



Clustering above Exponential Families with Tempered Exponential Measures

Ehsan Amid, Richard Nock and Manfred K. Warmuth

Exponential families & k -means clustering

Why this work: k -means = popular clustering, partition of space $\mathcal{X} \subseteq \mathbb{R}^d$ by finding set of centers $\mathcal{C} \doteq \{\mathbf{c}_j\}_{j \in [k]}$ minimizing loss to m -sample, $\mathbb{E}_{i \sim [m]} [\min_{j \in [k]} D(\boldsymbol{\theta}_i \| \mathbf{c}_j)]$

Exponential families & k -means clustering

Why this work: k -means = popular clustering, partition of space $\mathcal{X} \subseteq \mathbb{R}^d$ by finding set of centers $\mathcal{C} \doteq \{\mathbf{c}_j\}_{j \in [k]}$ minimizing loss to m -sample, $\mathbb{E}_{i \sim [m]} [\min_{j \in [k]} D(\boldsymbol{\theta}_i \| \mathbf{c}_j)]$

If D a Bregman divergence (Mahalanobis, Itakura-Saito, KL, etc.) then population minimizers trivial to compute **and** we equivalently minimise an information-theoretic loss between *distributions* in exponential families \Rightarrow embeds k -means in broad data generating processes

Exponential families & k -means clustering

Why this work: k -means = popular clustering, partition of space $\mathcal{X} \subseteq \mathbb{R}^d$ by finding set of centers $\mathcal{C} \doteq \{\mathbf{c}_j\}_{j \in [k]}$ minimizing loss to m -sample, $\mathbb{E}_{i \sim [m]} [\min_{j \in [k]} D(\boldsymbol{\theta}_i \| \mathbf{c}_j)]$

If D a Bregman divergence (Mahalanobis, Itakura-Saito, KL, etc.) then population minimizers trivial to compute **and** we equivalently minimise an information-theoretic loss between *distributions* in exponential families \Rightarrow embeds k -means in broad data generating processes

Can even be generalized above exponential families, to deformed and q -exponential families, keeping the Bregman divergence formulation of loss

Exponential families & k -means clustering

Why this work: k -means = population

set of

by finding
 $(\theta_i || c_j)$

If D a B
minimiz
loss bet
generat

ulation
-theoretic
oad data

But

Can even q -exponential families, keeping the Bregman divergence formulation of loss

q -exponential families, keeping the Bregman divergence formulation of loss

Exponential families & k -means clustering

Why this work: k -means = popular

set of

If D a B

minimiz

loss bet

generat

Can ever

q -expon

Universal modeling with Bregman divergences leads to **universal drawbacks**, such as lack of robustness to outliers (a Bregman divergence *lacks* robustness)

by finding $(\theta_i \| c_j)]$

ulation

-theoretic

oad data

... to deformed and

keeping the Bregman divergence formulation of loss

Exponential families & k -means clustering



Why this work: k -means = popular

set of

If D a B

minimiz

loss bet

generat

Can ever

q -expon

Our objective: get a *generalisation* of the **complete** framework, with optional additional properties for clustering such as robustness

by finding
 $(\theta_i \| c_j)]$

ulation

-theoretic

load data

... to deformed and

... keeping the Bregman divergence formulation of loss

Exponential families & k -means clustering



Why this work: k -means = popular

set of

If D a B

minimiz


loss bet

generat

Can ever

q -exponential families, keeping the Breg

by finding
 $(\theta_i \| c_j)]$

Our objective: get a *generalisation* of
the *complete*  *Generalizing*
optional addition
clustering suc

- Distributions (exponential families)
- The information theoretic distortion between distributions (KL divergence)
- Parameter distortions (Bregman divergences)
- The information-geometric / information theoretic link between clustering parameters and distributions

+ Get additional properties (robustness)

From Exponential families to Tempered Exponential Measures

Axiomatic characterization

Set of probability measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{p}}[\phi] \doteq \int \phi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x}) d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+constraint to maximize entropy

$$H(P) \doteq - \int p \log p d\xi$$

\Rightarrow get an exponential family

$$p_{\boldsymbol{\theta}}(\mathbf{x}) \propto \exp(\boldsymbol{\theta}^{\top} \phi(\mathbf{x}) - G(\boldsymbol{\theta}))$$

\uparrow
 $\boldsymbol{\theta} = \nabla G^{-1}(\mathbf{h})$
natural parameter

\uparrow
cumulant

From Exponential families to **Tempered Exponential Measures**

Axiomatic characterization

Set of probability measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{p}}[\phi] \doteq \int \phi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x}) d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+constraint to maximize entropy

$$H(P) \doteq - \int p \log p d\xi$$

⇒ get an exponential family

$$p_{\boldsymbol{\theta}}(\mathbf{x}) \propto \exp(\boldsymbol{\theta}^\top \phi(\mathbf{x}) - G(\boldsymbol{\theta}))$$

↑
 $\boldsymbol{\theta} = \nabla G^{-1}(\mathbf{h})$
natural parameter

↑
cumulant

Set of **unnormalized** measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{p}}[\phi] \doteq \int \phi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

$$t \in [0, 1], t^* \doteq 1/(2 - t)$$

From Exponential families to **Tempered Exponential Measures**

Axiomatic characterization

Set of probability measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{P}}[\phi] \doteq \int \phi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x}) d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+constraint to maximize entropy

$$H(P) \doteq - \int p \log p d\xi$$

\Rightarrow get an exponential family

$$p_{\boldsymbol{\theta}}(\mathbf{x}) \propto \exp(\boldsymbol{\theta}^\top \phi(\mathbf{x}) - G(\boldsymbol{\theta}))$$

$$\uparrow$$
$$\boldsymbol{\theta} = \nabla G^{-1}(\mathbf{h})$$

natural parameter


$$\uparrow$$

cumulant

Set of **unnormalized** measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{P}}[\phi] \doteq \int \phi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+maximize a generalized Tsallis entropy


$$H_t(\tilde{P}) \doteq - \int (\tilde{p} \log_t \tilde{p} - \log_{t-1} \tilde{p}) d\xi$$

$$\text{tempered log} \longrightarrow \log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1)$$

