

Information k means and application to digital images

Gautier Appert

ENSAE Paris Tech.
joint work with Olivier Catoni

Representation of image with multi-level bags of labels

First step: creating a bag of labels

- Divide each image $\{X_1, \dots, X_n\}$ into non overlapping patches $\{B_i, i \in I\}$
 \implies cluster all patches $\{B_j, \text{for all } j\}$.

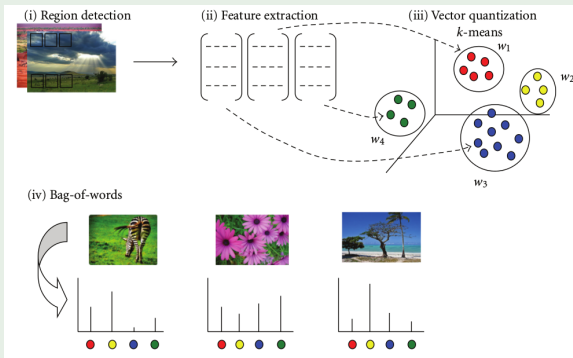


Figure : Figure 1 in Bag-of-Words Representation in Image Annotation: A Review. Chih-Fong Tsai

- Each image X is then represented as $\mathbb{P}_{W|X} = \frac{1}{m} \sum_{i=1}^m \delta_{w_i}$.

Representation of image with multi-level bags of labels

Multi-level

- Clustering labels with respect to a distortion measure is difficult \implies Instead use contextual modeling : cluster labels w and w' if they share the same context, but do not appear together. Define the context C of W in image X as $C = \mathbb{P}_{W|X} \circ f_{W,\Delta}^{-1}$, where $f_{W,\Delta}$ is the function that sends W to the outer state Δ . Note that $C \in \mathcal{M}_+^1(\mathcal{W} \cup \{\Delta\})$ is a random measure, and a function of the couple of random variables (X, W) .
- Cluster/agreggate words $\{w, w'\}$ **iff** $\mathbb{P}_{C|W=w} \simeq \mathbb{P}_{C|W=w'}$.
- Which means that we are looking for a classification function $\ell : \mathcal{W} \rightarrow \mathcal{Z}$ such that

$$\mathbb{P}_{C|W} = \mathbb{P}_{C|\ell(W)} \iff C \perp\!\!\!\perp W \mid \ell(W).$$

- When this is the case, $\mathbb{P}_{W|X}$ can be recovered from $\mathbb{P}_{\ell(W)|X}$ and $\mathbb{P}_{W|\ell(W)}$.

Euclidian k -means: theoretical and empirical loss

Let $X \in \mathbb{R}^d$ with $\mathbb{P}_X \left[\|X\|_2^2 \right] < \infty$ and let $\ell : \mathcal{X} \rightarrow \{1, \dots, k\}$ be the labelling function. The set \mathcal{X} can be \mathbb{R}^d or the index set $\{1, \dots, n\}$.

theoretical loss

$$\begin{aligned} & \inf_{\ell} \inf_{\mu_1, \dots, \mu_k} \mathbb{P}_X \left[\|X - \mu_{\ell(X)}\|_2^2 \right] \\ &= \inf_{\ell} \mathbb{P}_X \left[\|X - \mathbb{E}[X | \ell(X)]\|_2^2 \right] \\ &= \inf_{\mu_1, \dots, \mu_k} \mathbb{P}_X \left[\min_{j \leq k} \|X - \mu_j\|_2^2 \right] \end{aligned}$$

Empirical loss

$$\begin{aligned} & \inf_{\ell} \inf_{\mu_1, \dots, \mu_k} \frac{1}{n} \sum_{i=1}^n \|X_i - \mu_{\ell(i)}\|_2^2 \\ &= \inf_{\ell} \frac{1}{n} \sum_{j=1}^k \sum_{i \in \ell^{-1}(j)} \|X_i - \bar{\mu}_j\|^2 \\ &= \inf_{\mu_1, \dots, \mu_k} \frac{1}{n} \sum_{i=1}^n \min_{j \leq k} \|X_i - \mu_j\|_2^2 \end{aligned}$$

Lloyd's algorithm finds a local minimum through an iterative scheme: allocate data points to the nearest centroid and recompute centers from this partition.

Geometric mean of a conditional probability measure

Definition (Geometric mean of a conditional probability measure.)

