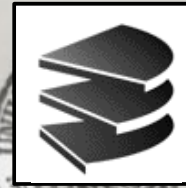




LUDWIG-  
MAXIMILIANS-  
UNIVERSITÄT  
MÜNCHEN



INSTITUTE FOR  
INFORMATICS  
DATABASE  
GROUP



# SIMILARITY ESTIMATION USING BAYES ENSEMBLES

Marisa Thoma

*joint work with*

Franz Graf, Tobias Emrich, Hans-Peter Kriegel and Matthias Schubert

LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN  
DATABASE GROUP



Scientific and Statistical Database Management  
22<sup>nd</sup> International Conference, SSDBM 2010  
Heidelberg, Germany, June/July 2010

- SIMILARITY MEASURES
  - Requirements and pitfalls
- BAYES ENSEMBLE DISTANCE (BED)
  - Bayes Estimates (BE)
  - Relevance weights
  - Feature space optimization
- EXPERIMENTS
- SUMMARY AND OUTLOOK

## PAIR-WISE OBJECT COMPARISON

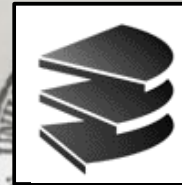
Objects  $x_1, x_2 \in \mathbb{R}^d$

$s(x_1, x_2)$  : Similarity of  $x_1$  and  $x_2$

$s(x_1, x_2) > s(x_1, x_3) \Rightarrow x_1$  more similar to  $x_2$  than to  $x_3$

## APPLICATIONS:

- Similarity / retrieval queries (ranking, range queries)
- Clustering
- Model Training



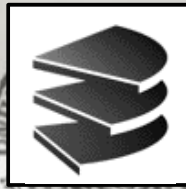
## PAIR-WISE OBJECT COMPARISON

Distance Measures:  $s(x_1, x_2) \sim 1 / d(x_1, x_2)$

e.g.  $L_p$ -norm: 
$$L_p(x_1, x_2) = \left( \sum_{i=1}^d |x_{1,i} - x_{2,i}|^p \right)^{1/p}$$

### PROBLEMS:

- Target similarity is ignored
- Correlated Dimensions
- Irrelevant Dimensions
- 1 Dimension has a large influence on small distances
- Distances can grow arbitrarily large



## PAIR-WISE OBJECT COMPARISON

Distance Measures:  $s(x_1, x_2) \sim 1 / d(x_1, x_2)$

e.g.  $L_p$ -norm:  $L_p(x_1, x_2) = \left( \sum_{i=1}^d |x_{1,i} - x_{2,i}|^p \right)^{1/p}$

## PROBLEMS:

- Target similarity is ignored
- Correlated Dimensions
- Irrelevant Dimensions
- 1 Dimension has a large influence on small distances
- Distances can grow arbitrarily large



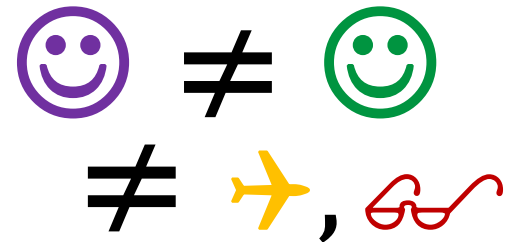
## PAIR-WISE OBJECT COMPARISON

Distance Measures:  $s(x_1, x_2) \sim 1 / d(x_1, x_2)$

e.g.  $L_p$ -norm: 
$$L_p(x_1, x_2) = \left( \sum_{i=1}^d |x_{1,i} - x_{2,i}|^p \right)^{1/p}$$

### PROBLEMS:

- Target similarity is ignored
- Correlated Dimensions
- Irrelevant Dimensions
- 1 Dimension has a large influence on small distances
- Distances can grow arbitrarily large



## PAIR-WISE OBJECT COMPARISON

Distance Measures:  $s(x_1, x_2) \sim 1 / d(x_1, x_2)$

e.g.  $L_p$ -norm:  $L_p(x_1, x_2) = \left( \sum_{i=1}^d |x_{1,i} - x_{2,i}|^p \right)^{1/p}$

