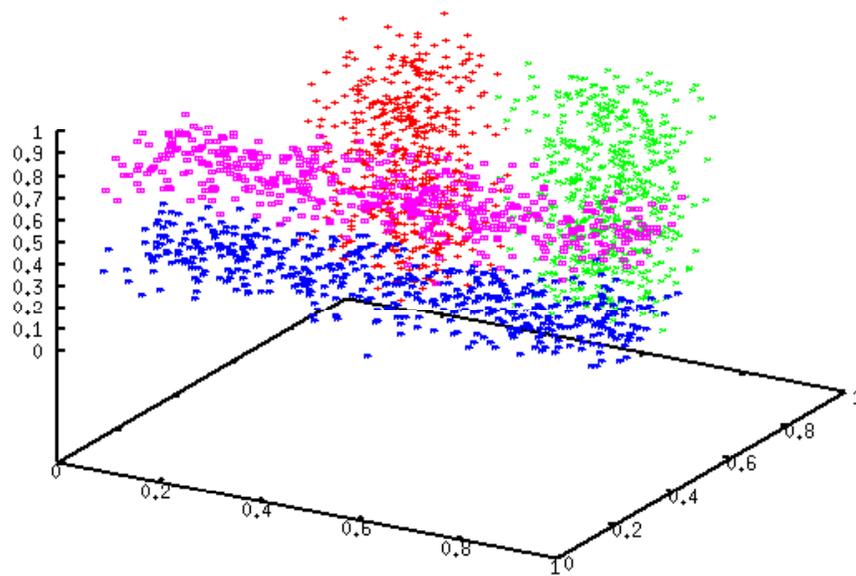
$$\lim_{d \rightarrow \infty} \text{var} \left(\frac{\|X_d\|}{E[\|X_d\|]} \right) = 0 \quad \Rightarrow \quad \frac{D_{\max} - D_{\min}}{D_{\min}} \rightarrow 0$$

- valid for a broad range of data distributions
- but only within one single distribution

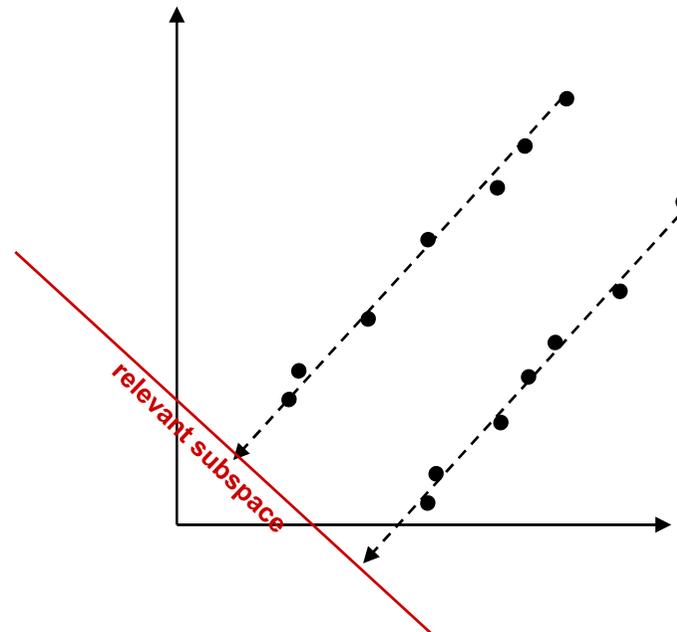
- Bennett et al. (1999): nearest-neighbor queries are still meaningful, if the search is limited to the same cluster and if the clusters are well separated.
- Separation of clusters relates to *relevant attributes* (helpful to discern between clusters) as opposed to *irrelevant attributes* (indistinguishable distribution of attribute values for different clusters).



The Curse of Dimensionality



- Redundant attributes: dependencies/correlations among attributes
 - can result in lower intrinsic dimensionality of a data set
 - bad discrimination of distances can still be a problem



- there are other effects of the “curse of dimensionality”
- we mainly aim at distinguishing these effects:
 - concentration effect within distributions
 - impediment of similarity search by irrelevant attributes
 - partly: impact of redundant/correlated attributes
- as a remedy for similarity assessment in high dimensional data, to use shared nearest neighbor (SNN) information has been proposed but never evaluated systematically
- here: evaluation of the effects on primary distances (Manhattan, Euclidean, fractional L_p ($L_{0.6}$ and $L_{0.8}$), cosine) and secondary distances (SNN)

- secondary distances are defined on top of primary distances
- shared nearest neighbor (SNN) information:
 - assess the set of s nearest neighbors for two objects x and y in terms of some primary distance (Euclidean, Manhattan, cosine...)
 - derive overlap of neighbors (common objects in the NN of x and y)

$$\text{SNN}_s(x, y) = |\text{NN}_s(x) \cap \text{NN}_s(y)|$$

- similarity measure

$$\text{simcos}_s(x, y) = \frac{\text{SNN}_s(x, y)}{s}$$

cosine of the angle between membership vectors for $\text{NN}(x)$ and $\text{NN}(y)$

- SNN has been used before in mining high-dimensional data, but alleged quality improvement has never been evaluated

Shared-Neighbor Distances

- distance measures based on SNN:

$$\text{dinv}_s(x, y) = 1 - \text{simcos}_s(x, y)$$

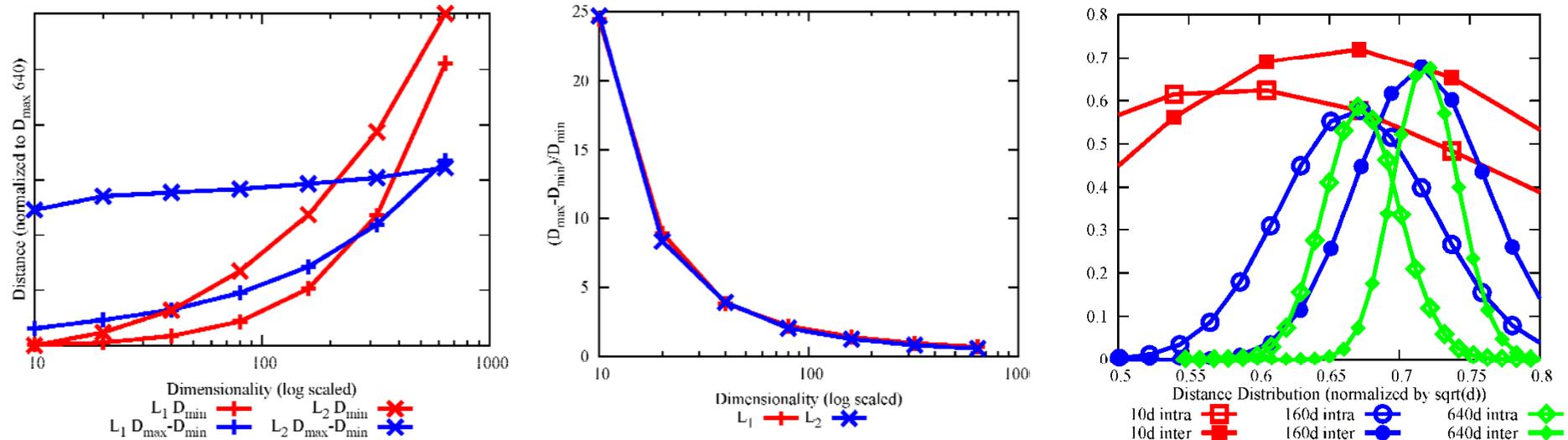
$$\text{dacos}_s(x, y) = \arccos(\text{simcos}_s(x, y))$$

$$\text{dln}_s(x, y) = -\ln(\text{simcos}_s(x, y))$$

- dinv : linear inversion
 - dacos penalizes slightly suboptimal similarities more strongly
 - dln more tolerant for relatively high similarity values but approaches infinity for very low similarity values
- for assessment of ranking quality, these formulations are equivalent as the ranking is unaffected
 - only dacos is a metric (if the underlying primary distance is a metric)

