

**A Completed Information Projection Interpretation of Expectation  
Propagation**

**John MacLaren Walsh  
Dept. of ECE  
Drexel University  
Philadelphia, PA 19104**

# Exponential Family Densities & Rudimentary Information Geometry

- form

$$p_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = \exp(\mathbf{t}(\boldsymbol{\theta}) \cdot \boldsymbol{\lambda} - \psi_{\mathbf{t}}(\boldsymbol{\lambda})), \quad \psi_{\mathbf{t}}(\boldsymbol{\lambda}) := \log \left( \int_{\Theta} \exp(\mathbf{t}(\boldsymbol{\theta}) \cdot \boldsymbol{\lambda} - \psi_{\mathbf{t}}(\boldsymbol{\lambda})) d\boldsymbol{\theta} \right) \quad (1)$$

- duality between partition function and neg-entropy

$$\boldsymbol{\eta} := \int_{\Theta} \mathbf{t}(\boldsymbol{\theta}) \exp(\mathbf{t}(\boldsymbol{\theta}) \cdot \boldsymbol{\lambda} - \psi_{\mathbf{t}}(\boldsymbol{\lambda})) d\boldsymbol{\theta}$$

This map between  $\boldsymbol{\lambda}$  and  $\boldsymbol{\eta}$  is one to one, and may be interpreted to be the gradient of the log partition function

$$\boldsymbol{\eta} = \nabla_{\boldsymbol{\lambda}} \psi_{\mathbf{t}}(\boldsymbol{\lambda})$$

This relation indicates a Legendre transformation connection between  $\boldsymbol{\lambda}$  and  $\boldsymbol{\eta}$ . In particular, due to the convexity of the log partition function, we can form its Fenchel conjugate as

$$h(\boldsymbol{\eta}) := - \inf_{\boldsymbol{\lambda}_1} \{ \psi_{\mathbf{t}}(\boldsymbol{\lambda}_1) - \boldsymbol{\eta} \cdot \boldsymbol{\lambda}_1 \} = \psi_{\mathbf{t}}(\boldsymbol{\lambda}) - \boldsymbol{\eta} \cdot \boldsymbol{\lambda}$$

## Expectation Propagation

- joint density

$$p_{\mathbf{r}, \boldsymbol{\theta}}(\mathbf{r}, \boldsymbol{\theta}) \propto \prod_{a=1}^M f_{a, \mathbf{r}}(\boldsymbol{\theta}_a), \quad \boldsymbol{\theta}_a \subseteq \boldsymbol{\theta} \quad (2)$$

- approximate

$$p_{\mathbf{r}, \boldsymbol{\theta}}(\mathbf{r}, \boldsymbol{\theta}) \approx \prod_{a=1}^M g_{a, \lambda_a(\mathbf{r})}(\boldsymbol{\theta}_a) \quad (3)$$

- refinement rules

$$g_{a, \lambda_a} = \arg \min_{g_{a, \lambda_a}} \mathcal{D}(v_a \| q)$$

$$v_a(\boldsymbol{\theta}) := \alpha f_{a, \mathbf{r}}(\boldsymbol{\theta}_a) \prod_{c \neq a} g_{c, \lambda_c}(\boldsymbol{\theta}_c), \quad q(\boldsymbol{\theta}) := \beta \prod_{c=1}^M g_{c, \lambda_c}(\boldsymbol{\theta}_c) \quad (4)$$

$$\nabla_{\lambda_a} \mathcal{D} = \mathbb{E}_q[\mathbf{t}_a(\boldsymbol{\theta}_a)] - \mathbb{E}_{v_a}[\mathbf{t}_a(\boldsymbol{\theta}_a)] \quad (5)$$

## Bregman Divergences

- differentiable convex function is lower bounded by 1st Taylor app.  $h$  of Legendre type [1] then

$$h(\boldsymbol{\chi}) \geq h(\boldsymbol{\varsigma}) + \nabla h(\boldsymbol{\varsigma}) \cdot (\boldsymbol{\chi} - \boldsymbol{\varsigma}) \quad (6)$$

with equality if and only if  $\boldsymbol{\chi} = \boldsymbol{\varsigma}$ .

