

Generalized Clustering via kernel embeddings

Stefanie Jegelka, Arthur Gretton, Bernhard Schölkopf,
Bharath K. Sriperumbudur, Ulrike von Luxburg



MAX-PLANCK-GESELLSCHAFT

Max Planck Institute
for Biological Cybernetics

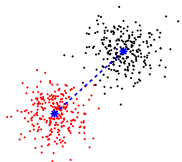
Tübingen, Germany



BIOLOGISCHE KYBERNETIK

Idea

Decompose sample into...

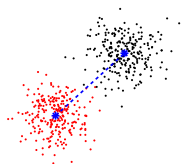


locally distinct clusters
separation by
first-order moments (means)

$$\text{Mixture Model: } P = \pi_1 P_1 + \pi_2 P_2$$

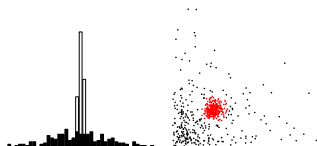
Idea

Decompose sample into...



locally distinct clusters
separation by
first-order moments (means)

generalize
→



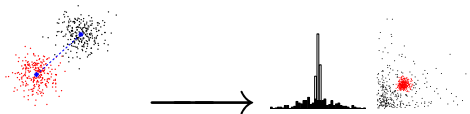
distinct distributions
separation by
higher-order moments

variance, kurtosis, ...

$$\text{Mixture Model: } P = \pi_1 P_1 + \pi_2 P_2$$

Decomposition

$$P = \pi_1 P_1 + \pi_2 P_2$$



first-order moments (means)

higher-order moments

$$\max_{P_1, P_2}$$

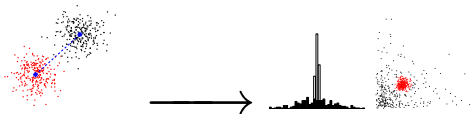
$$D(P_1, P_2) + \lambda \Omega(P_1, P_2)$$

max. discrepancy

favor "simplicity"

Decomposition

$$P = \pi_1 P_1 + \pi_2 P_2$$



first-order moments (means)

higher-order moments

\max_{P_1, P_2}

$D(P_1, P_2)$

+

$\lambda \Omega(P_1, P_2)$

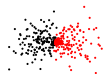
max. discrepancy

favor "simplicity"

Which Discrepancy Measure?

means

$$|\mathbb{E}_{x \sim P_1}[x] - \mathbb{E}_{y \sim P_2}[y]| = |\mu_1 - \mu_2|$$

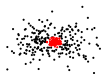
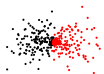


Which Discrepancy Measure?

means $|\mathbb{E}_{x \sim P_1}[x] - \mathbb{E}_{y \sim P_2}[y]| = |\mu_1 - \mu_2|$

variance $|\mathbb{E}_{x \sim P_1}[(x - \mu_1)^2] - \mathbb{E}_{y \sim P_2}[(y - \mu_2)^2]|$

d th moment $|\mathbb{E}_{x \sim P_1}[x^d] - \mathbb{E}_{y \sim P_2}[y^d]|$



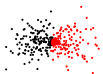
Which Discrepancy Measure?

means $|\mathbb{E}_{x \sim P_1}[x] - \mathbb{E}_{y \sim P_2}[y]| = |\mu_1 - \mu_2|$

variance $|\mathbb{E}_{x \sim P_1}[(x - \mu_1)^2] - \mathbb{E}_{y \sim P_2}[(y - \mu_2)^2]|$

d th moment $|\mathbb{E}_{x \sim P_1}[x^d] - \mathbb{E}_{y \sim P_2}[y^d]|$

general $|\mathbb{E}_{x \sim P_1}[g(x)] - \mathbb{E}_{y \sim P_2}[g(y)]|$



which "discriminative" g ?



Kernel Framework: Maximum Mean Discrepancy

$$P = \pi_1 P_1 + \pi_2 P_2$$

$$\max_{P_1, P_2} \underbrace{D(P_1, P_2)}_{\text{discrepancy}} + \lambda \underbrace{\Omega(P_1, P_2)}_{\text{"simplicity"}}$$

$$\text{MMD}(P_1, P_2) = \sup_{g \in \mathcal{H}, \|g\|_{\mathcal{H}} \leq 1} \left| \mathbb{E}_{x \sim P_1} g(x) - \mathbb{E}_{y \sim P_2} g(y) \right|$$

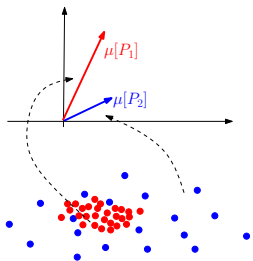
- implementation ?