- Artificial data sets: $n = 10.000$ items, $c = 100$ clusters, up to $d = 640$ dimensions, cluster sizes randomly determined.
- Relevant attribute values normally distributed, irrelevant attribute values uniformly distributed.
- Data sets:
 - *All-Relevant*: all dimensions relevant for all clusters
 - *10-Relevant*: first 10 dimensions are relevant for all clusters, the remaining dimensions are irrelevant
 - *Cyc-Relevant*: i th attribute is relevant for the j th cluster when $i \bmod c = j$, otherwise irrelevant (here: $c = 10$, $n = 1000$)
 - *Half-Relevant*: for each cluster, an attribute is chosen to be relevant with probability 0.5, and irrelevant otherwise
 - *All-Dependent*: derived from *All-Relevant* introducing correlations among attributes
 $X_{i \in AllDependent}, Y_{i \in AllRelevant}: X_i = Y_i (1 \leq i \leq 10), X_i = \frac{1}{2} (X_{i-10} + Y_i) (i > 10)$
 - *10-Dependent*: derived from *10-Relevant* introducing correlations among attributes

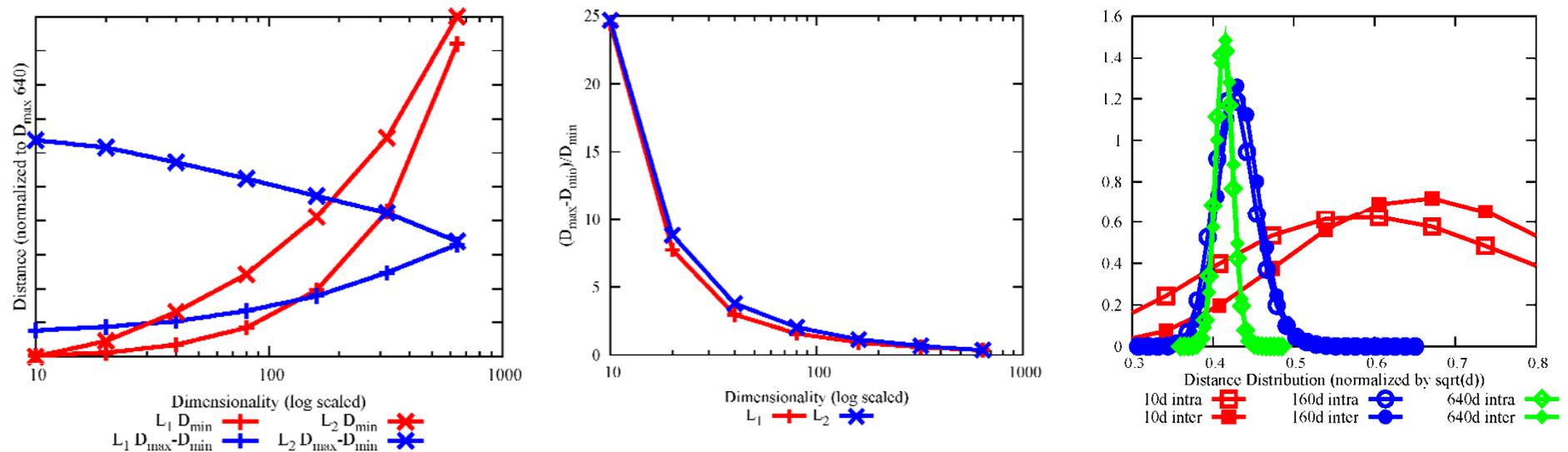
Data sets show properties of the “curse of dimensionality”



All-Relevant

$$\lim_{d \rightarrow \infty} \text{var} \left(\frac{\|X_d\|}{E[\|X_d\|]} \right) = 0 \Rightarrow \frac{D_{\max} - D_{\min}}{D_{\min}} \rightarrow 0$$

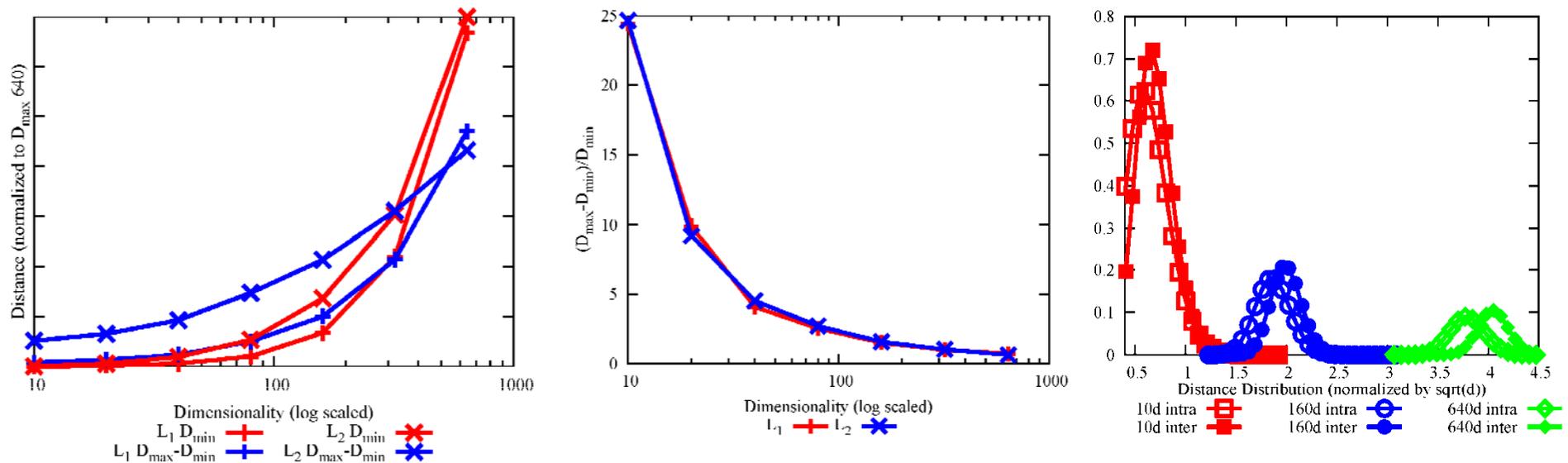
Data sets show properties of the “curse of dimensionality”



10-Relevant

$$\lim_{d \rightarrow \infty} \text{var} \left(\frac{\|X_d\|}{E[\|X_d\|]} \right) = 0 \Rightarrow \frac{D_{\max} - D_{\min}}{D_{\min}} \rightarrow 0$$