Concave, $\lim_{t \rightarrow 1} \log_t = \log$

From Exponential families to **Tempered Exponential Measures**

Axiomatic characterization

Set of probability measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{P}}[\phi] \doteq \int \phi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x}) d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+constraint to maximize entropy

$$H(P) \doteq - \int p \log p d\xi$$

⇒ get an exponential family

$$p_{\boldsymbol{\theta}}(\mathbf{x}) \propto \exp(\boldsymbol{\theta}^\top \phi(\mathbf{x}) - G(\boldsymbol{\theta}))$$

$$\uparrow$$

$$\boldsymbol{\theta} = \nabla G^{-1}(\mathbf{h})$$

natural parameter

$$\uparrow$$

cumulant

Set of **unnormalized** measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{P}}[\phi] \doteq \int \phi(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+maximize a generalized Tsallis entropy

$$H_t(\tilde{P}) \doteq - \int (\tilde{p} \log_t \tilde{p} - \log_{t-1} \tilde{p}) d\xi$$

$$\log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1)$$

➡ ⇒ get a **tempered exponential measure**

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \frac{\exp_t(\boldsymbol{\theta}^\top \phi(\mathbf{x}))}{\exp_t(G_t(\boldsymbol{\theta}))}$$

$$\uparrow$$

$$\boldsymbol{\theta} = \nabla G_t^{-1}(\mathbf{h})$$

cumulant

Google Research

Alternative expression

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}) \ominus_t G_t(\boldsymbol{\theta}))$$

$$z \ominus_t x \doteq \frac{z - x}{1 + (1 - t)x}$$

Tempered Exponential Measures

Set of **unnormalized** measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{p}}[\boldsymbol{\phi}] \doteq \int \boldsymbol{\phi}(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+maximize a generalized Tsallis entropy

$$H_t(\tilde{P}) \doteq - \int (\tilde{p} \log_t \tilde{p} - \log_{t-1} \tilde{p}) d\xi$$

$$\log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1)$$

\Rightarrow get a **tempered exponential measure**

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \frac{\exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}))}{\exp_t(G_t(\boldsymbol{\theta}))}$$

\uparrow
 $\boldsymbol{\theta} = \nabla G_t^{-1}(\mathbf{h})$
 \uparrow cumulant

Alternative expression

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}) \ominus_t G_t(\boldsymbol{\theta}))$$

$$z \ominus_t x \doteq \frac{z - x}{1 + (1 - t)x}$$

Cumulant in closed form:

$$G_t(\boldsymbol{\theta}) = (\log_t)^* \int (\exp_t)^*(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x})) d\xi$$

$$(\log_t)^*(z) \doteq t^* \log_{t^*} \left(\frac{z}{t^*} \right)$$

$$(\exp_t)^*(z) \doteq t^* \exp_{t^*} \left(\frac{z}{t^*} \right)$$

$$\exp_t(z) \doteq [1 + (1 - t)z]_+^{1/(1-t)} \quad \text{tempered exp}$$

$$[\cdot]_+ \doteq \max\{0, \cdot\} \quad \lim_{t \rightarrow 1} \exp_t = \exp$$

Tempered Exponential Measures

Set of **unnormalized** measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{p}}[\boldsymbol{\phi}] \doteq \int \boldsymbol{\phi}(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+maximize a generalized Tsallis entropy

$$H_t(\tilde{P}) \doteq - \int (\tilde{p} \log_t \tilde{p} - \log_{t-1} \tilde{p}) d\xi$$

$$\log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1)$$

\Rightarrow get a **tempered exponential measure**

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \frac{\exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}))}{\exp_t(G_t(\boldsymbol{\theta}))}$$

\uparrow $\boldsymbol{\theta} = \nabla G_t^{-1}(\mathbf{h})$
 \uparrow cumulant

Alternative expression

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}) \ominus_t G_t(\boldsymbol{\theta}))$$

$$z \ominus_t x \doteq \frac{z - x}{1 + (1 - t)x}$$

Cumulant in closed form:

$$G_t(\boldsymbol{\theta}) = (\log_t)^* \int (\exp_t)^*(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x})) d\xi$$

Total mass in closed form:

$$\int \tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) d\xi = 1 + (1 - t)(G_t(\boldsymbol{\theta}) - \boldsymbol{\theta}^\top \mathbf{h})$$

Tempered Exponential Measures

Set of **unnormalized** measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{p}}[\boldsymbol{\phi}] \doteq \int \boldsymbol{\phi}(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+maximize a generalized Tsallis entropy

$$H_t(\tilde{P}) \doteq - \int (\tilde{p} \log_t \tilde{p} - \log_{t-1} \tilde{p}) d\xi$$

$$\log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1)$$

\Rightarrow get a **tempered exponential measure**

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \frac{\exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}))}{\exp_t(G_t(\boldsymbol{\theta}))}$$

\uparrow $\boldsymbol{\theta} = \nabla G_t^{-1}(\mathbf{h})$
 \uparrow cumulant

Alternative expression

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}) \ominus_t G_t(\boldsymbol{\theta}))$$

$$z \ominus_t x \doteq \frac{z - x}{1 + (1 - t)x}$$

Cumulant in closed form:

$$G_t(\boldsymbol{\theta}) = (\log_t)^* \int (\exp_t)^*(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x})) d\xi$$