- *Bregman Divergence* [1],  $B_h$  associated with  $h$ :

$$B_h(\boldsymbol{\chi}, \boldsymbol{\varsigma}) := h(\boldsymbol{\chi}) - h(\boldsymbol{\varsigma}) - \nabla h(\boldsymbol{\varsigma}) \cdot (\boldsymbol{\chi} - \boldsymbol{\varsigma})$$

has some of the properties of a distance. In particular, we see from (6) that

$$B_h(\boldsymbol{\chi}, \boldsymbol{\varsigma}) \geq 0 \quad B_h(\boldsymbol{\chi}, \boldsymbol{\varsigma}) = 0 \Leftrightarrow \boldsymbol{\chi} = \boldsymbol{\varsigma}$$

- non-symmetric, triangle inequality in only a subset of cases
- KL Divergence: choose  $h$  as negentropy

## Method of Alternating Bregman Projections

Find points in 2 convex sets  $\mathcal{P}$  and  $\mathcal{Q}$  which min. the Bregman divergence  $B_h$  between these two sets. The projection algorithm that is often employed in this case is the *method of alternating projections* [2, 3, 4] which may be described via the iteration

$$\boldsymbol{\chi}^{(k)} := \overleftarrow{\mathbf{p}}_{\mathcal{P}} \boldsymbol{\varsigma}^{(k)}, \quad \boldsymbol{\varsigma}^{(k+1)} := \overrightarrow{\mathbf{p}}_{\mathcal{Q}} \boldsymbol{\chi}^{(k)}$$

## Dykstra's Algorithm with Cyclic Bregman Projections

$$\boldsymbol{\chi}^{(k+1)} := \overleftarrow{\mathbf{p}}_{\mathcal{C}_{k \bmod s}} \nabla h^* \left( \nabla h(\boldsymbol{\chi}^{(k)}) + \boldsymbol{\tau}^{(k+1-s)} \right) \quad (7)$$

$$\boldsymbol{\tau}^{(k+1)} := \nabla h(\boldsymbol{\chi}^{(k)}) + \boldsymbol{\tau}^{(k+1-s)} - \nabla h(\boldsymbol{\chi}^{(k+1)}) \quad (8)$$

where we initialize  $\boldsymbol{\tau}^{(-s+1)}, \dots, \boldsymbol{\tau}^{(0)} = \mathbf{0}$ . This algorithm, under some assumptions, can be shown [1] to solve the best approximation problem, in which one is seeking the point in  $\mathcal{C} := \bigcap_{i=0}^{s-1} \mathcal{C}_i$  which minimizes the Bregman divergence  $B_h$  in the first argument from the initial point  $\boldsymbol{\chi}^{(0)}$ .

*method of cyclic Bregman projections* [5][6] Instead of choosing the convex set to project on cyclicly, one may also choose it randomly [7, 8].

## Two Sets Related to EP

- Make as many copies of the space as factors
- One set: densities which are supported on all copies being equal

$$\mathcal{Q} := \{ \mathbf{b} \in \mathcal{B} | \mathbb{P}_{\mathbf{b}} [\mathbf{x}^1 = \dots = \mathbf{x}^M] = 1 \} \quad (9)$$

- Another set: product of densities of approximating family form

$$\mathcal{P} := \{ \mathbf{b} | \mathbf{b} = \exp(\boldsymbol{\lambda} \cdot \hat{\mathbf{t}}(\mathbf{x}) - \psi_{\hat{\mathbf{t}}}(\boldsymbol{\lambda})) , \boldsymbol{\lambda} \in \mathbb{R}^{M_V} \} \quad (10)$$