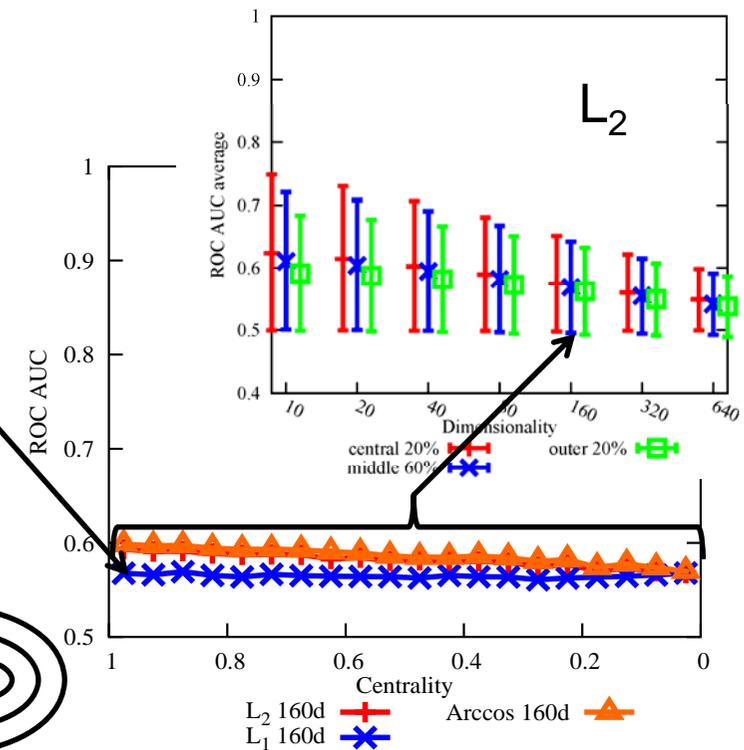
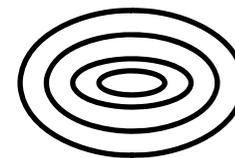
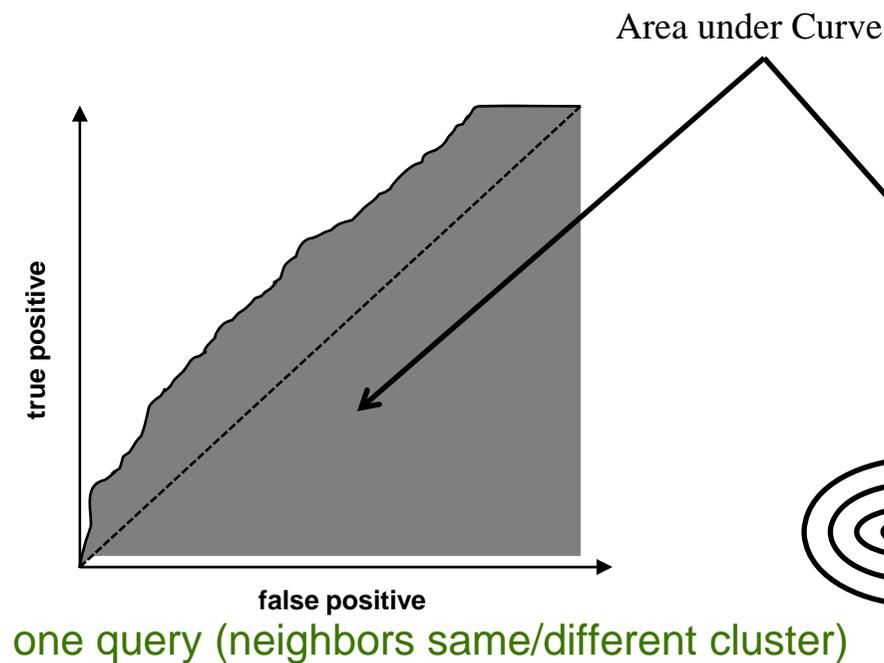
Data sets show properties of the “curse of dimensionality”



All-Dependent

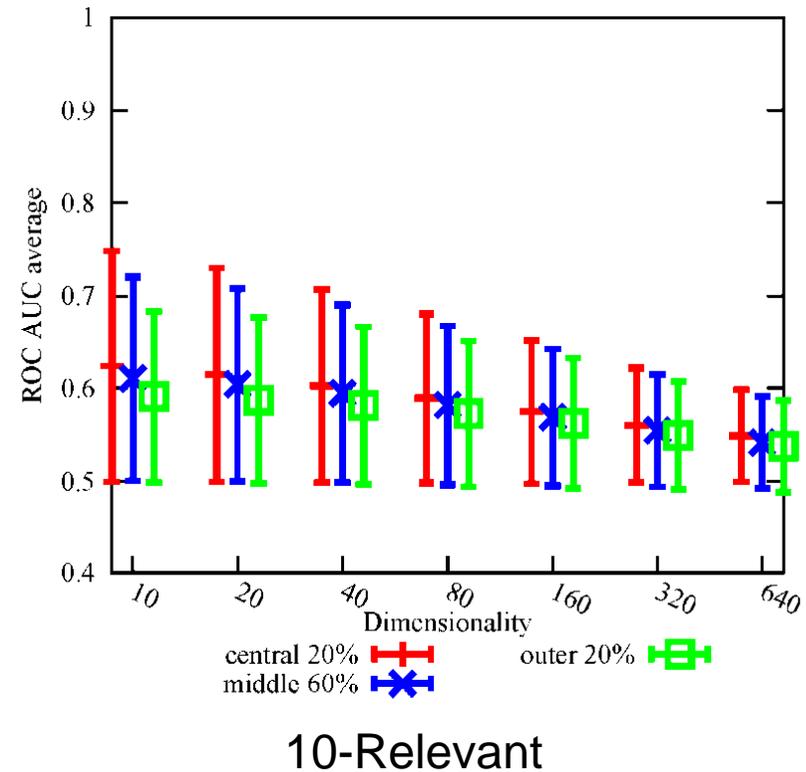
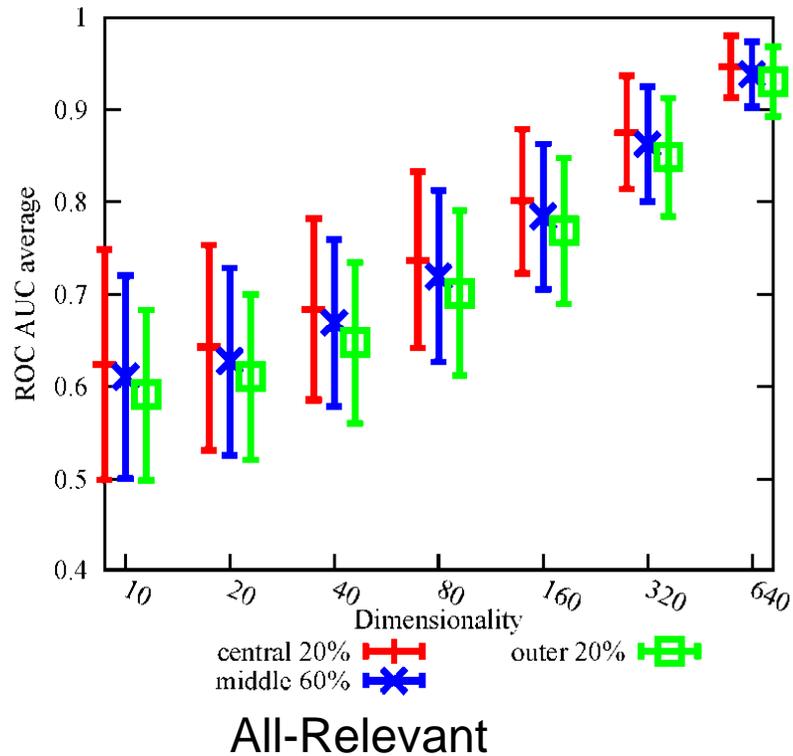
$$\lim_{d \rightarrow \infty} \text{var} \left(\frac{\|X_d\|}{E[\|X_d\|]} \right) = 0 \Rightarrow \frac{D_{\max} - D_{\min}}{D_{\min}} \rightarrow 0$$

- Using each item in turn as a query, neighborhood ranking reported in terms of the Area under curve (AUC) of the Receiver Operating Characteristic (ROC)