Total mass in closed form:

$$\int \tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) d\xi = 1 + (1 - t)(G_t(\boldsymbol{\theta}) - \boldsymbol{\theta}^\top \mathbf{h})$$

$$\tilde{p}_{t|\boldsymbol{\theta}}^{1/t^*} = \text{co-density}$$

Tempered Exponential Measures

Set of **unnormalized** measures satisfying a constraint on their expectation

$$\tilde{\mathcal{P}}_{t|\mathbf{h}} \doteq \left\{ \tilde{p} \left| \begin{array}{l} \mathbb{E}_{\tilde{p}}[\boldsymbol{\phi}] \doteq \int \boldsymbol{\phi}(\mathbf{x}) \tilde{p}(\mathbf{x}) d\xi = \mathbf{h}, \\ \int \tilde{p}(\mathbf{x})^{1/t^*} d\xi = 1, \\ \tilde{p}(\mathbf{x}) \geq 0, \forall \mathbf{x} \in \mathcal{X}. \end{array} \right. \right\}$$

+maximize a generalized Tsallis entropy

$$H_t(\tilde{P}) \doteq - \int (\tilde{p} \log_t \tilde{p} - \log_{t-1} \tilde{p}) d\xi$$

$$\log_t(z) \doteq \frac{1}{1-t} (z^{1-t} - 1)$$

\Rightarrow get a **tempered exponential measure**

$$\tilde{p}_{t|\boldsymbol{\theta}}(\mathbf{x}) \propto \frac{\exp_t(\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}))}{\exp_t(G_t(\boldsymbol{\theta}))}$$

\uparrow $\boldsymbol{\theta} = \nabla G_t^{-1}(\mathbf{h})$
 \uparrow cumulant

Google Research

Information Geometric Distortions (gen. Bregman divs)

Information-theoretic distortion between two TEMs, generalizing (reverse) KL divergence:

$$F_t(\tilde{P}_{t|\hat{\theta}} \| \tilde{P}_{t|\theta}) \doteq \int f \left(\frac{d\tilde{p}_{t|\hat{\theta}}}{d\xi} \oslash_t \frac{d\tilde{p}_{t|\theta}}{d\xi} \right) \cdot d\tilde{p}_{t|\theta}$$

$$f \doteq -\log_t$$
$$x \oslash_t y \doteq [x^{1-t} - y^{1-t} + 1]_+^{\frac{1}{1-t}}$$

Information Geometric Distortions (gen. Bregman divs)

Information-theoretic distortion between two TEMs, generalizing (reverse) KL divergence:

$$F_t(\tilde{P}_{t|\hat{\theta}} \| \tilde{P}_{t|\theta}) \doteq \int f \left(\frac{d\tilde{p}_{t|\hat{\theta}}}{d\xi} \oslash_t \frac{d\tilde{p}_{t|\theta}}{d\xi} \right) \cdot d\tilde{p}_{t|\theta} \quad \begin{aligned} f &\doteq -\log_t \\ x \oslash_t y &\doteq [x^{1-t} - y^{1-t} + 1]_+^{\frac{1}{1-t}} \end{aligned}$$

Theorem: for any 2 members of the same TEM family, $F_t(\tilde{P}_{t|\hat{\theta}} \| \tilde{P}_{t|\theta}) = B_{G_t}(\hat{\theta} \| \theta)$, with

$$B_{G_t}(\hat{\theta} \| \theta) \doteq \frac{G_t(\hat{\theta}) - G_t(\theta) - (\hat{\theta} - \theta)^\top \nabla G_t(\theta)}{1 + (1-t)G_t(\hat{\theta})} \quad \} \text{ Bregman divergence}$$

Clustering: population minimizers

Given training sample $\{\boldsymbol{\theta}_i\}_{i=1}^m$, we seek its left and right population minimizers, i.e. having $L_l(\boldsymbol{\theta}) \doteq \mathbb{E}_i[B_{G_t}(\boldsymbol{\theta} \parallel \boldsymbol{\theta}_i)]$; $L_r(\boldsymbol{\theta}) \doteq \mathbb{E}_i[B_{G_t}(\boldsymbol{\theta}_i \parallel \boldsymbol{\theta})]$, we want to compute

$$\underset{\uparrow}{\boldsymbol{\theta}_l} \doteq \arg \min_{\boldsymbol{\theta}} L_l(\boldsymbol{\theta}) \quad ; \quad \underset{\uparrow}{\boldsymbol{\theta}_r} \doteq \arg \min_{\boldsymbol{\theta}} L_r(\boldsymbol{\theta})$$

left population minimizer right population minimizer

Clustering: population minimizers

Given training sample $\{\theta_i\}_{i=1}^m$, we seek its left and right population minimizers, i.e. having $L_l(\theta) \doteq \mathbb{E}_i[B_{G_t}(\theta \parallel \theta_i)]$; $L_r(\theta) \doteq \mathbb{E}_i[B_{G_t}(\theta_i \parallel \theta)]$, we want to compute

$$\underset{\substack{\uparrow \\ \text{left population minimizer}}}{\theta_l} \doteq \arg \min_{\theta} L_l(\theta) \quad ; \quad \underset{\substack{\uparrow \\ \text{right population minimizer}}}{\theta_r} \doteq \arg \min_{\theta} L_r(\theta)$$

Theorem: we have

$$\theta_l = \nabla G_t^{-1}(\alpha_* \cdot \mathbb{E}_i \nabla G_t(\theta_i))$$

$$\alpha_* > 0$$

Precise interval to search (Cf paper)