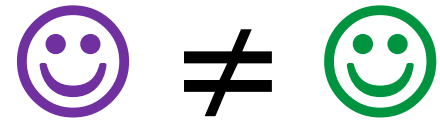
### PROBLEMS:

- Target similarity is ignored
- Correlated Dimensions
- Irrelevant Dimensions

*Can be solved via **Mahalanobis distance**:*

$$D_{\text{Mah}}(x_1, x_2) = \left( (x_1 - x_2)^T A (x_1 - x_2) \right)^{1/2}$$

- 1 Dimension has a large influence on small distances



- Distances can grow arbitrarily large







## PAIR-WISE OBJECT COMPARISON

Distance Measures:  $s(x_1, x_2) \sim 1 / d(x_1, x_2)$

e.g.  $L_p$ -norm: 
$$L_p(x_1, x_2) = \left( \sum_{i=1}^d |x_{1,i} - x_{2,i}|^p \right)^{1/p}$$

### PROBLEMS:

- Target similarity is ignored
  - Correlated Dimensions
  - Irrelevant Dimensions
- Can be solved via **Mahalanobis distance**:*
- $$D_{\text{Mah}}(x_1, x_2) = \left( (x_1 - x_2)^T A (x_1 - x_2) \right)^{1/2}$$

- 1 Dimension has a large influence on small distances   $\neq$  
- Distances can grow arbitrarily large  $\neq$   , 



## DISSIMILARITY PROBABILITIES OF OBJECT PAIRS:

Bayes classifier for a feature  $i$ :

$x_1, x_2$  dissimilar (DIS):  $P(\text{DIS} \mid (x_{1,i} - x_{2,i}))$

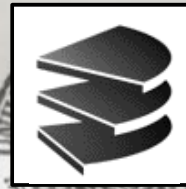
## BAYES ESTIMATE:

Priors for similar ( $p_{\text{SIM}}$ ) and dissimilar objects ( $p_{\text{DIS}}$ ):

$$\text{BE}_i(x_1, x_2) = \frac{p_{\text{DIS}} \cdot P((x_{1,i} - x_{2,i}) \mid \text{DIS})}{p_{\text{DIS}} \cdot P((x_{1,i} - x_{2,i}) \mid \text{DIS}) + p_{\text{SIM}} \cdot P((x_{1,i} - x_{2,i}) \mid \text{SIM})}$$

## BAYES ENSEMBLE DISTANCE:

$$\text{BED}(x_1, x_2) = \frac{1}{d} \cdot \sum_{i=1}^d \text{BE}_i(x_1, x_2)$$



BAYES ENSEMBLE DISTANCE:

$$\text{BED}(x_1, x_2) = \frac{1}{d} \cdot \sum_{i=1}^d \text{BE}_i(x_1, x_2)$$

NAÏVE BAYES CLASSIFIER:

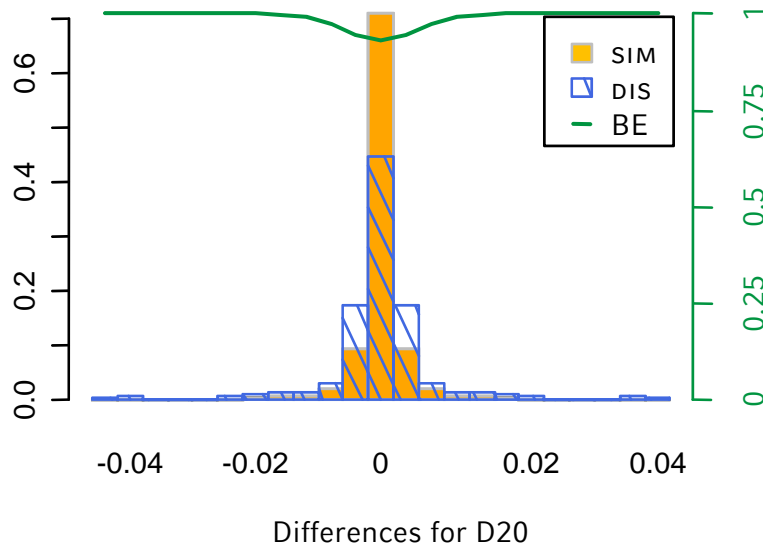
$$\text{NB}(x_1, x_2) = \frac{1}{\text{scale}} \cdot \prod_{i=1}^d \text{BE}_i(x_1, x_2)$$