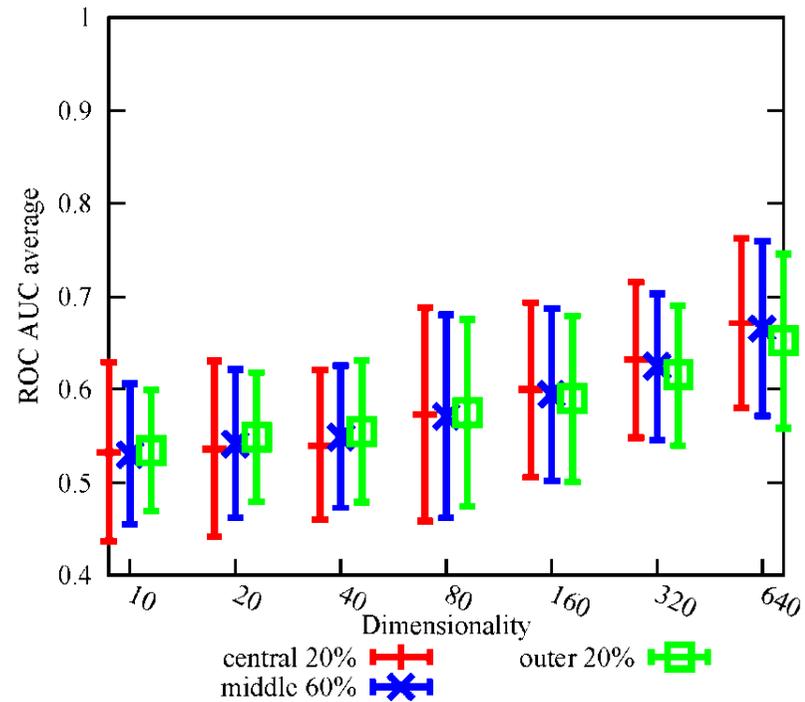


average over all items as a query

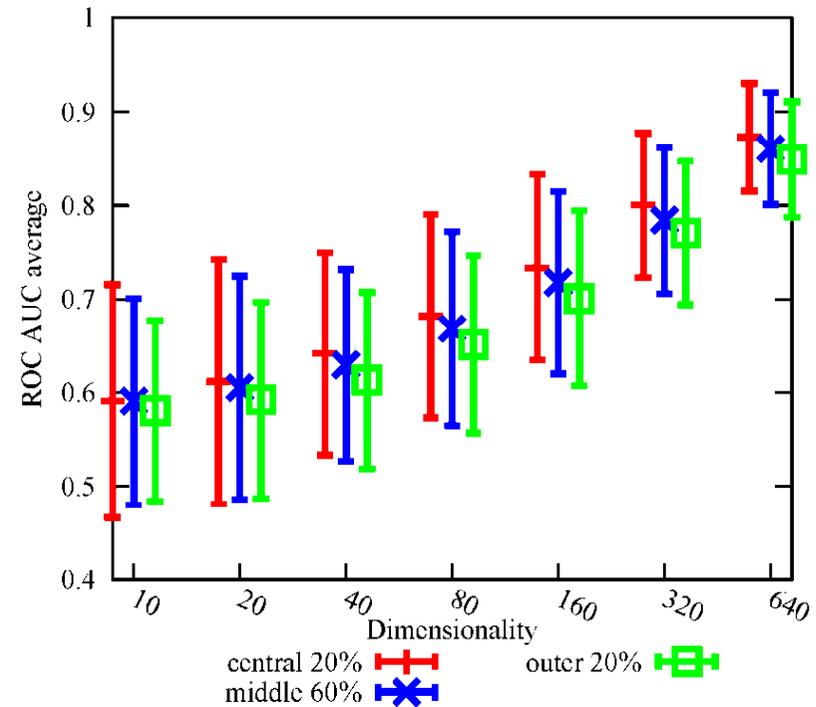
Euclidean distance



Euclidean distance

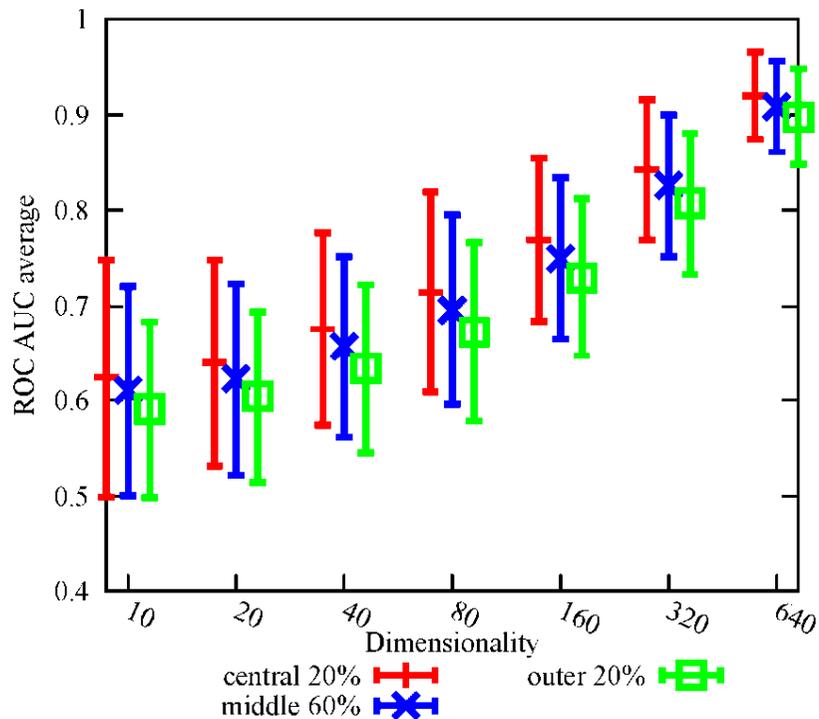


Cyc-Relevant

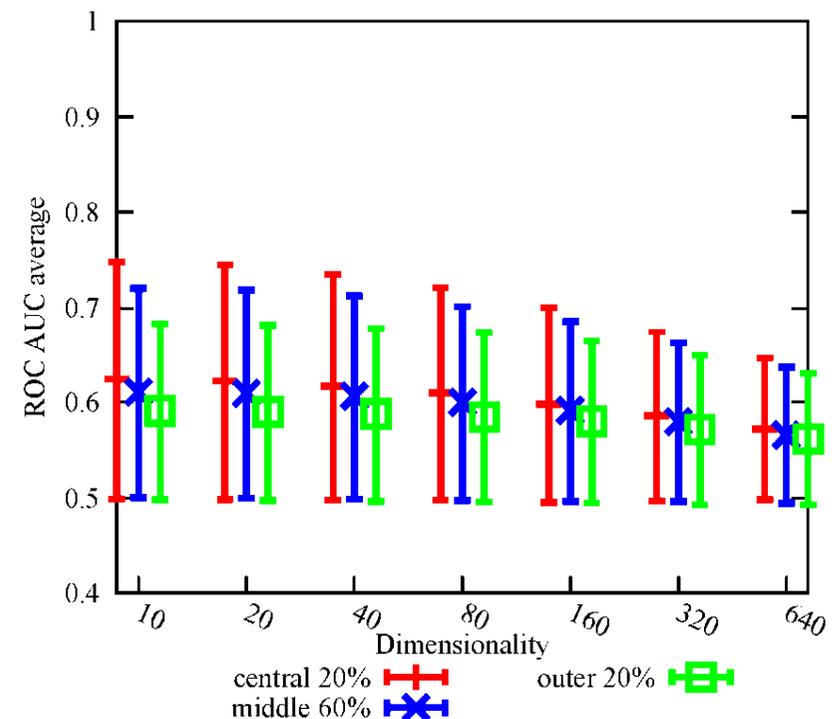


Half-Relevant

Euclidean distance



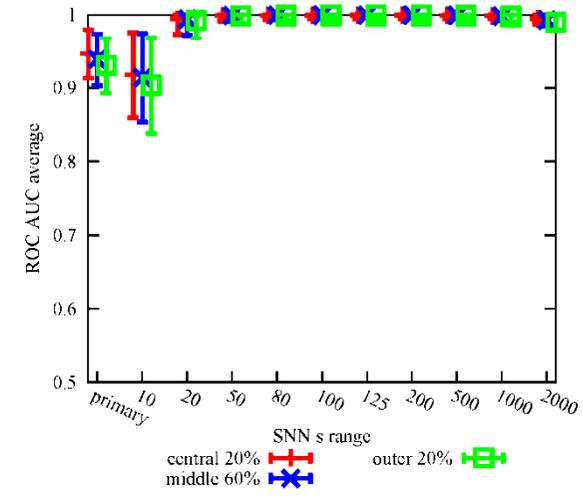
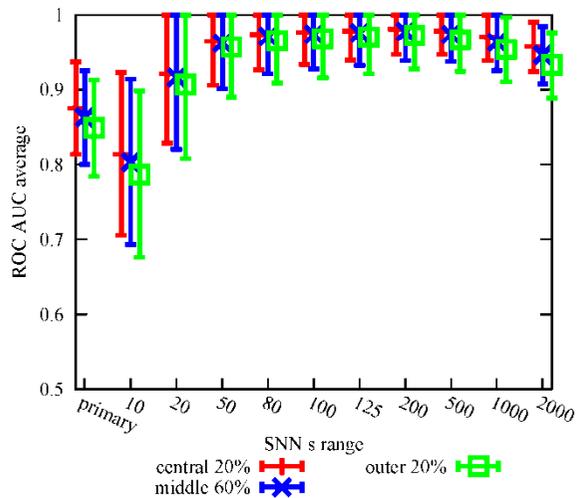
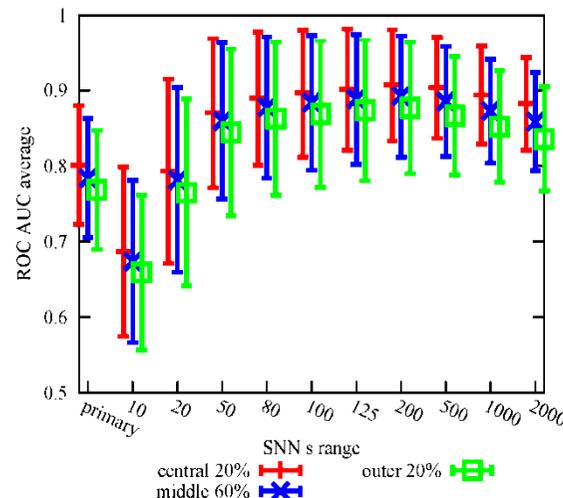