Closed forms available in particular cases

$$\theta_r = \mathbb{E}_i \left[\frac{1}{\exp_t^{1-t}(G_t(\theta_i))} \cdot \theta_i \right]$$

Clustering: population minimizers... and robustness

Given training sample $\{\theta_i\}_{i=1}^m$, we seek its left and right population minimizers, i.e. having $L_l(\theta) \doteq \mathbb{E}_i[B_{G_t}(\theta|\theta_i)]$; $L_r(\theta) \doteq \mathbb{E}_i[B_{G_t}(\theta_i|\theta)]$, we want to compute

$$\underset{\substack{\uparrow \\ \text{left population minimizer}}}{\theta_l} \doteq \arg \min_{\theta} L_l(\theta) \quad ; \quad \underset{\substack{\uparrow \\ \text{right population minimizer}}}{\theta_r} \doteq \arg \min_{\theta} L_r(\theta)$$

Theorem: we have

$$\theta_l = \nabla G_t^{-1}(\alpha_* \cdot \mathbb{E}_i \nabla G_t(\theta_i)) \quad \theta_r = \mathbb{E}_i \left[\frac{1}{\exp_t^{1-t}(G_t(\theta_i))} \cdot \theta_i \right]$$

Robustness: add outlier θ_* with weight ϵ . The center moves as $\theta_{l/r}^{\text{new}} - \theta_{l/r}^{\text{old}} = \epsilon \cdot z(\theta_*)$

If the influence function, $z(\cdot)$, has bounded norm, then the center is **robust**

Clustering: population minimizers... and robustness

Given training sample $\{\boldsymbol{\theta}_i\}_{i=1}^m$, we seek its left and right population minimizers, i.e. having $L_l(\boldsymbol{\theta}) \doteq \mathbb{E}_i[B_{G_t}(\boldsymbol{\theta}||\boldsymbol{\theta}_i)]$; $L_r(\boldsymbol{\theta}) \doteq \mathbb{E}_i[B_{G_t}(\boldsymbol{\theta}_i||\boldsymbol{\theta})]$, we want to compute

$$\underset{\substack{\uparrow \\ \text{left population minimizer}}}{\boldsymbol{\theta}_l} \doteq \arg \min_{\boldsymbol{\theta}} L_l(\boldsymbol{\theta}) \quad ; \quad \underset{\substack{\uparrow \\ \text{right population minimizer}}}{\boldsymbol{\theta}_r} \doteq \arg \min_{\boldsymbol{\theta}} L_r(\boldsymbol{\theta})$$

Theorem: we have

$$\boldsymbol{\theta}_l = \nabla G_t^{-1}(\alpha_* \cdot \mathbb{E}_i \nabla G_t(\boldsymbol{\theta}_i)) \quad \boldsymbol{\theta}_r = \mathbb{E}_i \left[\frac{1}{\exp_t^{1-t}(G_t(\boldsymbol{\theta}_i))} \cdot \boldsymbol{\theta}_i \right]$$

Robustness: add outlier $\boldsymbol{\theta}_*$ with weight ϵ . The center moves as $\boldsymbol{\theta}_{l/r}^{\text{new}} - \boldsymbol{\theta}_{l/r}^{\text{old}} = \epsilon \cdot \mathbf{z}(\boldsymbol{\theta}_*)$

If the influence function, $\mathbf{z}(\cdot)$, has bounded norm, then the center is **robust**

Theorem: left robust *iff* robust for $t = 1$; right robust *iff*

Google Research

$$(G_t(\boldsymbol{\theta}) = \Omega(\|\boldsymbol{\theta}\|)) \wedge (t \neq 1)$$

TEMs: two examples with all details

TEM	Support	λ	θ	\hbar	$G_t^*(\hbar)$
1D t -exponential	$\left[0, \frac{3-2t}{(1-t)\lambda}\right]$	λ	$\frac{-\lambda}{3-2t}$	$t^* \left(\frac{3-2t}{\lambda}\right)^{2-t^*}$	$-t^* \cdot \left(\log_{\frac{1}{2-t^*}} \left(\frac{\hbar}{t^*}\right) - 1\right)$
1D t -Gaussian ($\mu = 0$)	$\left[-\frac{1}{\sqrt{1-t}}, \frac{1}{\sqrt{1-t}}\right]$	σ^2	$-\frac{t^*}{2\sigma^2}$	$(c_{t^*} \sqrt{2})^{1-t^*} \sigma^{3-t^*}$	$-\frac{t^*}{2} \cdot \left(\log_{t^{**}} (2c_{t^*}^2 \hbar) - 1\right)$

TEMs: two examples with all details

TEM	Support	λ	θ	\hbar	$G_t^*(\hbar)$
1D t -exponential	$\left[0, \frac{3-2t}{(1-t)\lambda}\right]$	λ	$\frac{-\lambda}{3-2t}$	$t^* \left(\frac{3-2t}{\lambda}\right)^{2-t^*}$	$-t^* \cdot \left(\log_{\frac{1}{2-t^*}} \left(\frac{\hbar}{t^*}\right) - 1\right)$
1D t -Gaussian ($\mu = 0$)	$\left[-\frac{1}{\sqrt{1-t}}, \frac{1}{\sqrt{1-t}}\right]$	σ^2	$-\frac{t^*}{2\sigma^2}$	$(c_{t^*} \sqrt{2})^{1-t^*} \sigma^{3-t^*}$	$-\frac{t^*}{2} \cdot \left(\log_{t^{**}} (2c_{t^*}^2 \hbar) - 1\right)$

TEM	$G_t(\theta)$	$B_{G_t}(\hat{\theta} \parallel \theta)$
1D t -exponential	$-\log_{2-t} \left((- \theta)^{\frac{1}{2-t}}\right)$	$t^* \cdot \left(\left(\frac{\hat{\theta}}{\theta}\right)^{2-t^*} - (2-t^*) \cdot \log_{t^*} \left(\frac{\hat{\theta}}{\theta}\right) - 1 \right)$
1D t -Gaussian ($\mu = 0$)	$(\log_t)^* \left(\frac{c_{t^*}}{\sqrt{-\theta}}\right)$	$\frac{t^*}{2} \cdot \left(\left(\sqrt{\frac{\hat{\theta}}{\theta}}\right)^{3-t^*} - (3-t^*) \cdot \log_{t^*} \sqrt{\frac{\hat{\theta}}{\theta}} - 1 \right)$

Two distinct generalisations of Itakura-Saito divergence !