Define the geometric mean function $\mathcal{G}(\cdot, \cdot)$ of a conditional probability measure $dP(t|s) = m(t|s) d\mu(t)$ with respect to the probability measure $dP(s)$ as

$$\begin{aligned}\mathcal{G}(dP(t|s), dP(s)) &= Z^{-1} \exp \left\{ \int \log [dP(t|s)] dP(s) \right\} \\ &\stackrel{\text{def}}{=} Z^{-1} \exp \left\{ \int \log [m(t|s)] dP(s) \right\} d\mu(t),\end{aligned}$$

where Z is a normalizing constant. Note that this is independent from the choice of μ : if $\nu \in \mathcal{M}_+^1$ is such that $\mu \ll \nu$, $dP(t|s) = \frac{d\mu}{d\nu}(t) m(t|s) d\nu(t)$ and

$$\mathcal{G}(dP(t|s), dP(s)) = Z^{-1} \exp \left\{ \int \log \left(\frac{d\mu}{d\nu}(t) m(t|s) \right) dP(s) \right\} d\nu(t).$$

Information k -means: theoretical and empirical loss

Let (Y, X) be a couple of random variables, assume that $\mathbb{P}_{Y|X}$ is known, whereas \mathbb{P}_X may be unknown.

Theoretical version

$$\begin{aligned} \inf_{\ell} \inf_{Q_{Y|\ell(X)}} \mathbb{P}_X \left[\mathcal{K}(Q_{Y|\ell(X)}, \mathbb{P}_{Y|X}) \right] &= \inf_{Q_{Y|j}, j \leq k} \mathbb{P}_X \left[\inf_{j \leq k} \mathcal{K}(Q_{Y|j}, \mathbb{P}_{Y|X}) \right] \\ &= \inf_{\ell} \mathbb{P}_X \left[\mathcal{K}(Q_{Y|\ell(X)}^*, \mathbb{P}_{Y|X}) \right] = \inf_{\ell} \mathbb{P}_X \left[\log(Z_{\ell(X)}^{-1}) \right], \end{aligned}$$

where $Q_{Y|\ell(X)}^*$ (the information k -means centers) and $Z_{\ell(X)}$ (the normalizing constants) are defined as

$$Q_{Y|\ell(X)}^* \stackrel{\text{def}}{=} \mathcal{G}(\mathbb{P}_{Y|X}, \mathbb{P}_{X|\ell(X)}) = Z_{\ell(X)}^{-1} \exp \left\{ \mathbb{P}_{X|\ell(X)} \left[\log \mathbb{P}_{Y|X} \right] \right\}.$$

Information k -means: theoretical and empirical loss

Consider a set of conditional probability distributions $R_{Y|i}$ for the random variable Y knowing $i \in [n]$, that we want to cluster.

Empirical loss

$$\begin{aligned} & \inf_{\ell: \{1, \dots, n\} \rightarrow \{1, \dots, k\}} \inf_{Q_{Y|\ell(i)}} \frac{1}{n} \sum_{i=1}^n \mathcal{K}(Q_{Y|\ell(i)}, R_{Y|i}) \\ &= \inf_{Q_{Y|j}} \frac{1}{n} \sum_{i=1}^n \inf_{j \in \{1, \dots, k\}} \mathcal{K}(Q_{Y|j}, R_{Y|i}) \\ &= \inf_{\ell} \frac{1}{n} \sum_{j=1}^k \inf_{Q_{Y|j}} \sum_{i \in \ell^{-1}(j)} \mathcal{K}(Q_{Y|j}, R_{Y|i}) \\ &= \inf_{\ell} \frac{1}{n} \sum_{i=1}^n \mathcal{K}(Q_{Y|\ell(i)}^*, R_{Y|i}) = \inf_{\ell} \sum_{j=1}^k \frac{|\ell^{-1}(j)|}{n} \log(Z_j^{-1}), \end{aligned}$$

$$\text{where } Q_{Y|j}^* = Z_j^{-1} \prod_{i \in \ell^{-1}(j)} R_{Y|i}^{1/|\ell^{-1}(j)|}.$$

Starting point

Due to the properties of the Kullback divergence, the following algorithm to compute an initial classification ℓ gives promising results.