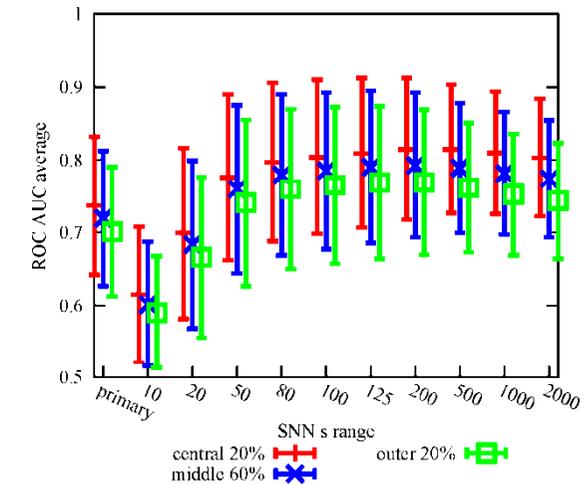
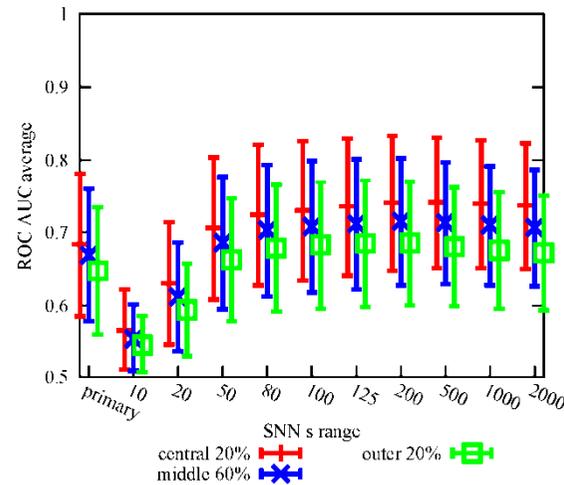
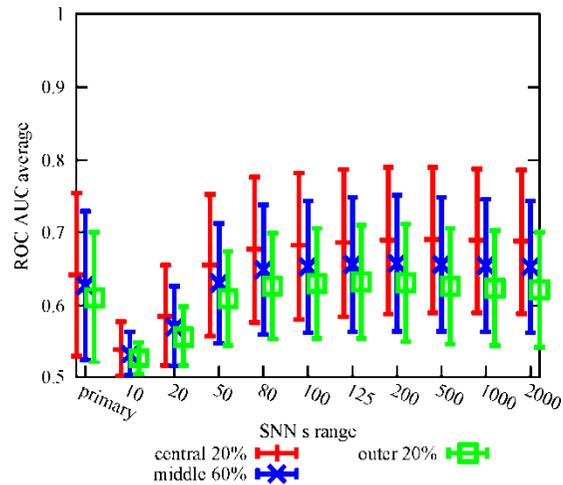
All-Dependent



10-Dependent

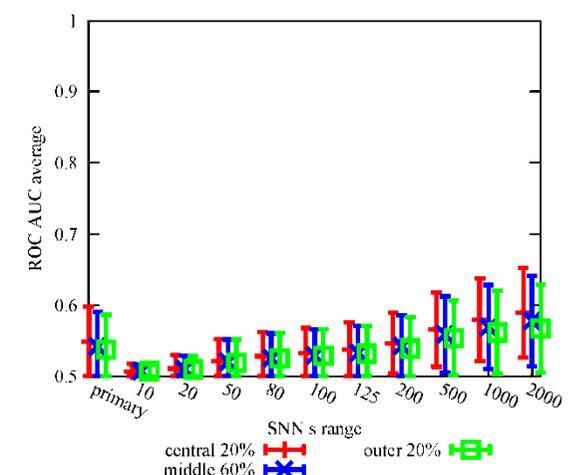
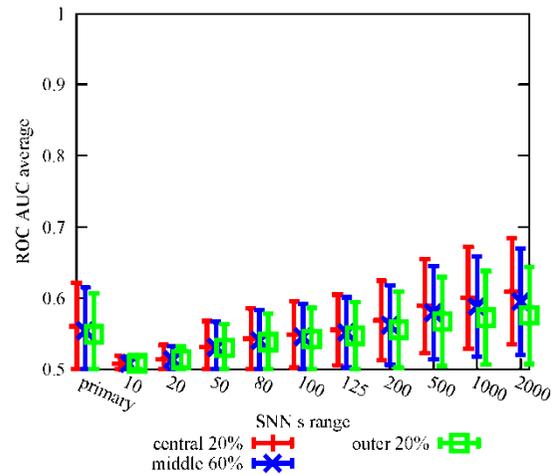
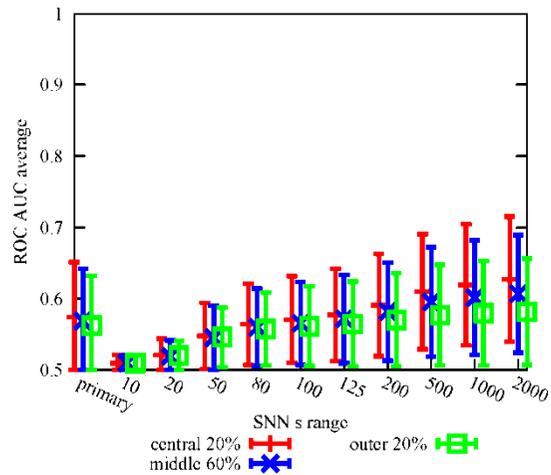
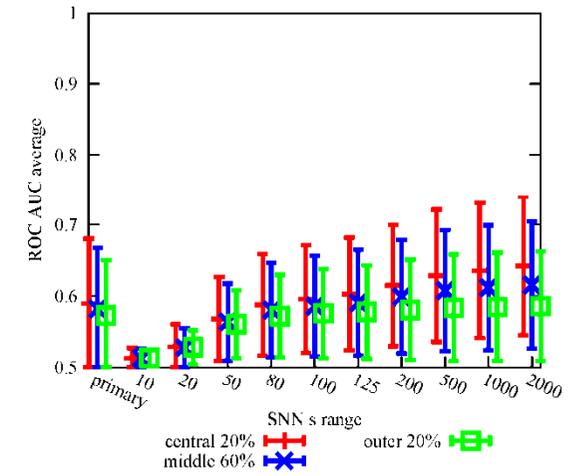
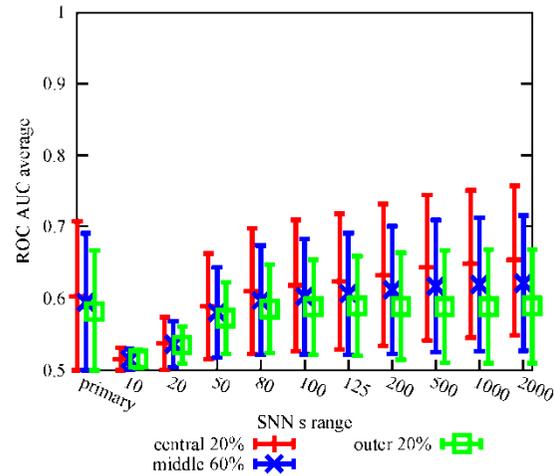
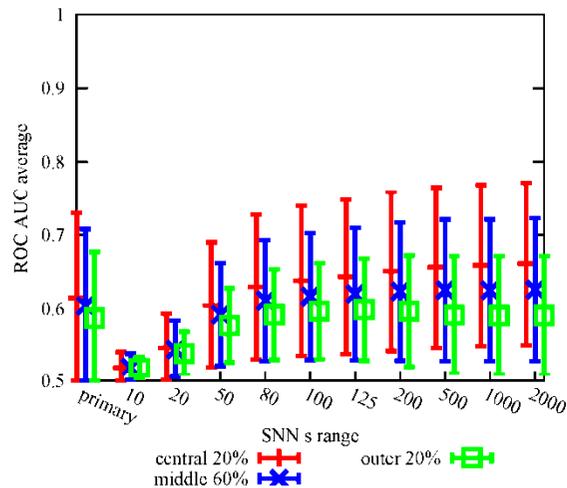
SNN based on Euclidean