- BED also results in scores in  $[0,1]$
- BED is more stable against outlier dimensions

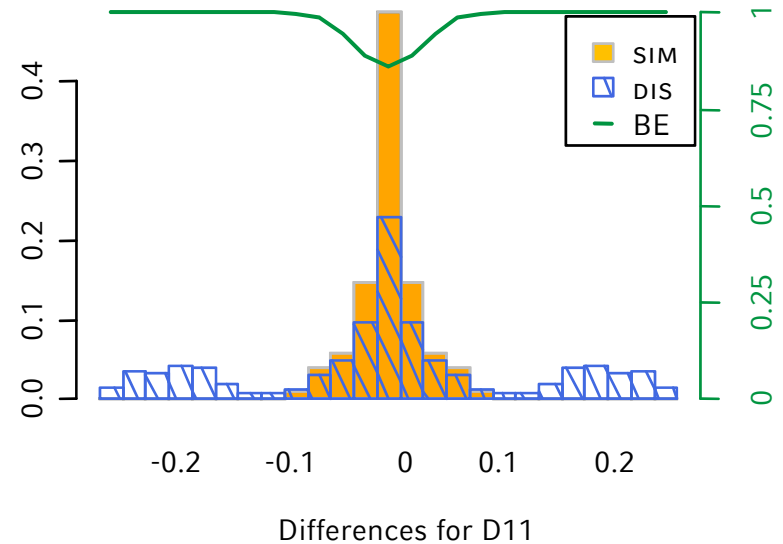
## EXAMPLES:

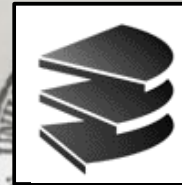
## COLOR HISTOGRAM DATA

### SUBOPTIMAL DIMENSION



### GOOD DIMENSION





INTRODUCE FEATURE WEIGHTS:

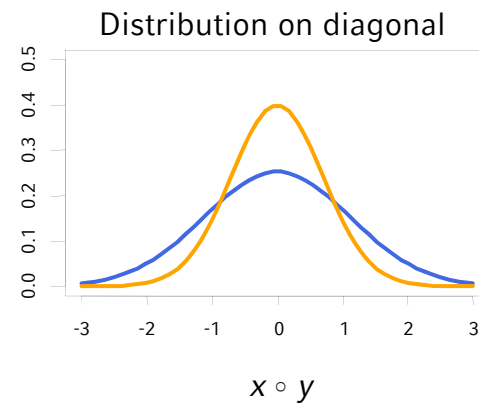
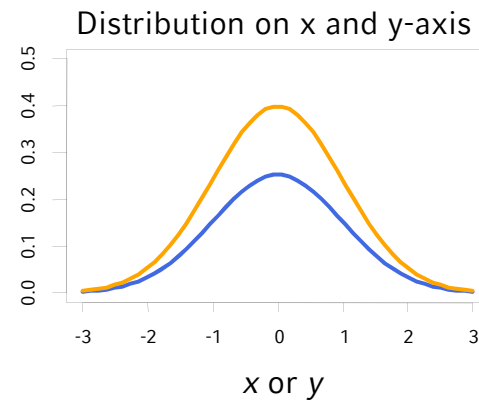
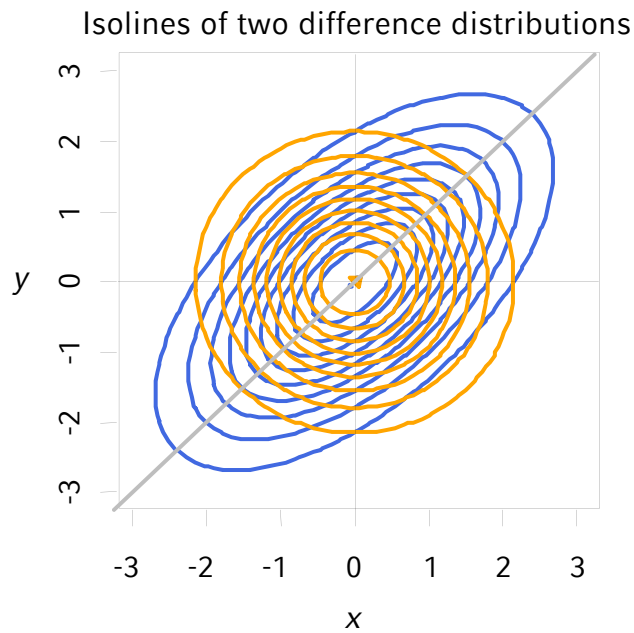
$$\text{BED}(x_1, x_2) = \left( \sum_{i=1}^d q_i \right)^{-1} \cdot \sum_{i=1}^d q_i \cdot \text{BE}_i(x_1, x_2)$$