- Starting point: one factor per copy

# Actual Sets $\mathcal{P}$ & $\mathcal{Q}$ for 2 Bits

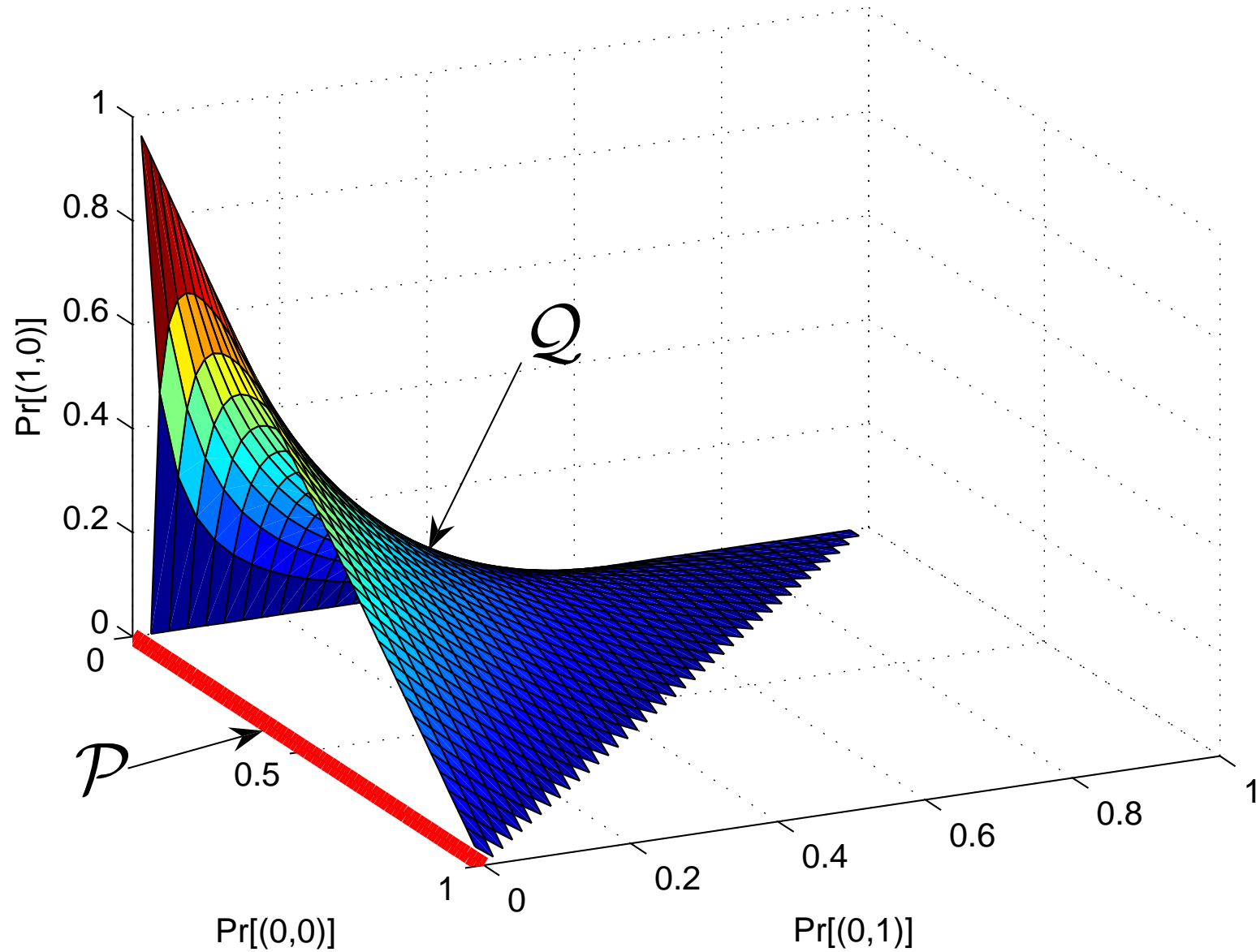


Figure 1: The sets  $\mathcal{E}_8^{\mathcal{P}}$  projects between.