All-Relevant
20/40/80/160/320/640 dimensions



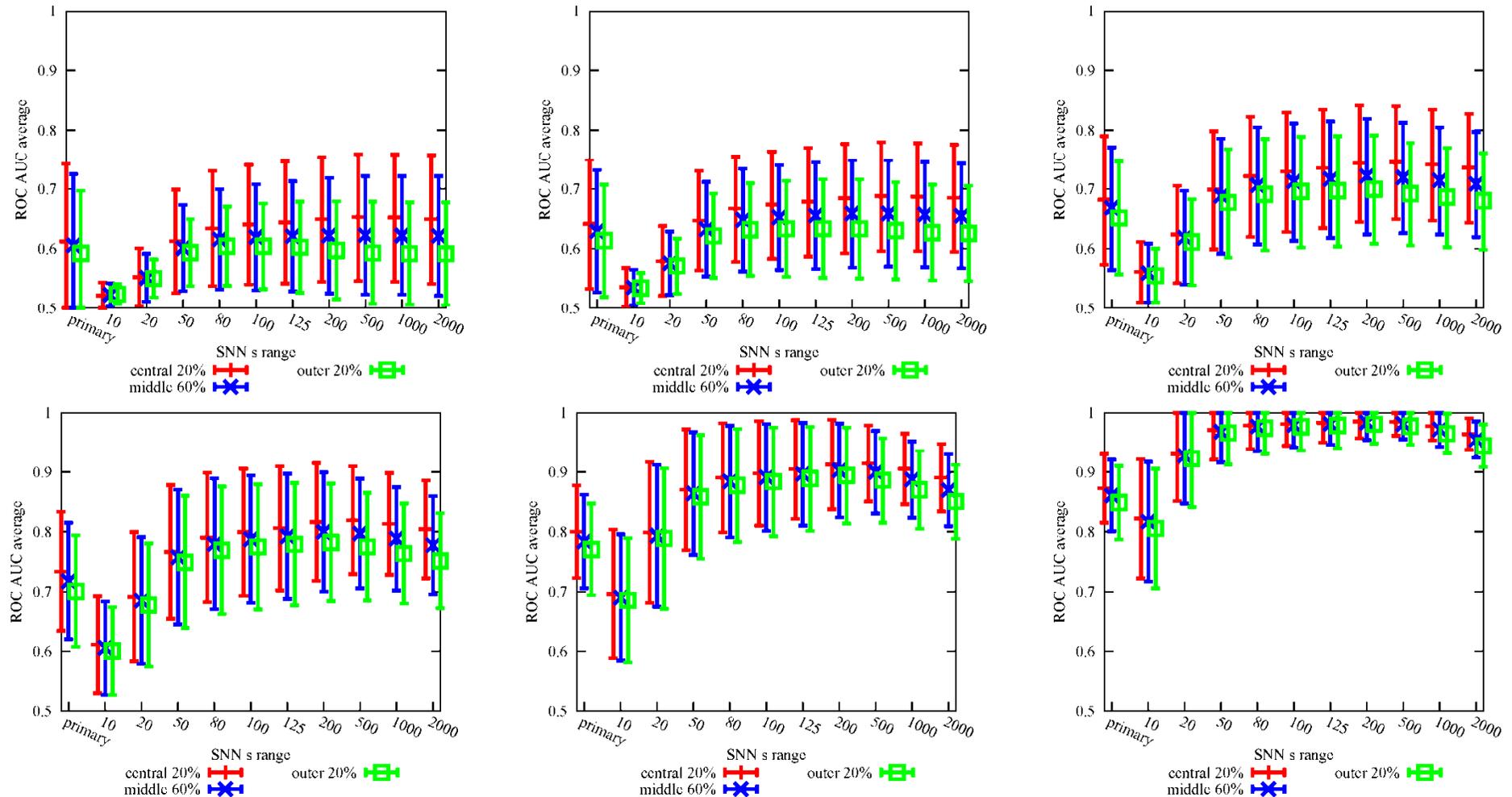
SNN based on Euclidean

10-Relevant
20/40/80/160/320/640 dimensions



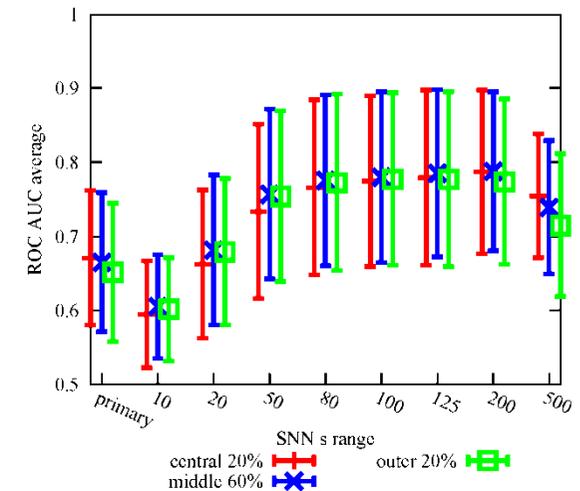
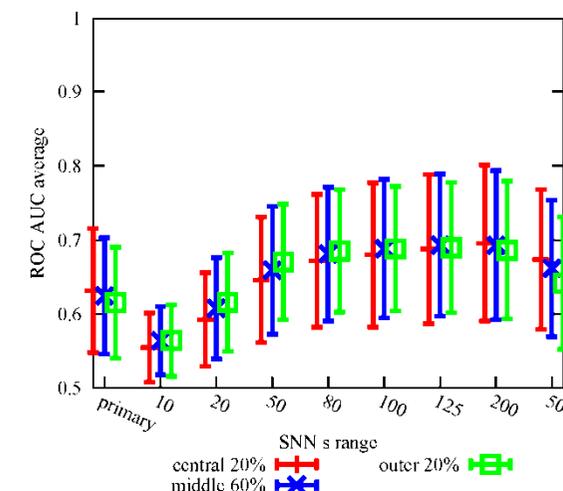
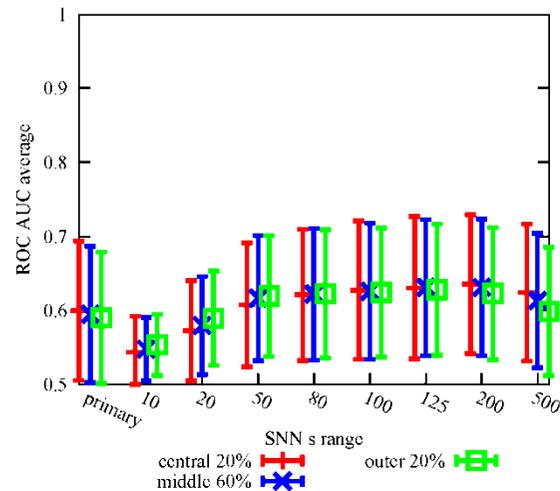
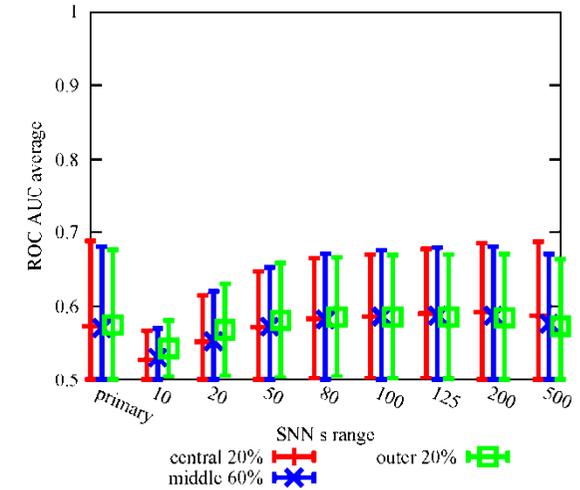