ASSUMING A GAUSSIAN DISTRIBUTION:

Use variance difference as quality measure

$$q_i = \sigma_{\text{DIS}}^2 - \sigma_{\text{SIM}}^2 = \text{avg}_{x_d \in \text{DIS}}(x_{d,i}^2) - \text{avg}_{x_s \in \text{SIM}}(x_{s,i}^2)$$

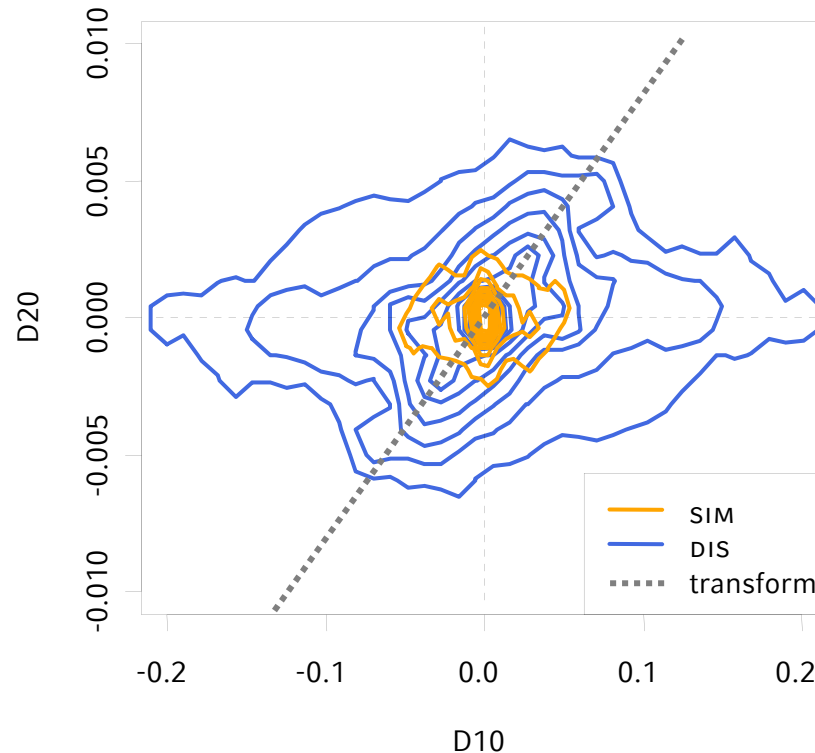
## EXPLOIT FEATURE CORRELATION

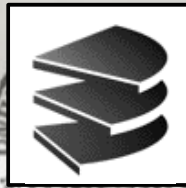


EXAMPLE:

COLOR HISTOGRAM DATA

FEATURE DIFFERENCE DISTRIBUTION FOR TWO DIMENSIONS





## MAXIMIZE VARIANCE DIFFERENCE

Transform to other space via  $W = (w_1, \dots, w_{d^*})$ :

$$\begin{aligned} \max \quad & w_i^T \cdot (\Sigma_{\text{DIS}} - \Sigma_{\text{SIM}}) \cdot w_i \\ \text{s.t.} \quad & w_i \perp w_j \quad \forall i, j \in \{1, \dots, d^*\} \end{aligned}$$

$$\Sigma_{\text{SIM}} = \sum_{x_s \in \text{SIM}} x_s^T \cdot x_s, \quad \Sigma_{\text{DIS}} = \sum_{x_d \in \text{DIS}} x_d^T \cdot x_d$$