Kernel Framework: Maximum Mean Discrepancy

$$P = \pi_1 P_1 + \pi_2 P_2$$

$$\max_{P_1, P_2} \underbrace{D(P_1, P_2)}_{\text{discrepancy}} + \lambda \underbrace{\Omega(P_1, P_2)}_{\text{"simplicity"}}$$

$$\text{MMD}(P_1, P_2) = \sup_{g \in \mathcal{H}, \|g\|_{\mathcal{H}} \leq 1} \left| \mathbb{E}_{x \sim P_1} g(x) - \mathbb{E}_{y \sim P_2} g(y) \right|$$



Embedding

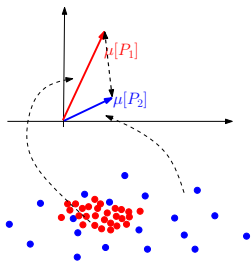
- represent each P_i by a mean function $\mu[P_i]$ in Hilbert space \mathcal{H}

Kernel Framework: Maximum Mean Discrepancy

$$P = \pi_1 P_1 + \pi_2 P_2$$

$$\max_{P_1, P_2} \underbrace{D(P_1, P_2)}_{\text{discrepancy}} + \lambda \underbrace{\Omega(P_1, P_2)}_{\text{"simplicity"}}$$

$$\begin{aligned} \text{MMD}(P_1, P_2) &= \sup_{g \in \mathcal{H}, \|g\|_{\mathcal{H}} \leq 1} \left| \mathbb{E}_{x \sim P_1} g(x) - \mathbb{E}_{y \sim P_2} g(y) \right| \\ &= \|\mu[P_1] - \mu[P_2]\|_{\mathcal{H}} \end{aligned}$$



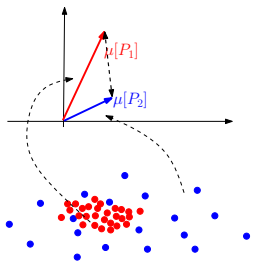
- represent each P_i by a mean function $\mu[P_i]$ in Hilbert space \mathcal{H}
- discrepancy steerable via kernel

Kernel Framework: Maximum Mean Discrepancy

$$P = \pi_1 P_1 + \pi_2 P_2$$

$$\max_{P_1, P_2} \underbrace{D(P_1, P_2)}_{\text{discrepancy}} + \lambda \underbrace{\Omega(P_1, P_2)}_{\text{"simplicity"}}$$

$$\begin{aligned} \pi_1 \pi_2 \text{MMD}(P_1, P_2)^2 &= \pi_1 \pi_2 \sup_{g \in \mathcal{H}, \|g\|_{\mathcal{H}} \leq 1} \left| \mathbb{E}_{x \sim P_1} g(x) - \mathbb{E}_{y \sim P_2} g(y) \right|^2 \\ &= \pi_1 \pi_2 \|\mu[P_1] - \mu[P_2]\|_{\mathcal{H}}^2 \end{aligned}$$

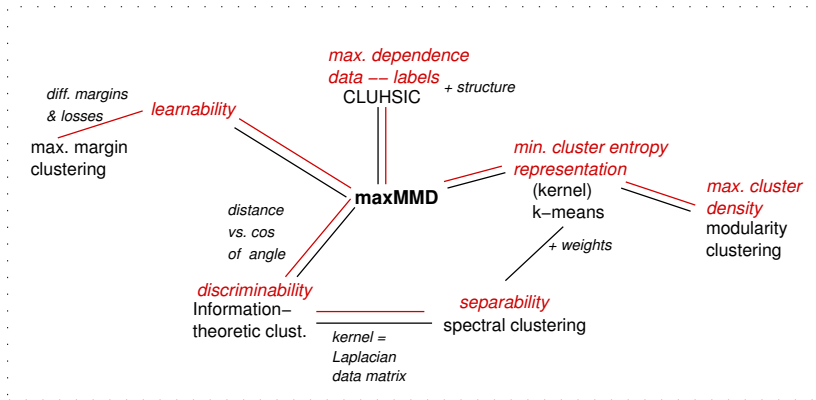


- represent each P_i by a mean function $\mu[P_i]$ in Hilbert space \mathcal{H}
- discrepancy steerable via kernel

Connections between clustering concepts – why?