TEMs: two examples with all details

TEM	Support	λ	θ	\hbar	$G_t^*(\hbar)$
1D t -exponential	$\left[0, \frac{3-2t}{(1-t)\lambda}\right]$	λ	$\frac{-\lambda}{3-2t}$	$t^* \left(\frac{3-2t}{\lambda}\right)^{2-t^*}$	$-t^* \cdot \left(\log_{\frac{1}{2-t^*}} \left(\frac{\hbar}{t^*}\right) - 1\right)$
1D t -Gaussian ($\mu = 0$)	$\left[-\frac{1}{\sqrt{1-t}}, \frac{1}{\sqrt{1-t}}\right]$	σ^2	$-\frac{t^*}{2\sigma^2}$	$(c_{t^*}\sqrt{2})^{1-t^*}\sigma^{3-t^*}$	$-\frac{t^*}{2} \cdot (\log_{t^{**}}(2c_{t^*}^2\hbar) - 1)$

TEM	$G_t(\theta)$	$B_{G_t}(\hat{\theta} \theta)$
1D t -exponential	$-\log_{2-t} \left((- \theta)^{\frac{1}{2-t}}\right)$	$t^* \cdot \left(\left(\frac{\hat{\theta}}{\theta}\right)^{2-t^*} - (2-t^*) \cdot \log_{t^*} \left(\frac{\hat{\theta}}{\theta}\right) - 1\right)$
1D t -Gaussian ($\mu = 0$)	$(\log_t)^* \left(\frac{c_{t^*}}{\sqrt{-\theta}}\right)$	$\frac{t^*}{2} \cdot \left(\left(\sqrt{\frac{\hat{\theta}}{\theta}}\right)^{3-t^*} - (3-t^*) \cdot \log_{t^*} \sqrt{\frac{\hat{\theta}}{\theta}} - 1\right)$

TEM	θ_l	θ_r
1D t -exponential	$-\mathbb{E}_i \left[\frac{1}{(-\theta_i)^{1-t^*}} \right] / \mathbb{E}_i \left[\frac{1}{(-\theta_i)^{2-t^*}} \right]$	$-\mathbb{E}_i \left[(-\theta_i)^{2-t^*} \right]$
1D t -Gaussian ($\mu = 0$)	$-\mathbb{E}_i \left[\frac{1}{(-\theta_i)^{\frac{1-t^*}{2}}} \right] / \mathbb{E}_i \left[\frac{1}{(-\theta_i)^{\frac{3-t^*}{2}}} \right]$	$-\frac{1}{(c_{t^*}\sqrt{t^*})^{1-t^*}} \cdot \mathbb{E}_i \left[(-\theta_i)^{\frac{3-t^*}{2}} \right]$

TEMs: two examples with all details

TEM	Support	λ	θ	\hbar	$G_t^*(\hbar)$
1D t -exponential	$\left[0, \frac{3-2t}{(1-t)\lambda}\right]$	λ	$\frac{-\lambda}{3-2t}$	$t^* \left(\frac{3-2t}{\lambda}\right)^{2-t^*}$	$-t^* \cdot \left(\log_{\frac{1}{2-t^*}} \left(\frac{\hbar}{t^*}\right) - 1\right)$
1D t -Gaussian ($\mu = 0$)	$\left[-\frac{1}{\sqrt{1-t}}, \frac{1}{\sqrt{1-t}}\right]$	σ^2	$-\frac{t^*}{2\sigma^2}$	$(c_{t^*}\sqrt{2})^{1-t^*} \sigma^{3-t^*}$	$-\frac{t^*}{2} \cdot (\log_{t^{**}}(2c_{t^*}^2 \hbar) - 1)$

TEM	$G_t(\theta)$	$B_{G_t}(\hat{\theta} \parallel \theta)$
1D t -exponential	$-\log_{2-t} \left((-\theta)^{\frac{1}{2-t}} \right)$	$t^* \cdot \left(\left(\frac{\hat{\theta}}{t^*} \right)^{2-t^*} - \left(\frac{\theta}{t^*} \right)^{2-t^*} \right)$
1D t -Gaussian ($\mu = 0$)	$(\log_t)^* \left(\frac{c_{t^*}}{\sqrt{-\theta}} \right)$	$\frac{t^*}{2} \cdot \left(\left(\frac{\hat{\theta}}{t^*} \right)^{\frac{3-t^*}{2}} - \left(\frac{\theta}{t^*} \right)^{\frac{3-t^*}{2}} \right)$

