

SIMILARITY ESTIMATION:

- Applications for similarity or dissimilarity terms:
 - Retrieval Queries (Ranking, Range Queries)
 - Clustering
 - Model Training
- Similarity estimates reflect actual data similarity, i.e.: $s(x_1, x_2) > s(x_1, x_3) \Rightarrow x_1$ more similar to x_2 than to x_3
- Commonly used: *Distance Measures*
- For $x_1, x_2 \in \mathbb{R}^d$: L_p -norms

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

- Pitfalls:
 - Target similarity is ignored
 - Irrelevant features / Correlated features
 - Large influence of single dimensions
 - Arbitrarily large distances

OUR SOLUTION:

- Split the problem into similar (SIM) and dissimilar (DIS) object pairs
- Control the influence of a dimension i as probability in $[0,1]$ using a *Bayes Estimate* (BE):

$$BE_i(x_1, x_2) = \frac{p_{DIS} \cdot P((x_{1,i} - x_{2,i}) | DIS)}{p_{DIS} \cdot P((x_{1,i} - x_{2,i}) | DIS) + p_{SIM} \cdot P((x_{1,i} - x_{2,i}) | SIM)}$$

- Global distance = *Bayes Ensemble Distance* (BED):

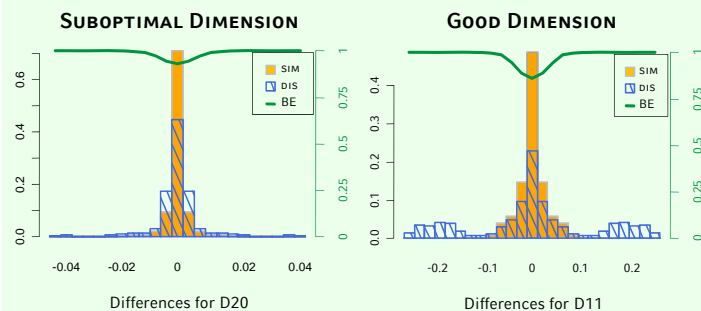
$$BED(x_1, x_2) = \frac{1}{d} \cdot \sum_{i=1}^d BE_i(x_1, x_2)$$

- More stable against outlier dimensions than a classical Naïve Bayes Classifier:

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

DIFFERENCE DISTRIBUTIONS:

- Similar (SIM) and dissimilar (DIS) object pairs can form differentiable difference distributions
- Example distributions for 32-dimensional color Histogram Data of 34 class image dataset



FEATURE QUALITY ASSESSMENT:

- Ensemble method: Introduce weights
- Relevance terms for variance difference of SIM and DIS

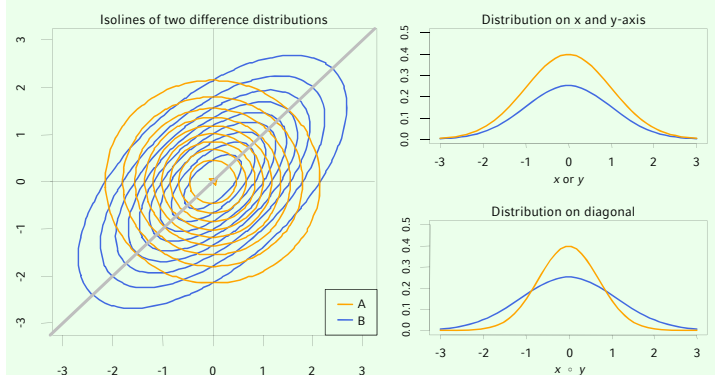
$$q_i = \sigma_{DIS}^2 - \sigma_{SIM}^2 = \text{avg}_{x_d \in DIS} (x_{d,i}^2) - \text{avg}_{x_s \in SIM} (x_{s,i}^2)$$

- Include into ensemble:

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

FEATURE SPACE IMPROVEMENT:

- Dimensionality reduction
- Exploitation of correlated features



- Distance covariances of SIM and DIS

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

- Target transformation $W = (w_1, \dots, w_{d^*})$ to dimension d^*
- Maximize the variance difference:

$$\max w_i^T \cdot (\Sigma_{DIS} - \Sigma_{SIM}) \cdot w_i$$

s.t. $w_i \perp w_j \quad \forall i, j \in \{1, \dots, d^*\}$

- Equivalent to solving EVD: $\lambda w = (\Sigma_{DIS} - \Sigma_{SIM}) \cdot w$

ALGORITHM:

INPUT: X with $x_i \in \mathbb{R}^d$, SIM, DIS, target dimension d^*

- (1) Derive Σ_{SIM} and Σ_{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: Bayes Ensemble Distance

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

CONCLUSIONS ON BEDs:

- Balanced, adaptive distance measure
- Easily interpretable
- Applicable to various datasets (discrete class labels, pair-wise similarity labels, regression target functions)