This is equivalent to solving the EVD:

$$\lambda w = (\Sigma_{\text{DIS}} - \Sigma_{\text{SIM}}) \cdot w$$

TRAIN\_BED ( $X$  with  $x_i \in \mathbb{R}^d$ , SIM, DIS,  $d^*$ ):

1. Derive  $\Sigma_{\text{SIM}}$  and  $\Sigma_{\text{DIS}}$
2. Compute feature transformation  $W \in \mathbb{R}^{d, d^*}$
3. Get weights  $q_i$  ( $i \in \{1, \dots, d^*\}$ ) for new feature space  $W^T X$  using the features' variance differences

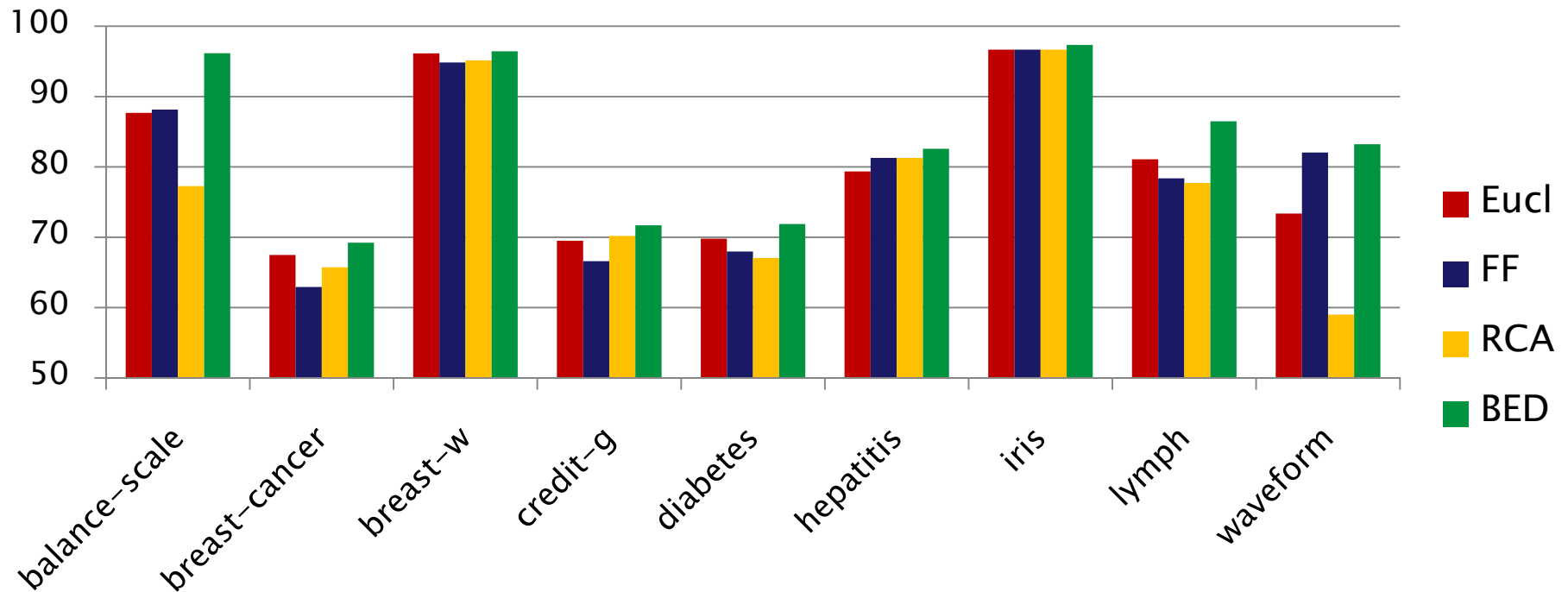
OUTPUT:

$$\text{BED}(x_1, x_2) = \left( \sum_{i=1}^{d^*} q_i \right)^{-1} \cdot \sum_{i=1}^{d^*} q_i \cdot \text{BE}_i(W^T x_1, W^T x_2)$$