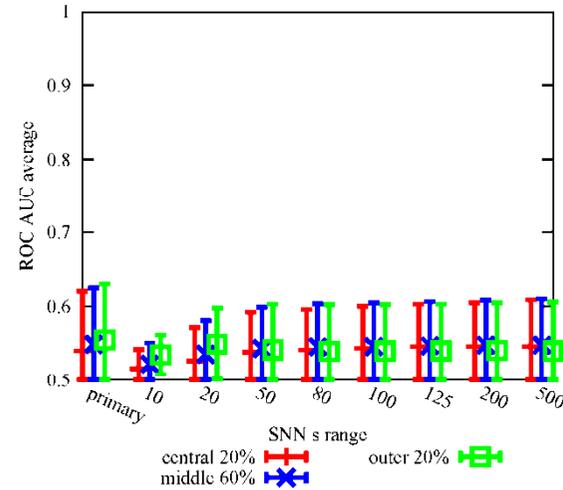
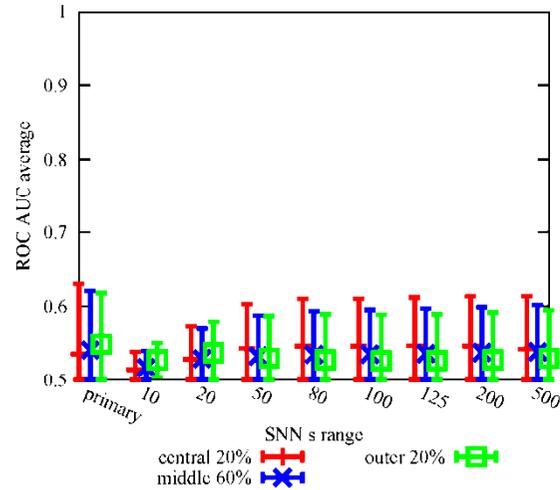
SNN based on Euclidean

Half-Relevant
20/40/80/160/320/640 dimensions

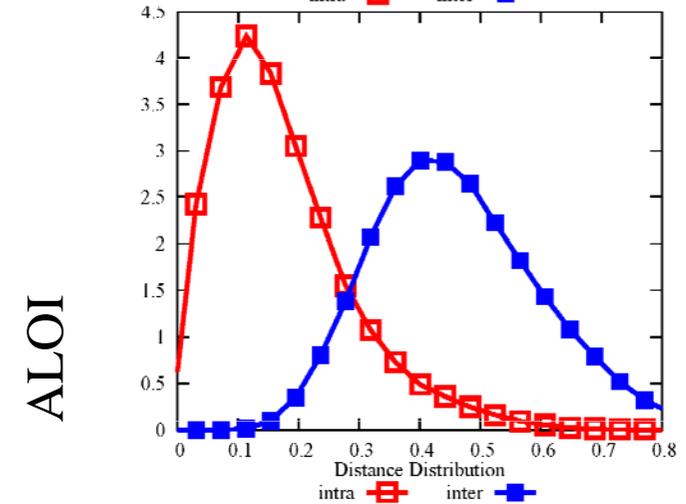
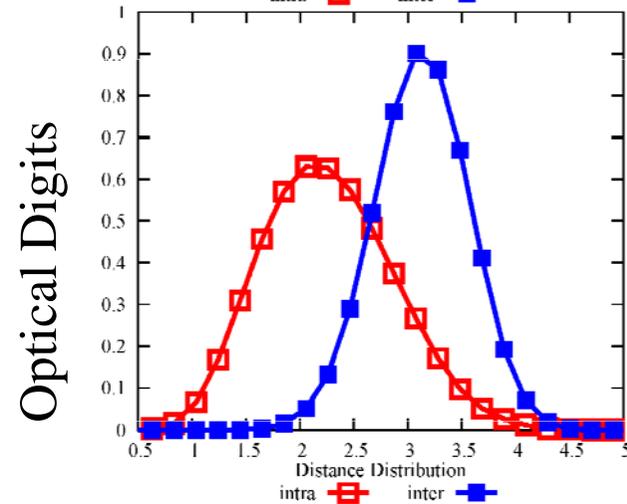
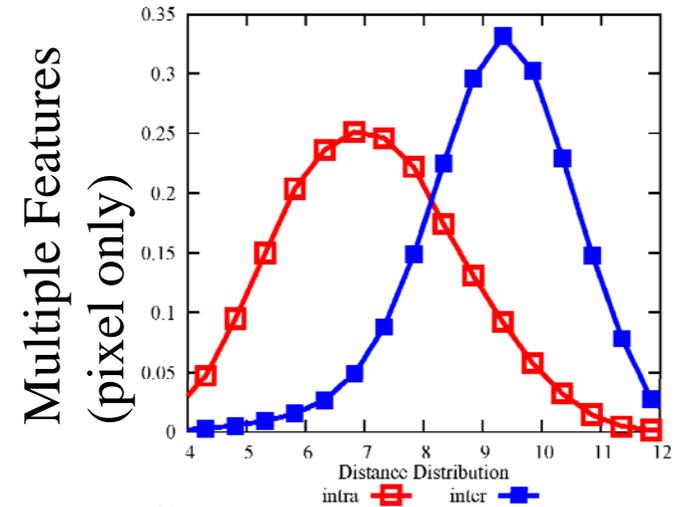
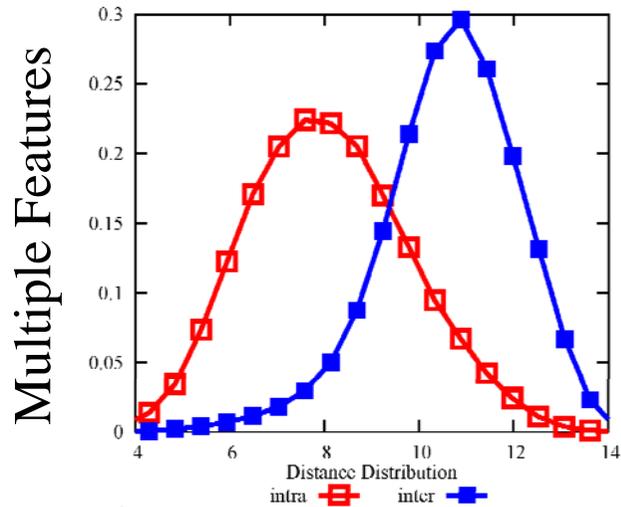


