

## PROBLEM STATEMENT

$k$ -means is a popular clustering method. +Crucial link w/ exponential families:

$objects\ clustered = expectation\ parameters;$

$distortion\ used = square\ Euclidean\ distance = KL\ divergence\ between\ Gaussians\ w/\ identity\ covariance;$

This can be generalized to exponential families, where the distortion becomes a general Bregman divergence. This provides a wide range of data generating distributions.

But some roadblocks remain in all cases, like the lack of robustness to outliers of population minimizers (e.g. the right population minimizer is always the arithmetic average, non robust). Also, the current generalizations of exponential families ( $q$ -exponential families or even deformed exponential families) fail at removing roadblocks.

In this work, we provide a new attempt at getting the same complete framework and gain additional (controllable) robustness for population minimizers, via a new generalization of exponential families to **non-normalized** measures.

## CONTRIBUTIONS

1. Generalizing exponential family to tempered exponential family while maintaining *divisive* normalizations
2. Information-theoretic distortions between distributions (tempered KL div.)
3. Parameter-based clustering distortions (generalizing Bregman divergences)
4. The information theoretic/geometric link between the distortions
5. Robustness of the population minimizers

## TEMPERED EXPONENTIAL MEASURES AND CO-DENSITIES

From left to right: tempered logarithm, tempered exponential, perspective transform of tempered exponential ( $t^* \doteq 1/(2-t)$ ,  $[\cdot]_+ \doteq \max\{0, \cdot\}$ , same can be defined for tempered log)

$$\log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1), \quad \exp_t(z) \doteq [1 + (1-t)z]_+^{1/(1-t)} \quad (\exp_t)^*(z) \doteq t^* \exp_{t^*} \left( \frac{z}{t^*} \right).$$

**Theorem (TEMs & their co-densities)** — let

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \mid \begin{array}{l} \mathbb{E}_{\tilde{p}}[\varphi] \doteq \int \varphi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right\}, \quad \begin{array}{l} \varphi : \mathcal{X} \rightarrow \mathbb{R}^d \text{ sufficient statistics,} \\ \mathbf{h} \in \mathbb{R}^d \text{ expectation parameter,} \\ \xi \text{ base measure.} \end{array}$$

and (Capital  $\tilde{P}$  the measure of density  $\tilde{p}$  wrt  $\xi$ )

$$H_t(\tilde{P}) \doteq - \int \psi_t(\tilde{p}(\mathbf{x})) d\xi, \quad \text{where } \psi_t(z) \doteq z \log_t z - \log_{t-1} z$$

a (generalized notion of) Tsallis entropy. Then for any  $t \in [0, 1]$  and  $\mathbf{h} \in \mathbb{R}^d$ , the solution  $\arg \max_{\tilde{\mathcal{P}}_{t|\mathbf{h}}} H_t$  has the **non-normalized** density

$$\tilde{p}_{t|\theta}(\mathbf{x}) = \frac{\exp_t(\theta^\top \varphi(\mathbf{x}))}{\exp_t(G_t(\theta))} \quad \text{where } G_t(\theta) = (\log_t)^* \int (\exp_t)^*(\theta^\top \varphi(\mathbf{x})) d\xi$$

$G_t =$  (convex) cumulant ensuring normalization of the **dual**  $\tilde{p}^{1/t^*}$  (**co-density**);  $\theta =$  natural parameter (correspondence  $\theta = \nabla G_t^{-1}(\mathbf{h})$  like exponential families).

Compared to other generalizations of exponential families (deformed &  $q$ -exp. families), key parameters in closed form, e.g. cumulant ( $G_t$ ), **total mass** of TEM ( $\int \tilde{p} d\xi$ ), Cf paper.

**Theorem (information theoretic / geometric link)** — let

$$F_t(\tilde{P}_{t|\hat{\theta}} \| \tilde{P}_{t|\theta}) \doteq \int f \left( \frac{d\tilde{p}_{t|\hat{\theta}}}{d\xi} \oslash_t \frac{d\tilde{p}_{t|\theta}}{d\xi} \right) \cdot d\tilde{p}_{t|\theta}, \quad \begin{array}{l} f \doteq -\log_t, \\ x \oslash_t y \doteq (x^{1-t} - y^{1-t} + 1)_+^{\frac{1}{1-t}} \text{ if } x, y \geq 0 \\ \text{(else undefined)} \end{array}$$

Then  $\forall \tilde{P}_{t|\hat{\theta}}, \tilde{P}_{t|\theta}$  of the **same TEM family**,

$$F_t(\tilde{P}_{t|\hat{\theta}} \| \tilde{P}_{t|\theta}) = B_{G_t}(\hat{\theta} \| \theta), \quad \text{where } B_{G_t}(\hat{\theta} \| \theta) \doteq \frac{\overbrace{G_t(\hat{\theta}) - G_t(\theta) - (\hat{\theta} - \theta)^\top \nabla G_t(\theta)}^{\text{Bregman divergence!}}}{1 + (1-t)G_t(\hat{\theta})}.$$

TEM (example)	Support	$\lambda$	$\theta$	$G_t(\theta)$	$B_{G_t}(\hat{\theta} \  \theta)$
$t$ -exponential	$[0, \frac{3-2t}{(1-t)\lambda}]$	$\lambda$	$-\frac{\lambda}{3-2t}$	$-\log_{2-t} \left( (-\theta)^{\frac{1}{2-t}} \right)$	$t^* \cdot \left( \left( \frac{\hat{\theta}}{\theta} \right)^{2-t^*} - (2-t^*) \cdot \log_{t^*} \left( \frac{\hat{\theta}}{\theta} \right) - 1 \right)$

which Bregman divergence for  $t \rightarrow 1$ ? ☺

## POPULATION MINIMIZERS

$\{\theta_i\}_{i=1}^m =$  training set of parameters endowed with an implicit distribution. Two losses for the so-called left and right population minimizers:

$$L_l(\theta) \doteq \mathbb{E}_i[B_{G_t}(\theta \| \theta_i)]; \quad L_r(\theta) \doteq \mathbb{E}_i[B_{G_t}(\theta_i \| \theta)], \quad \mathbb{E}_i[\cdot] = \text{average over the training sample.}$$

The left and right population minimizers, respectively  $\theta_l$  and  $\theta_r$ , are defined as the arguments minimizing the corresponding loss. **Robustness** to outliers postulates that adding a new data point  $\theta_*$  with weight  $\varepsilon$  shifts  $\theta_{l/r}^{\text{new}} - \theta_{l/r}^{\text{old}} = \varepsilon \cdot z(\theta_*)$  with  $z(\cdot)$  (influence function) of *bounded norm*.

**Theorem (population minimizers)** — The left and right population minimizers satisfy respectively:

$$\theta_l = \nabla G_t^{-1}(\alpha_* \cdot \mathbb{E}_i \nabla G_t(\theta_i)), \quad \text{for some } \alpha_* > 0 \text{ (closed form available in some cases, Cf paper)} \quad ; \quad \theta_r = \mathbb{E}_i \left[ \frac{1}{\exp_t^{1-t}(G_t(\theta_i))} \cdot \theta_i \right]$$

Furthermore, for  $t \in [0, 1]$ ,  $\theta_l$  is robust iff robust for  $t = 1$  and  $\theta_r$  is robust if  $G_t(\theta) = \Omega(\|\theta\|)$  and  $t \neq 1$ .

## EXPERIMENTS (SEE PAPER FOR NUMERICAL RESULTS)