## EP as a Hybrid Algorithm:

- Desired solution is projection:

$$\vec{\mathbf{p}}_{\mathcal{P}} \circ \overleftarrow{\mathbf{p}}_{\mathcal{Q}} \left( \frac{\int_{\Theta^M} \mathbf{s}(\mathbf{x}) \prod_{a=1}^M f_a(\mathbf{x}^a) d\mathbf{x}}{\int_{\Theta^M} \prod_{a=1}^M f_a(\mathbf{x}^a) d\mathbf{x}} \right) \quad (11)$$

$$\mathbb{E}_g[\mathbf{t}(\mathbf{x}^a)] = \mathbb{E}_{\boldsymbol{\theta}|\mathbf{r}}[\mathbf{t}(\boldsymbol{\theta})] \quad \forall a \in \{1, \dots, M\}$$

- EP iteratively tries to find it

$$\boldsymbol{\rho}_0, \boldsymbol{\tau}_0 = \mathbf{0}, \quad \boldsymbol{\chi}_0 = \frac{\int_{\Theta^M} \mathbf{s}(\mathbf{x}) \prod_{a=1}^M f_a(\mathbf{x}^a) d\mathbf{x}}{\int_{\Theta^M} \prod_{a=1}^M f_a(\mathbf{x}^a) d\mathbf{x}}, \quad \mathbf{k} \in \{0, 1, \dots, \}$$

$$\boldsymbol{\varsigma}_{\mathbf{k}} := \vec{\mathbf{p}}_{\mathcal{P}} \circ \nabla h^* (\nabla h(\boldsymbol{\chi}_{\mathbf{k}}) + \boldsymbol{\tau}_{\mathbf{k}}), \quad \boldsymbol{\tau}_{\mathbf{k}+1} := \nabla h(\boldsymbol{\chi}_{\mathbf{k}}) + \boldsymbol{\tau}_{\mathbf{k}} - \nabla h(\boldsymbol{\varsigma}_{\mathbf{k}})$$

$$\boldsymbol{\chi}_{\mathbf{k}+1} := \vec{\mathbf{p}}_{\mathcal{P}} \circ \overleftarrow{\mathbf{p}}_{\mathcal{Q}} \circ \nabla h^* (\nabla h(\boldsymbol{\varsigma}_{\mathbf{k}}) + \boldsymbol{\rho}_{\mathbf{k}}), \quad \boldsymbol{\rho}_{\mathbf{k}+1} := \nabla h(\boldsymbol{\varsigma}_{\mathbf{k}}) + \boldsymbol{\rho}_{\mathbf{k}} - \nabla h(\boldsymbol{\chi}_{\mathbf{k}+1})$$

- processing for left projection followed by right projection
- Can be viewed as a hybrid between alt. breg. proj. and Dykstra's w/ cyclic proj.