SNN based on Euclidean

Cyc-Relevant
20/40/80/160/320/640 dimensions

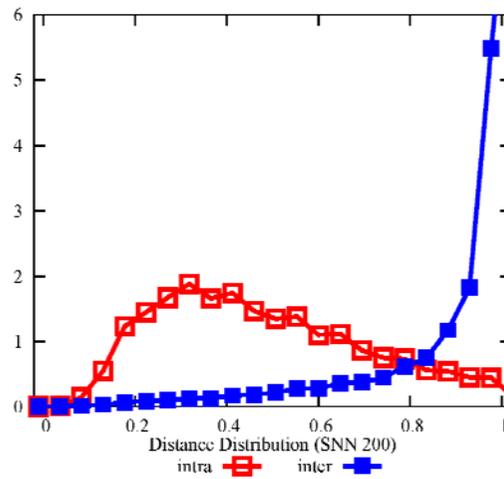


some real data sets: distributions of Euclidean distances

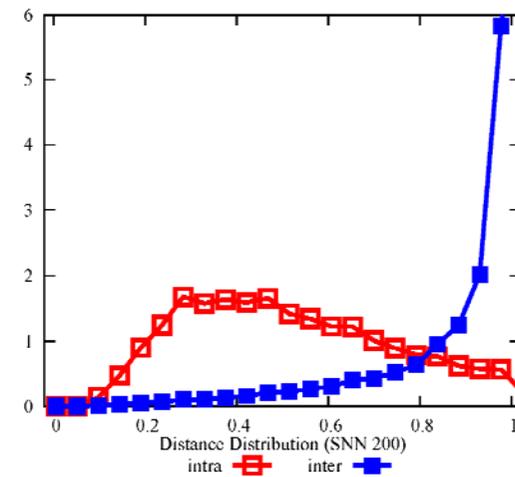


some real data sets: distributions of SNN distances (Euclidean)

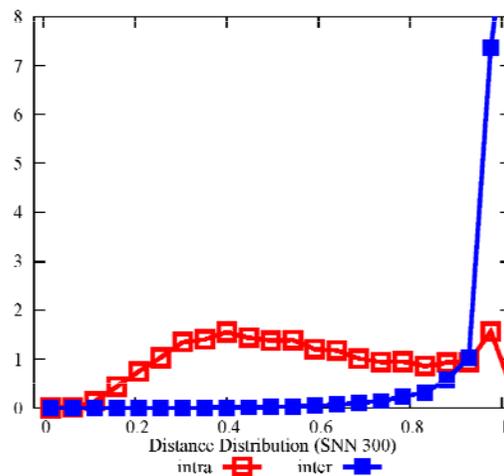
Multiple Features



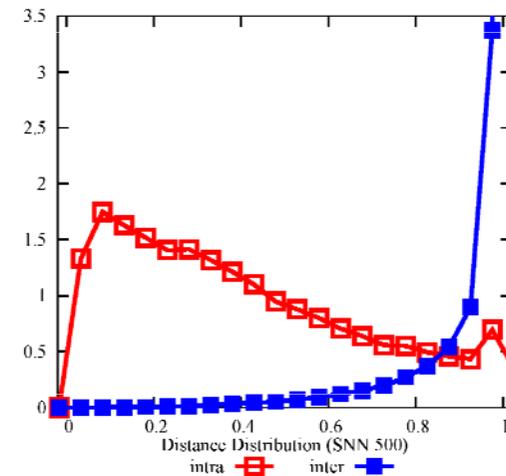
Multiple Features
(pixel only)



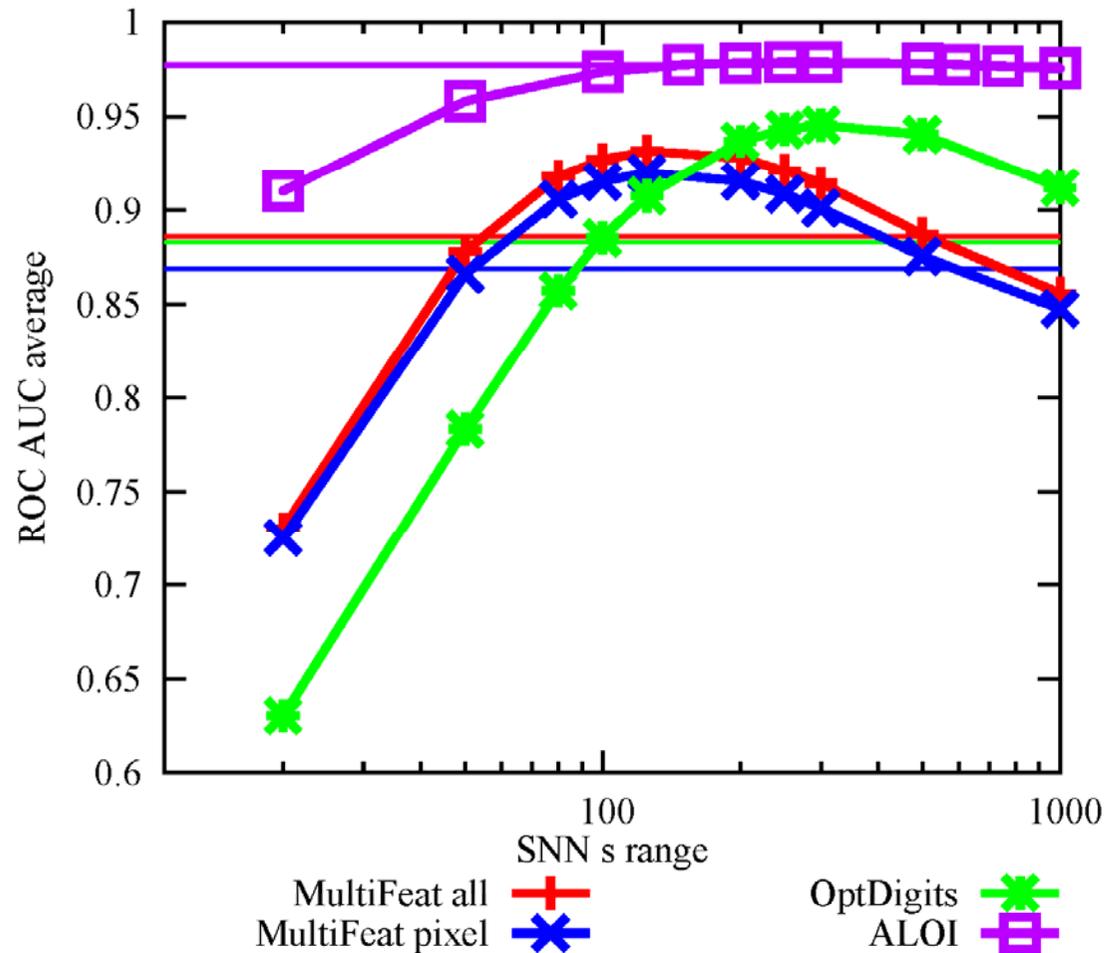
Optical Digits



ALOI



some real data sets: ranking quality



- The *curse of dimensionality* does not count in general as an excuse for everything – depends on the number and nature of distributions in a data set
- the nature of each particular problem needs to be studied in its own – part of the curse: it's always different than expected
- SNN information *can* improve neighborhood ranking for even very low quality of neighborhood queries
 - if the primary distance already performs good, the improvement by SNN in many cases seems actually to be more significant
 - open questions:
 - good choice of neighborhood size s : relationship between s and size of natural clusters?
 - k NN query based on SNN_s : relationship between k and s ?

supplementary material:

<http://www.dbs.ifi.lmu.de/research/SNN/>