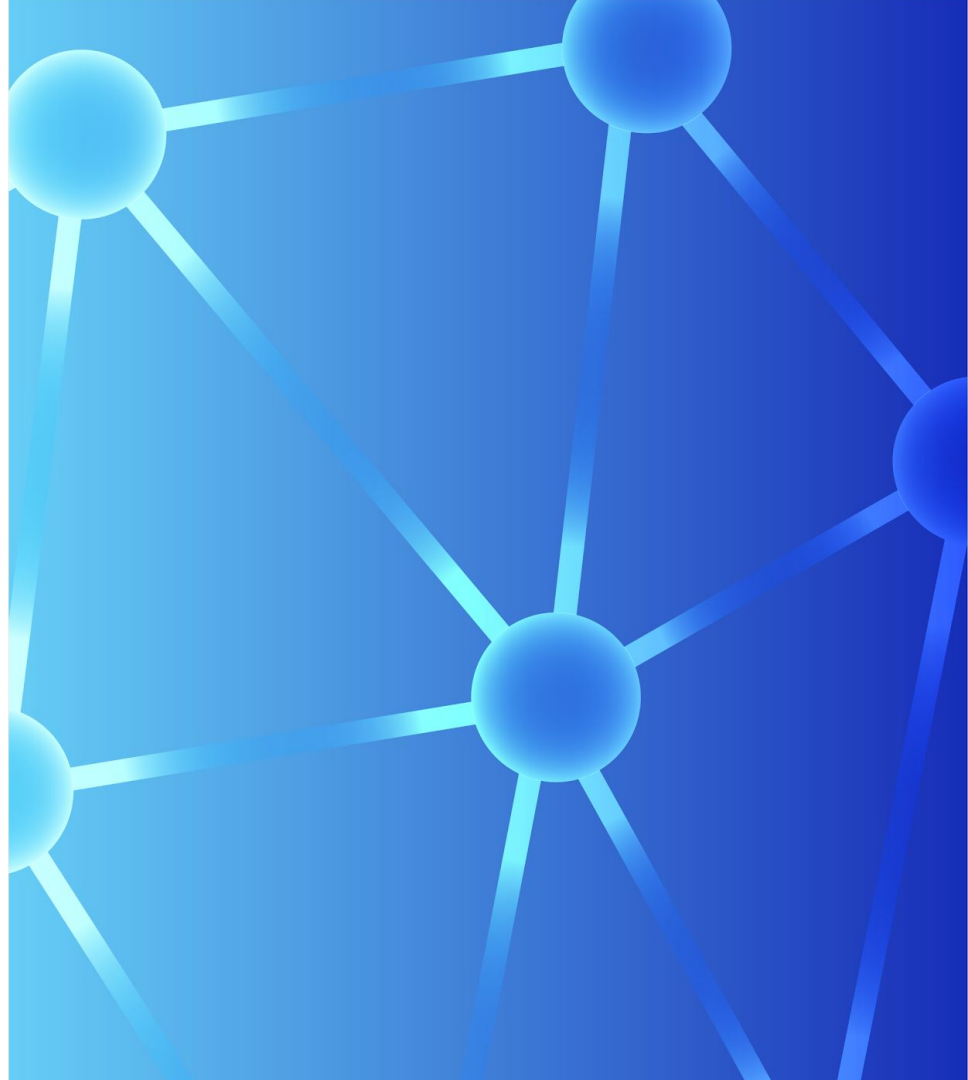
Right population minimizer: **not the arithmetic average**, unless $t=1$ (Bregman divergences)

TEM	θ_l	θ_r
1D t -exponential	$-\mathbb{E}_i \left[\frac{1}{(-\theta_i)^{1-t^*}} \right] / \mathbb{E}_i \left[\frac{1}{(-\theta_i)^{2-t^*}} \right]$	$-\mathbb{E}_i \left[(-\theta_i)^{2-t^*} \right]$
1D t -Gaussian ($\mu = 0$)	$-\mathbb{E}_i \left[\frac{1}{(-\theta_i)^{\frac{1-t^*}{2}}} \right] / \mathbb{E}_i \left[\frac{1}{(-\theta_i)^{\frac{3-t^*}{2}}} \right]$	$-\frac{1}{(c_{t^*}\sqrt{t^*})^{1-t^*}} \cdot \mathbb{E}_i \left[(-\theta_i)^{\frac{3-t^*}{2}} \right]$

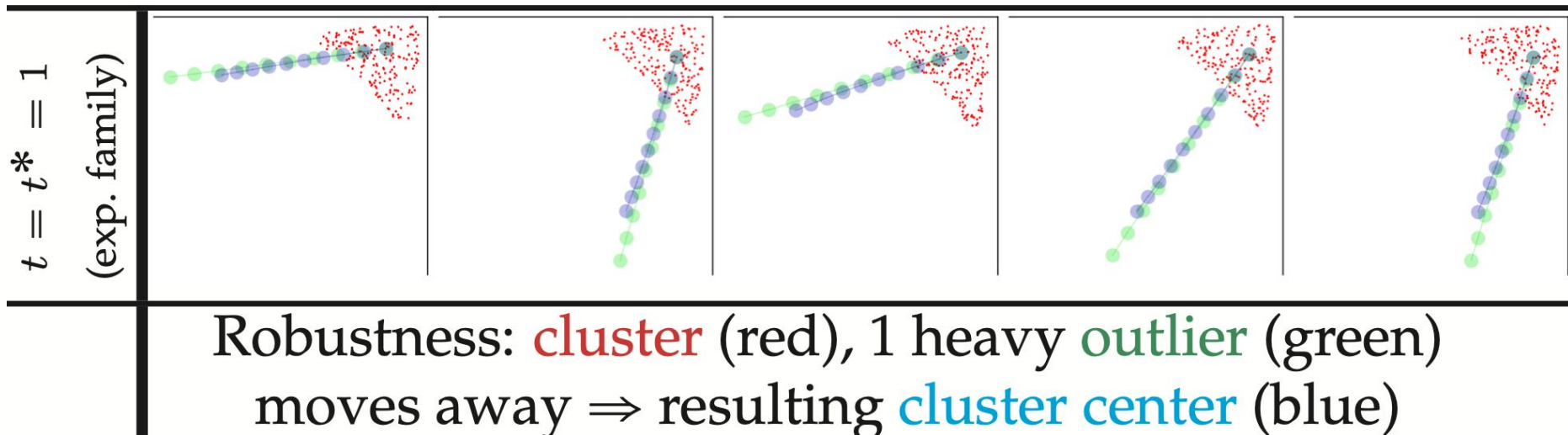
Experiments

(more in paper)