- Start from $k = 1$ and $\ell^{-1}(1) = \{1, \dots, n\}$.
- Switch from k to $k + 1$ by removing iteratively from $\ell^{-1}(k)$ $\arg \max_{i \in \ell^{-1}(k)} \mathcal{K}(Q_{Y|i}^*, R_{Y|i})$ to put it in $\ell^{-1}(k + 1)$, until $\log(Z_k^{-1}) \leq \eta$.
- Continue if $\log(Z_{k+1}^{-1}) > \eta$.

Link with Information projection

Definition (Information projection.)

Let P be a probability distribution, and let \mathcal{Q} be set of probability distribution. The information projection or I-projection of P onto \mathcal{Q} is defined as

$$Q^* \in \arg \min_{Q \in \mathcal{Q}} \mathcal{K}(Q, P).$$

Information k -means seen as an information projection

- Consider the model

$$\begin{aligned} \mathcal{Q} &= \left\{ Q_{Y,X} : Q_X = P_X, Q_{Y|X} = Q_{Y|\ell(X)}, \ell(X) \in \{1, \dots, k\} \right\}, \\ \inf_{Q_{Y,X} \in \mathcal{Q}} \mathcal{K}(Q_{Y,X}, P_{Y,X}) \\ &= \inf_{Q_{Y,X} \in \mathcal{Q}} Q_X \left[\mathcal{K}(Q_{Y|X}, P_{Y|X}) \right] + \mathcal{K}(Q_X, P_X) \\ &= \inf_{\ell, Q_{Y|\ell(X)}} P_X \left[\mathcal{K}(Q_{Y|\ell(X)}, P_{Y|X}) \right]. \end{aligned}$$

Information k -means become Euclidian with Gaussian distribution

Information k -means generalize Euclidian k -means

- Take $\mathbb{P}_{Y|X} = \mathcal{N}_p(X, \Sigma)$.
- One obtains

$$\begin{aligned} dQ_{Y|\ell(X)}^*(y) &\propto \exp \left\{ \mathbb{P}_{X|\ell(X)} \left[\log \left(\frac{d\mathbb{P}_{Y|X}}{d\lambda}(y) \right) \right] \right\} d\lambda(y) \\ &\propto \exp \left\{ -\frac{1}{2} (y^\top \Sigma^{-1} y - 2y^\top \Sigma^{-1} \mathbb{E}[X|\ell(X)]) \right\} \\ &\propto \mathcal{N}_p(\mathbb{E}[X|\ell(X)], \Sigma) \end{aligned}$$

- Then $\mathcal{K}(Q_{Y|\ell(X)}^*, \mathbb{P}_{Y|X}) = \|X - \mathbb{E}[X|\ell(X)]\|_{\Sigma^{-1}}^2$.
- $\inf_{\ell} \mathbb{P}_X \left[\mathcal{K}(Q_{Y|\ell(X)}^*, \mathbb{P}_{Y|X}) \right] = \inf_{\ell} \mathbb{P}_X \left[\|X - \mathbb{E}[X|\ell(X)]\|_{\Sigma^{-1}}^2 \right]$

PAC-Bayesian Margin bounds on information k -means

Information k -means loss in the case of discrete \mathbb{P}_Y

Let $Y \in \mathcal{Y}$ with $|\mathcal{Y}| < \infty$.

- $\mathcal{L}(Q) = \mathbb{P}_X \left[\min_{i \leq k} \mathcal{K}(Q_{Y|i}, \mathbb{P}_{Y|X}) \right]$
- Put $q_i = \frac{dQ_{Y|i}}{d\nu}$, $p_X = \frac{d\mathbb{P}_{Y|X}}{d\nu} \implies \mathcal{L}(q) = \mathbb{P}_X \left[\min_{i \leq k} \mathcal{K}(q_i, p_X) \right]$.
- Recall $\mathcal{K}(q_i, p_X) = \langle q_i, \log(q_i) - \log(p_X) \rangle$.
- Put $\theta_i = (-q_i, \langle q_i, \log(q_i) \rangle)^\top$,
 $\theta_{i,j} = (q_i - q_j, \langle q_i, \log(q_i) \rangle - \langle q_j, \log(q_j) \rangle)^\top \in \mathbb{R}^{|\mathcal{Y}|+1}$ and
 $W = (\log p_X, 1)^\top \in \mathbb{R}^{|\mathcal{Y}|+1}$.
- Hence, $\mathcal{K}(q_i, p_X) = \langle \theta_i, W \rangle$ and
 $\mathcal{K}(q_i, p_X) < \mathcal{K}(q_j, p_X) \iff \langle \theta_{i,j}, W \rangle \geq 0$.