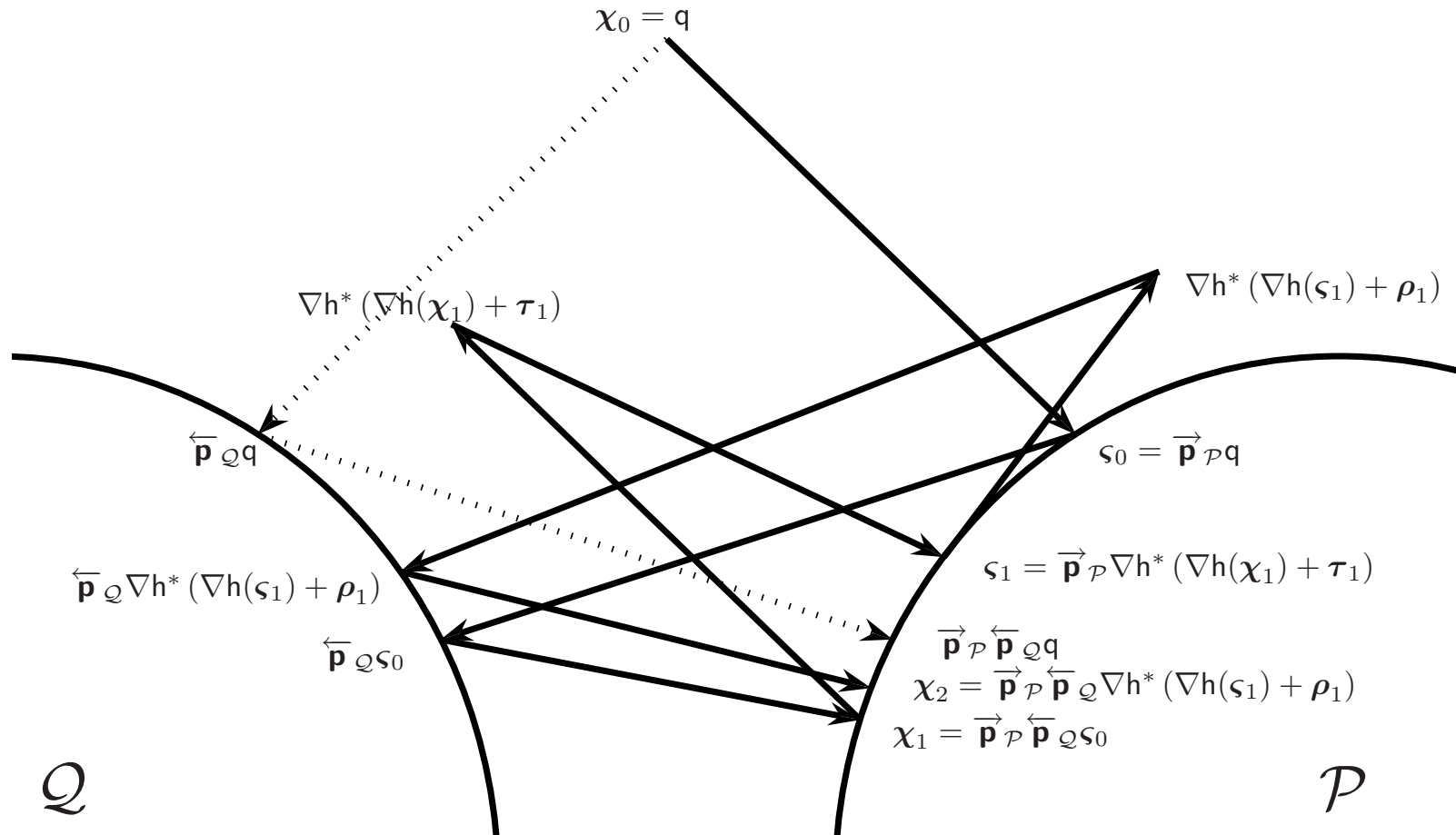


Figure 2: EP (solid arrows), and the composite projection problem it iteratively solves (dotted arrows), but with  $\|\cdot\|_2^2$  as the Bregman divergence and different sets.

### What does all this mean:

- Alternatively, could say that EP replaces a left proj. w/ a right proj. and log convex for convex from a convergent algorithm (Dijkstra's w/ cyclic proj.)
- Later could be root of occasionally good convergence behavior.
- favorite toy open problem of presenter, yours too now?

## Relationship to Prior Work (what was innovated)

- Innovated connection with Dykstra allowed for possible explanation of extrinsic information extraction within context of projection algorithm

$$\mathbf{g}_{\mathbf{a}, \lambda_{\mathbf{a}}} = \arg \min_{\mathbf{g}_{\mathbf{a}, \lambda_{\mathbf{a}}}} \mathcal{D}(\mathbf{v}_{\mathbf{a}} \parallel \mathbf{q})$$

$$\mathbf{v}_{\mathbf{a}}(\boldsymbol{\theta}) := \alpha f_{\mathbf{a}, \mathbf{r}}(\boldsymbol{\theta}_{\mathbf{a}}) \prod_{\mathbf{c} \neq \mathbf{a}} \mathbf{g}_{\mathbf{c}, \lambda_{\mathbf{c}}}(\boldsymbol{\theta}_{\mathbf{c}}), \quad \mathbf{q}(\boldsymbol{\theta}) := \beta \prod_{\mathbf{c}=1}^{\mathbf{M}} \mathbf{g}_{\mathbf{c}, \lambda_{\mathbf{c}}}(\boldsymbol{\theta}_{\mathbf{c}}) \quad (12)$$

- previous expositions had it as an intervening step amidst other projection

## So You Don't Go Home Hungry: Nonlinear Block Gauss Seidel Connection

- write down entire system for EP stationary point
- view Refinement as solution for subset of variables for subset of equations

