

Proposed - Energy Stats vs. Clustering.

Wednesday, March 8, 2017

8:49 AM

①

k-means objective function:

$$E = \frac{1}{2} \sum_{k=1}^K \sum_{x_i \in C_k} \|x_i - m_k\|^2$$

$$= \frac{1}{2} \sum_{k=1}^K \sum_{i=1}^n z_{ki} \|x_i - m_k\|^2$$

$$z_{ki} = \begin{cases} 1 & \text{if } x_i \in C_k \\ 0 & \text{otherwise} \end{cases}$$

Notice that $\sum_{k=1}^K z_{ki} = 1$, $\sum_{i=1}^n z_{ki} = n_k$. Introducing

$$X = (x_1 | x_2 | \dots | x_n)_{D \times N} \quad M = (m_1 | m_2 | \dots | m_K)_{D \times K}$$

$$Z = z_{ij} \in \mathbb{R}_{K \times N}$$

Then we have

$$E = \|X - MZ\|_F^2$$

We can minimize this over M yielding:

$$\min_Z \|X - XZ^\top (ZZ^\top)^{-1} Z\|_F^2$$

s.t. $z_{ij} \in \{0, 1\}$, $\sum_k z_{ki} = 1$

Define $\tilde{z} = z^T(z z^T)^{-1} z$ then (2)

$$\begin{aligned}\|x - x \tilde{z}\|_F^2 &= \text{Tr}((x - x \tilde{z})^T(x - x \tilde{z})) \\ &= \text{Tr}(x^T x - x^T x \tilde{z} - \tilde{z}^T x^T x + \tilde{z}^T x^T x \tilde{z}) \\ &= \text{Tr}(x^T x - 2x^T x \tilde{z} + \tilde{z}^T x^T x \tilde{z})\end{aligned}$$

Grom matrix: $K \equiv x^T x$.

$$\text{Tr}(K - 2K \tilde{z} + \tilde{z}^T K \tilde{z})$$

$$\begin{aligned}\|x - Mz\|^2 &\\ &= \text{Tr}((I - M)^T K (I - M))\end{aligned}$$

So k-means can be written as an Integer Quadratic Optimization Problem:

$$\min \text{Tr}(K - 2K \tilde{z} + \tilde{z}^T K \tilde{z})$$

$$\text{s.t. } z_i \in \{0, 1\}, \sum_n z_{ni} = 1, \tilde{z} = z^T(z z^T)^{-1} z.$$

Another formulation is

$$K = x^T x$$

$$\boxed{\begin{aligned}\max \text{Tr}(G^T K G) \\ \text{s.t. } G \geq 0, G^T G = I, G G^T J = J\end{aligned}}$$



$$\text{where } G_{ik} = \begin{cases} \frac{1}{\sqrt{n}} & \text{if } x_i \in C_k \\ 0 & \text{otherwise} \end{cases}$$

We can still manipulate and write this in different forms. In any case, the problem should be expressed in terms of K , the grom matrix.

We can also write $y \in \mathbb{R}^T (\mathbb{R} \mathbb{R}^T)^{-1}$. Then $M = Xy$.
 Thus y is a doubly stochastic matrix, thus $y \geq 0$.

$$\left\{ \begin{array}{l} \min_{Y, Z} \|X - XYZ\|_F^2 \\ \text{s.t. } Y \geq 0 \quad \vec{1}^T y_i = 1 \quad H(y_i) \gg 0 \\ \quad Z \geq 0 \quad \vec{1}^T z_i = 1 \quad H(z_i) = 0 \end{array} \right. \xrightarrow{\Delta} \text{Tr}(I - YZ)^T K(I - YZ) \quad \begin{array}{l} \downarrow \\ X^T X \\ \text{gram matrix} \end{array}$$

\underbrace{\hspace{10em}}_{\text{entropy constraints}}

In archetypal analysis we have

$$\left\{ \begin{array}{l} \min_{Y, Z} \|X - XYZ\|_F^2 \\ \text{s.t. } Y \geq 0 \quad \vec{1}^T y_i = 1 \\ \quad Z \geq 0 \quad \vec{1}^T z_i = 1 \end{array} \right.$$

Non-negative matrix factorization:

$$\left\{ \begin{array}{l} \min_{Y, Z} \|X - XYZ\|_F^2 \\ \text{s.t. } Y \geq 0 \\ \quad Z \geq 0 \end{array} \right.$$

There are few options to solve these optimization problems.