PAC-Bayesian Margin bounds on information k -means

upper bound on the Loss

- Using the fact that

$$\min_{i \leq k} a_i = \sum_{i=1}^k a_i \prod_{j=1}^{i-1} \mathbb{1}(a_i < a_j) \prod_{j=i+1}^k \mathbb{1}(a_i \leq a_j).$$

- We can rewrite

$$\mathcal{L}(q) \leq \sum_{i=1}^k \mathbb{P}_X \left[\langle \theta_i, W \rangle \prod_{j \neq i} \mathbb{1}(\langle \theta_{i,j}, W \rangle \geq 0) \right].$$

PAC-Bayesian Margin bounds on information k -means

Put a perturbation and a margin

- Gaussian perturbation $\rho_\theta = \mathcal{N}\{\theta, \beta^{-1} \mathbf{I}_{|y|+1}\}$.
- Margin $M = \gamma \|W\|$.

lemma

$$\mathcal{L}(q) \leq \sum_{i=1}^k \Phi\left(\gamma\sqrt{\beta}\right)^{-(k-1)} \\ \times \mathbb{P}_X \left[\langle \theta_i, W \rangle \prod_{j \neq i} \int \mathbb{1}(\langle \theta'_{i,j}, W \rangle + \gamma \|W\| \geq 0) \, d\rho_{\theta_{i,j}}(\theta'_{i,j}) \right]$$

- Looks like some kind of classification problem with margin $M = \gamma \|W\|$. [Catoni, Lecture notes, 2014].
- Estimation of the mean of $\langle \theta_i, W \rangle$. [Catoni, Giulini 2017].

PAC-Bayesian Margin bounds on information k -means

Upper bound of the information k -means loss

Introduce $g_1(t) = \frac{1}{t}(\exp(t) - 1)$ and $g_2 = \frac{1}{t^2}(\exp(t) - 1 - t)$. With probability at least $1 - \varepsilon$,

$$\begin{aligned}\mathcal{L}(q) \leq & \Phi\left(\gamma\sqrt{\beta}\right)^{-(k-1)} \left\{ \sum_{i=1}^k \hat{\mathbb{P}}_X^n \left[\langle \theta_i, Z \rangle \bar{H}(W, \theta_{-i}) \right] \right. \\ & + \frac{\lambda a}{2} \mathbb{P}_X \left(\langle \theta_i, W \rangle^2 \bar{H} \right) + \frac{\lambda b}{2\beta} \mathbb{P}_X \left(\|W\|^2 \bar{H} \right) \\ & \left. + \frac{\alpha^p}{p+1} \mathbb{P}_X \left(|\langle \theta_i, W \rangle| \|W\|^p \right) + \frac{k\beta}{2n\lambda} \sum_{i=1}^k \left\{ \|\theta_i\|^2 + \sum_{j \neq i} \|\theta_{ij}\|^2 \right\} \frac{k \log(\varepsilon^{-1})}{n\lambda} \right\},\end{aligned}$$

$$\text{where } a = g_2 \left(\frac{\lambda \|\theta_i\|}{\alpha} \right), \quad b = g_1 \left(\frac{\lambda^2}{2\beta\alpha^2} \right) \exp \left(\frac{\lambda \|\theta_i\|}{\alpha} \right),$$

$$\bar{H} = \prod_{j \neq i} \Phi \left(\sqrt{\beta} (\gamma + \|W\|^{-1} \langle \theta_{ij}, W \rangle) \right) \text{ and } Z = \frac{\min(\lambda \|W\|, 1) W}{\lambda \|W\|}.$$

PAC-Bayesian Margin bounds on Euclidian k -means

Euclidian k – means loss

- Consider $\mathcal{L}(\mu) = \mathbb{P}_X \left(\min_{i \leq k} \|X - \mu_i\|^2 \right)$