- better understanding of methods – better understanding of results
- new algorithms for old objectives by transfer

Connections



Generalization of K-means

assignment of point x_i
to cluster j

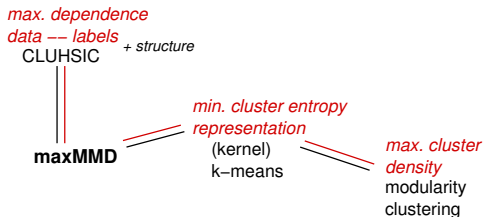
x_i mapped
into \mathcal{H}

mean representative of
"cluster" P_j

$$D(P_1, P_2) = - \sum_{i=1}^n \sum_{j=1,2} \alpha_{ij} \|\varphi(x_i) - \mu[P_j]\|^2 + \text{const}$$

max distance of **means** → max discrepancy of **moments**
 min **variance** → min **"entropy"**

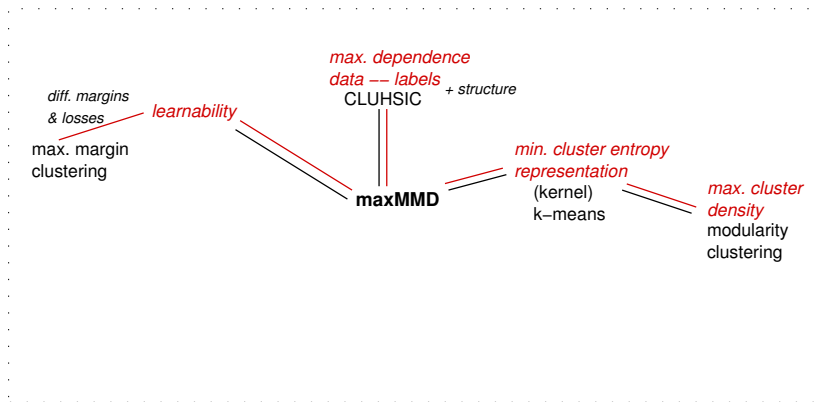
Connections



capture structure:

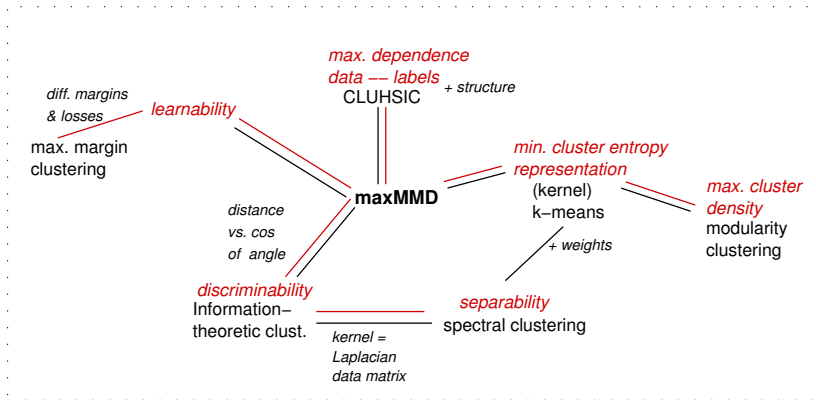
maximize statistical dependence between labels and data

Connections



capture structure:
minimize Bayes risk

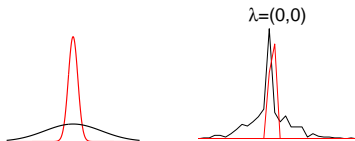
Connections



"Regularization"

$$P = \pi_1 P_1 + \pi_2 P_2$$

$$\max_{P_1, P_2} \underbrace{D(P_1, P_2)}_{\text{discrepancy}} + \lambda \underbrace{\Omega(P_1, P_2)}_{\text{"simplicity"}}$$



Implementation

- non-convex optimization problem:
 - without regularization: kernel k-means
 - with regularization: solver

- empirical performance: comparable to spectral clustering (Normalized cut), kernel k-means

Summary

- generalized clustering: by discrepancy of distributions
- flexible framework
- incorporates several clustering concepts