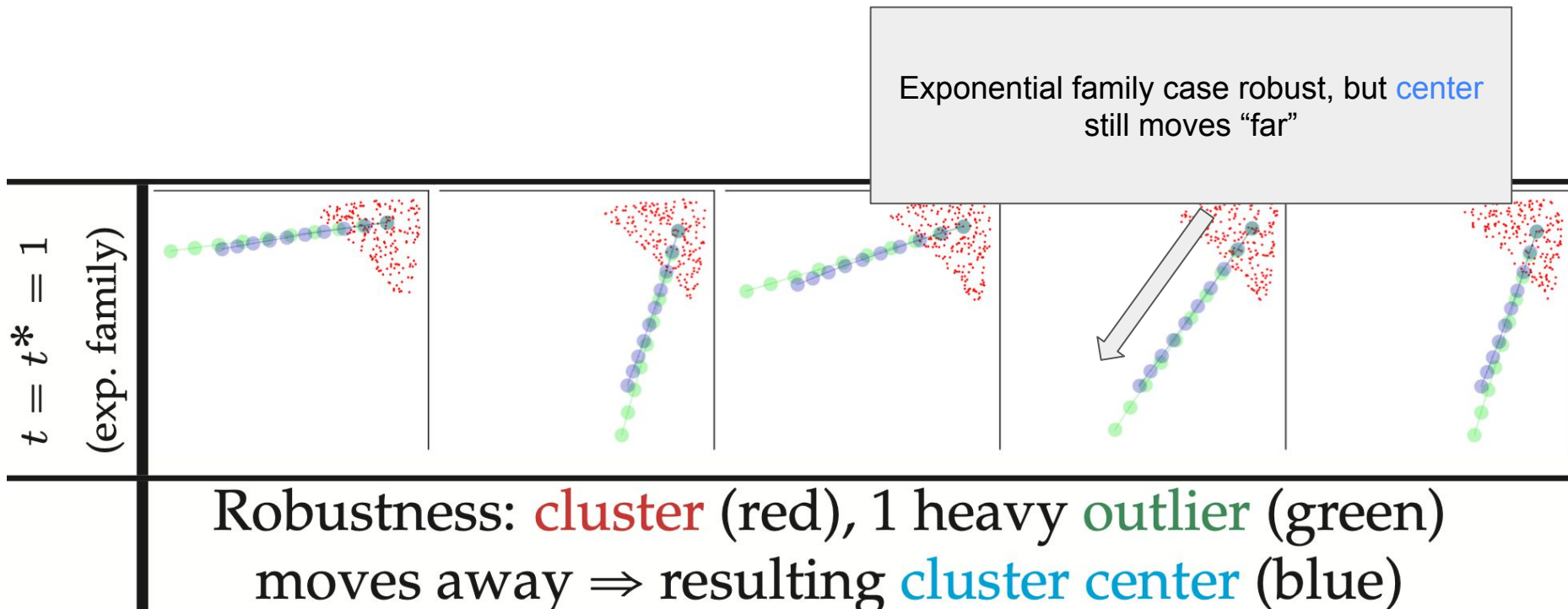
Google Research



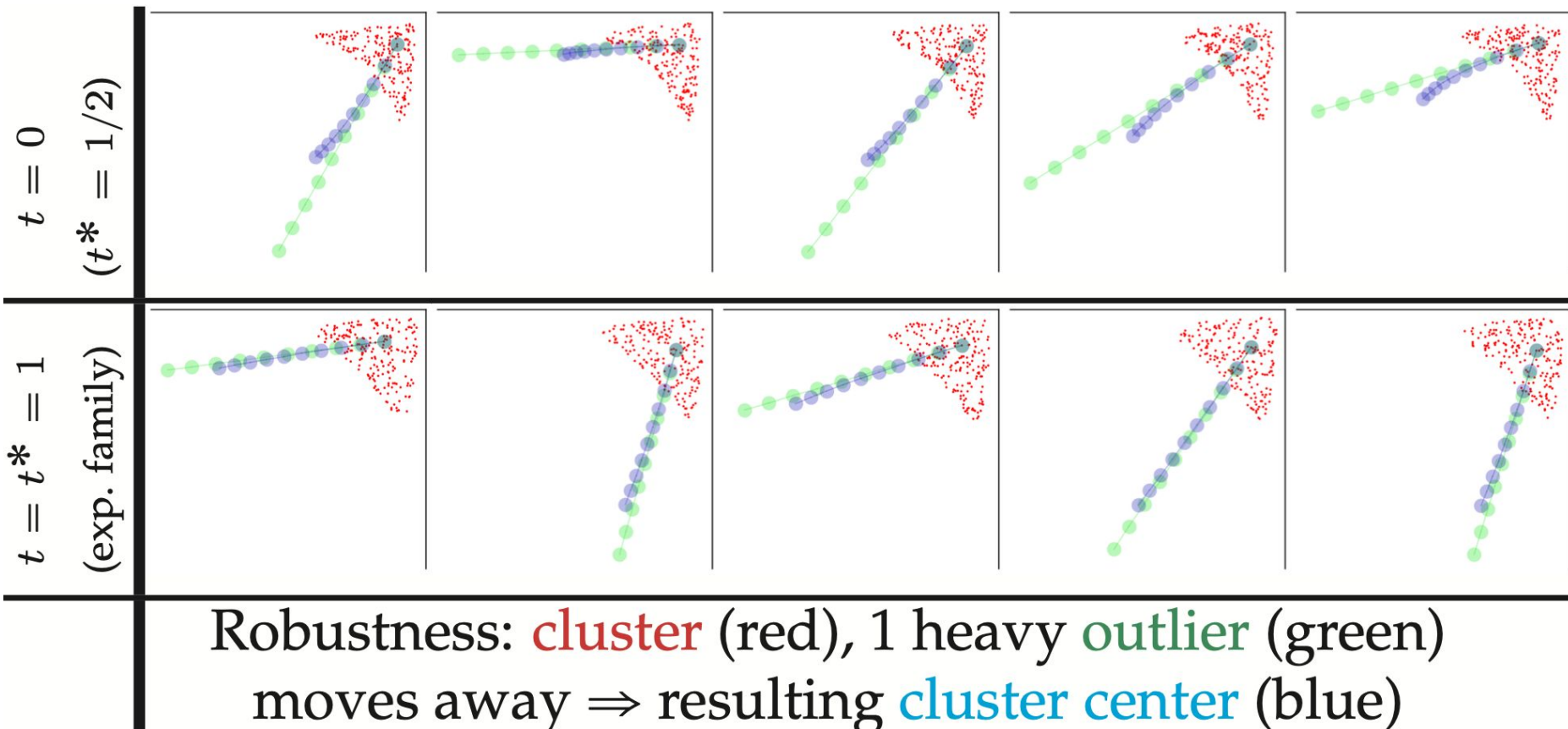
Robustness (t -exponential)