Change of notation

- Put $\theta_i = (-\mu_i, \|\mu_i\|^2)^\top$, $\theta_{i,j} = (\mu_i - \mu_j, \|\mu_j\|^2 - \|\mu_i\|^2)^\top \in \mathbb{R}^{p+1}$ and $W = (2X, 1)^\top \in \mathbb{R}^{p+1}$.
- Hence (polarization identity), $\|X - \mu_i\|^2 = \langle \theta_i, W \rangle + \|X\|^2$ and $\|X - \mu_i\|^2 < \|X - \mu_j\|^2 \iff \langle \theta_{i,j}, W \rangle \geq 0$.

$$\mathcal{L}(\mu) \leq \sum_{i=1}^k \mathbb{P}_X \left[\langle \theta_i, W \rangle \prod_{j \neq i} \mathbb{1}(\langle \theta_{i,j}, W \rangle \geq 0) \right] + \mathbb{P}_X \left[\|X\|^2 \right].$$

- Same bound as before for $\mathcal{L}(\mu) - \mathbb{P}_X \left[\|X\|^2 \right]$.

Application to images: a small example

How to create patches with a random support ?

- Let $(X_i, i \in I) \in \mathbb{R}^I$ be a random image, where $|I| < \infty$ is the number of pixels.
- Represent the pixel location i by a random variable S , putting

$$\mathbb{P}_{S,V|X} = \frac{1}{|I|} \sum_{i \in I} \delta_i \otimes \delta_{X_i}, \quad \text{where } (S, V) \in I \times \mathbb{R}.$$

- Add noise, introducing $V' = V + \xi$ such that $\mathbb{E}(V'|X, S) = V$.
- Put $U = (S, V')$ and change the representation of X to $\mathbb{P}_{U|X}$.
- Use an auxiliary set of images represented by the distribution $Q_{\theta,U} \in \mathcal{M}_+^1[\Theta \times (I \times \mathbb{R})]$, where $|\text{supp}(Q_{\theta,U})| < \infty$. Take for instance the empirical distribution of n independent copies of X , or

more precisely $Q_{\theta,U} = \frac{1}{n} \sum_{j=1}^n \delta_j \otimes \mathbb{P}_{S,V|X=X_j}$, where

$$(X_j, 1 \leq j \leq n) \sim \mathbb{P}_X^{\otimes n}.$$

- Solve
$$\inf_{\ell_\theta: \text{supp}(Q_{U|\theta}) \rightarrow \{1, \dots, k\}} \inf_{Q_{X|\theta, \ell_\theta(U)}} Q_{\theta,U} \left[\mathcal{K}(Q_{X|\theta, \ell_\theta(U)}, \mathbb{P}_{X|U}) \right].$$

Application to images: a small example

- Define the patch process as $\mathbb{P}_{T|X} = Q_{\theta, \ell_{\theta}(U)|X}$.
- In the exact case where $\mathbb{P}_{X|U} = Q_{X|\theta, \ell_{\theta}(U)}$,

$$\frac{d\mathbb{P}_{U|X, U \in \text{supp}(Q_U)}}{d\mathbb{P}_{U|U \in \text{supp}(Q_U)}}(u) = Z_X^{-1} Q_{\theta|U=u} \left[\frac{dQ_{\theta, \ell_{\theta}(U)|X}}{dQ_{\theta, \ell(U)}}(\theta, \ell_{\theta}(u)) \right].$$

- Cluster the patches solving

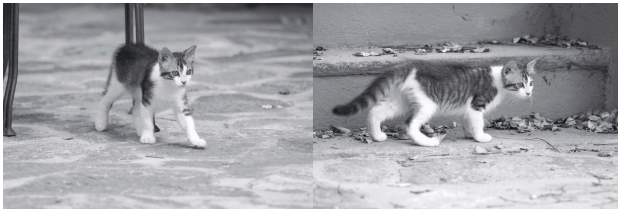
$$\inf_{\ell: \text{supp}(R_T) \rightarrow \{1, \dots, k\}} \inf_{R_{X|\ell(T)}} R_T [\mathcal{K}(R_{X|\ell(T)}, \mathbb{P}_{X|T})]$$

- Define a new representation as $\mathbb{P}_{W|X} = R_{\ell(T)|X}$.
- In the exact case where $R_{X|\ell(T)} = \mathbb{P}_{X|T}$ and $\text{supp}(R_T) = \text{supp}(\mathbb{P}_T)$,

$$\frac{d\mathbb{P}_{T|X}}{d\mathbb{P}_T}(t) = Z_X^{-1} \frac{dR_{\ell(T)|X}}{dR_{\ell(T)}}(\ell(t)),$$

showing that the previous representation $\mathbb{P}_{T|X}$ can be recovered exactly from the next one $\mathbb{P}_{W|X} = R_{\ell(T)|X}$ and the marginal distributions \mathbb{P}_T and $R_{\ell(T)}$.