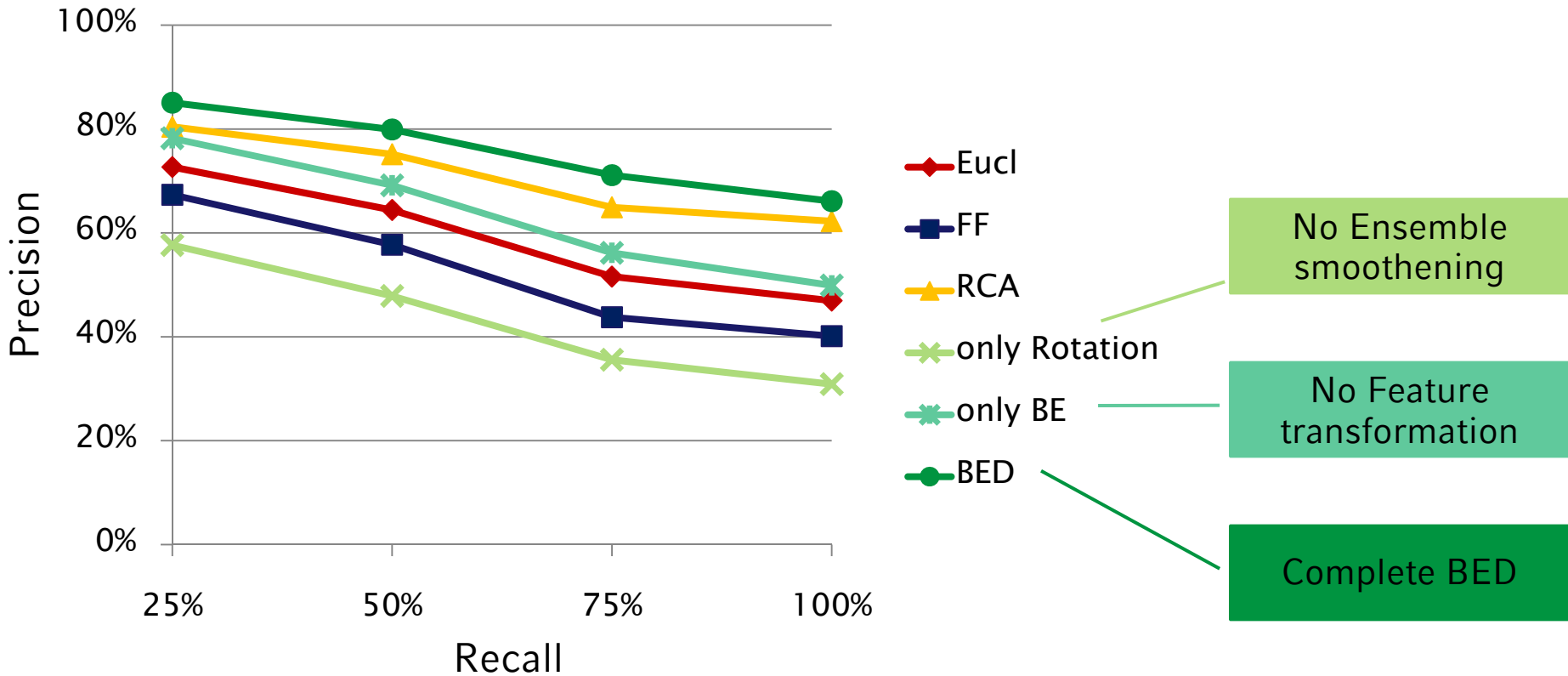


## UCI CLASSIFICATION DATASETS:

Accuracy Comparison to Euclidian Distance (Eucl), Fisher Faces (FF) and Relevant Component Analysis (RCA)



## IMAGE RETRIEVAL DATASET (PRECISION RECALL CURVES)



## BEDS NOW:

- Balanced, adaptive distance measure
- Easily interpretable
- Flexible w.r.t. data input

## IMPROVEMENTS:

- Runtime / Indexing ideas
- Local adaptivity
- Alternative weighting scheme



Thank you.