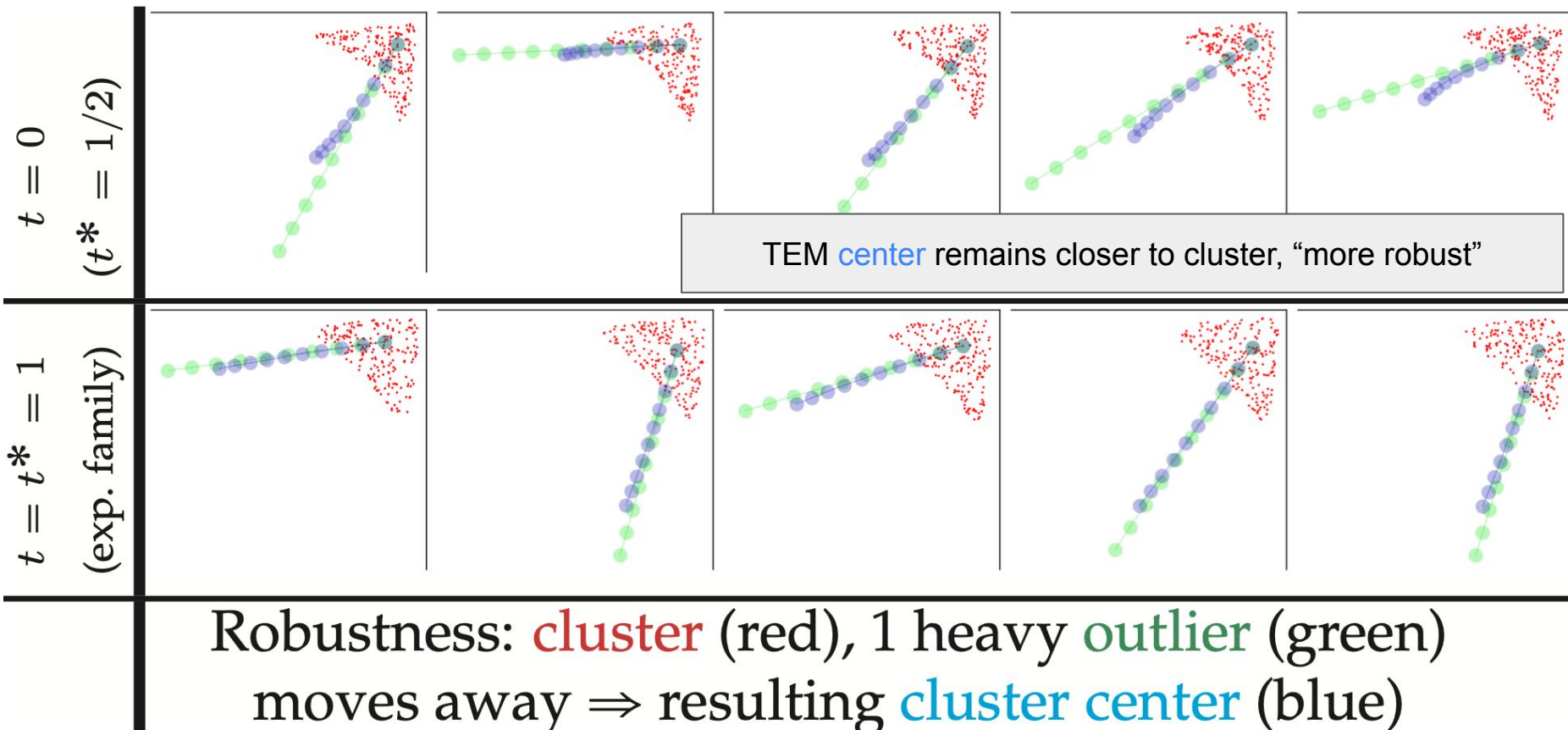
Robustness (t -exponential)



Robustness (t -exponential)



Robustness (t -exponential)



Thank You



**Ehsan
Amid**



**Richard
Nock**



**Manfred K.
Warmuth**