(4)

1. Frank-Wolfe Algorithm.

- (i) Initialize y, z
- (ii) Fix z , solve for y
- (iii) Fix y , solve for z
- (iv) Repeat until converge

This will require the gradient of the objective function, which will depend on the Gram matrix K :

$$\nabla(\dots) \sim Ky(z z^T) - K z^T - \dots$$

2. SDP Relaxation

Standard Primal Dual SDP

$$\begin{array}{ll} \min Q X & \rightarrow \\ \text{s.t. } x_{ii} = 1 & \max \Lambda \\ x \geq 0 & \text{s.t. } Q \geq \Lambda \\ & \Lambda \text{ diagonal} \end{array}$$

This involves an augmented Lagrangian. It is also possible to impose a rank constraint in the problem, and end up with an SDP with rank constraint. One can relax the constraint, solve it, then enforce the rank on a "rounding scheme" in the final solution.

(5)

3. Low Rank SDP.

We can employ the same technique used in that paper which provides the low-rank approx already in the optimization prob. itself. They did that for weighted k-means.

4. Alternating Direct Methods.

We can check how to solve these kind of optimization problems using alternating methods with augmented Lagrangians. A very nice opt's would be ADMM. This would be appealing and very novel. ADMM does not need gradient, & it would be very robust also.

How Energy Statistics Enters the Story?

In the above proposed formulation, all we need to do is exploit the kernel trick on k !

- From Sejdinovic et al, Ann. Stats (2013):

let $\rho: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a semi-metric on \mathcal{X} . Then $K(x, x') = \frac{1}{2}(\rho(x, x_0) + \rho(x', x_0) - \rho(x, x'))$ is positive definite iff ρ is of negative type.

(6)

Now the energy distance between two CDF's,

$$E_\alpha(F, G) = 2 \left[E \|x - y\|^\alpha - E \|x - x'\|^\alpha - E \|y - y'\|^\alpha \right]$$

where $x, x' \stackrel{iid}{\sim} F$, $y, y' \stackrel{iid}{\sim} G$ is of negative type.

We can embed F, G into a RKHS with kernel function given by the energy distance above.

They proved an equivalence between energy distance and MMD in RKHS. Thus we can replace the gram matrix K in the optimization problem to a kernel defined in terms of the energy distance!

If $x_i \sim \mu_i$, $x_j \sim \mu_j$

$$\|x_i - x_j\| = D(\mu_i, \mu_j)$$

To compute this from data
"empirical meas".

(7)

Then replace $k_{ij} \rightarrow k(x_i, x_j)$.

- Understand this better, on how to incorporate energy distance as a kernel.

Sketch of the problem:

$$\begin{aligned} & \min_{\mathbf{w}} \mathbf{Tr}(\mathbf{w}^T \mathbf{K} \mathbf{w}) \\ \text{s.t. } & \dots \end{aligned}$$

\downarrow

$$\text{replace } k = \mathcal{E}_\alpha(x_i, x_j)$$

$$\begin{aligned} & \text{SDP, with kernel} \\ & \text{energy dist.} \end{aligned}$$

→ can leave this

← use low rank SDP

← augmented Lagrange techniques

← Non-linear prog. stuff.

← ADMM.

we'll explore
role of \mathcal{L} through
most validate, and
Comparing
with ℓ_1 -
methods.

Clustering algorithm
based on energy dist. $\mathcal{E}(\cdot, \cdot)$ does
not assume
F, G.
It's gonna be non-parametric!