## Convergence Theorem by Applying NLBGS Theory

Relevant references include [9], [10], and [11]. Our next theorem is an application of a theorem from [11]

$$(\boldsymbol{\lambda}_a + \boldsymbol{\gamma}_a) - \boldsymbol{\Lambda} \left( \frac{\int_{\Theta} \mathbf{t}_a(\boldsymbol{\theta}_a) \exp \left( \sum_{c=1}^M \mathbf{t}_c(\boldsymbol{\theta}_c) \cdot \boldsymbol{\lambda}_c \right) d\boldsymbol{\theta}}{\int_{\Theta} \exp \left( \sum_{c=1}^M \mathbf{t}_c(\boldsymbol{\theta}_c) \cdot \boldsymbol{\lambda}_c \right) d\boldsymbol{\theta}} \right) = \mathbf{0} \quad (13)$$

$$(\boldsymbol{\lambda}_a + \boldsymbol{\gamma}_a) - \boldsymbol{\Lambda} \left( \frac{\int_{\Theta} \mathbf{u}_a(\mathbf{y}_a) f_a(\mathbf{y}_a) \exp (\mathbf{u}_a(\mathbf{y}_a) \cdot \boldsymbol{\gamma}_a) d\mathbf{x}_a}{\int_{\Theta} f_a(\mathbf{y}_a) \exp (\mathbf{u}_a(\mathbf{y}_a) \cdot \boldsymbol{\gamma}_a) d\mathbf{x}_a} \right) = \mathbf{0} \quad (14)$$

**Thm. 1** (Convergence of Expectation Propagation Algorithms (Single Parameter Space)): If, when regarded as a function of  $[\boldsymbol{\lambda}^T, \boldsymbol{\gamma}^T]^T$ , the vector function set equal to zero in the system of equations (13) and (14) is an m-function, continuous, and surjective (onto)  $\mathbb{R}^{2v}$ , the expectation propagation algorithm converges to the unique solution of (13) and (14) and thus the unique interior critical point of the constrained optimization problem.

# References

- [1] H.H. Bauschke and A.S. Lewis, “Dykstra’s algorithm with Bregman projections: a convergence proof,” *Optimization*, , no. 48, pp. 409–427, 2000.
- [2] I. Csiszár and G. Tusnády, “Information geometry and alternating minimization procedures,” *Statistics and Decisions, Supplement Issue*, pp. 205–237, 1984.
- [3] H.H. Bauschke and J.M. Borwein, “Dykstra’s alternating projection algorithm for two sets,” *J. Approx. Theory*, , no. 79(3), pp. 418–443, 1994.
- [4] H.H. Bauschke, P. L. Combettes, and D. Noll, “Joint minimization with alternating bregman proximity operators,” <http://mip.ups-tlse.fr/noll/PAPERS/heinz.pdf>.
- [5] L. M. Bregman, “Proof of the convergence of sheleikhovskii’s method for a problem with transportation constraints,” *USSR J. Comp. Math. and Math. Phys.*, vol. 7, no. 1, pp. 147–156, 1967.
- [6] L.G. Gubin, B.T. Polyak, and E.V. Raik, “The method of projections for finding the common point of convex sets,” *USSR J. Comp. Math. and Math. Phys.*, vol. 7, no. 6, pp. 1211–1228, 1967.
- [7] H.H. Bauschke and J.M. Borwein, “Legendre functions and the method of random Bregman projections,” *J. Conv. Anal.*, , no. 4(1), pp. 27–67, 1997.
- [8] H.H. Bauschke and D. Noll, “The method of forward projections,” *J. Nonlin. and Conv. Anal.*, vol. 3, no. 2, pp. 191–205, 2002.
- [9] W. C. Rheinboldt, “On m-functions and their application to nonlinear gauss-seidel iterations and network flows,” *Journal Mathematical Analysis and Applications*, pp. 274–307, 1971.
- [10] J. J. Moré, “Nonlinear generalizations of matrix diagonal dominance with application to gauss-seidel iterations,” *SIAM Journal on Numerical Analysis*, vol. 9, pp. 357–378, June 1972.
- [11] W. C. Rheinboldt, “On classes of n-dimensional nonlinear mappings generalizing several types of matrices,” in *Proc. Symposium on the Numerical Solution of Partial Differential Equations. II.* 1970, pp. 501–546, Academic Press.