Application to images: a small example



Application to images: a small example



Figure : Extracted patches : 500×500 from two images 1000×1500 . .

the training sample corresponds to the extracted patches.

Application on images: small example

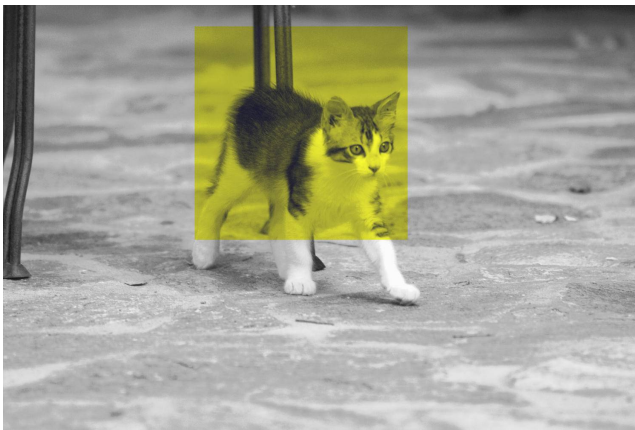


Figure : Selected image.

Application on images: small example



Figure : Clustering with information k -means.

Application on images: small example



Figure : Clustering with information k -means.

Application on images: small example



Figure : Clustering with information k -means.

Application on images: small example

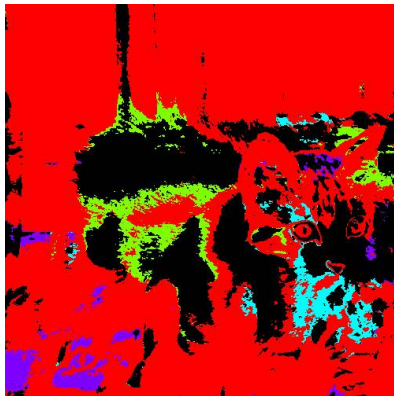


Figure : Clustering with information k -means.

Application on images: small example

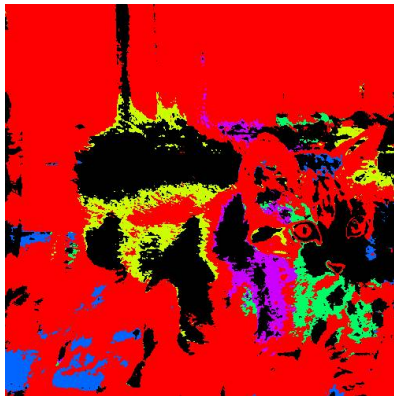


Figure : Clustering with information k -means.

Application on images: small example

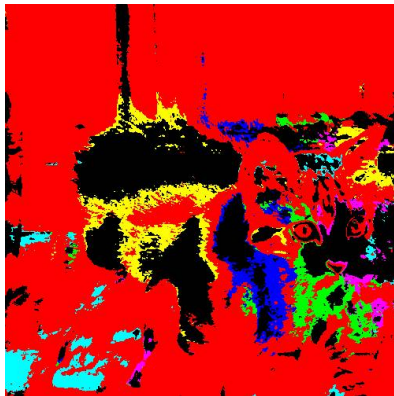


Figure : Clustering with information k -means.

Application on images: small example

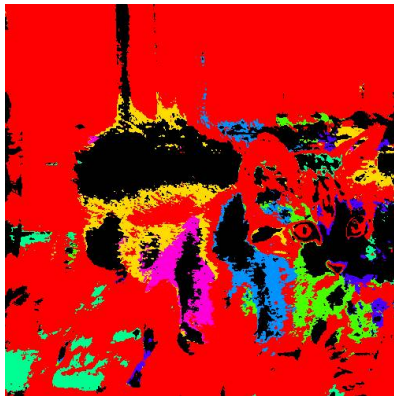


Figure : Clustering with information k -means.

Application on images: small example



Figure : Clustering with information k -means.