

# Supplementary Materials

## A Operator Splitting Methods for Convex Clustering

### A.1 ADMM for Convex Clustering

In this section, we derive and give the full form of the ADMM presented in Algorithm 1 for the convex clustering problem (1). We begin by noting that, in typical applications, most of the weights  $w_{ij}$  are zero and hence do not enter into the optimization problem. We can omit the  $\binom{n}{2}$ -term sum and instead write the convex clustering problem (1) as

$$\arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \left( \sum_{((i,j),w) \in \mathcal{E}} w_{ij} \|\mathbf{U}_{i \cdot} - \mathbf{U}_{j \cdot}\|_q \right)$$

where  $\mathcal{E}$  is the set of directed edges with non-zero weights  $w$  connecting  $i$  to  $j$ .

In this form, it is clear that the convex clustering problem is amenable to operator splitting methods; in particular, Chi and Lange (2015) showed that the Alternating Direction Method of Multipliers (ADMM) (Glowinski and Marroco, 1975; Gabay and Mercier, 1976; Boyd et al., 2011) works particularly well for this problem. Algorithm A1 differs from the ADMM derived in Chi and Lange (2015) in two significant ways: firstly, we only consider edges with non-zero weights, thereby greatly reducing storage requirements of the algorithm; and secondly, we implement the algorithm in “matrix-form” rather than in a fully vectorized form. These differences make the resulting algorithm both easier to derive and to read, as well as more able to take advantage of highly optimized numerical linear algebra libraries.

We note that while we are solving a matrix-valued problem, it is not a semi-definite program, and the additional complexity typically associated with matrix-valued optimization does not apply here. Because we are optimizing over the space of all matrices of a certain size,

the underlying problem is essentially Euclidean in geometry and standard (vector-valued) optimization techniques can be applied, replacing the (squared) Euclidean norm with the (squared) Frobenius norm and the standard Euclidean inner product with the Frobenius inner product as necessary.

The derivation of the ADMM for convex clustering [1](#) is relatively straight-forward. We begin by introducing an auxiliary variable  $\mathbf{V}$  containing the pairwise differences between connected rows of  $\mathbf{U}$ . The problem then becomes

$$\arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{U} - \mathbf{X}\|_F^2 + \lambda \underbrace{\left( \sum_{(e_l, w_l) \in \mathcal{E}} w_l \|\mathbf{V}_l\|_q \right)}_{P(\mathbf{V}; \mathbf{w}, q)} \quad \text{subject to } \mathbf{D}\mathbf{U} - \mathbf{V} = 0.$$

From here, we use the scaled form of the ADMM as given by a matrix version of Equations (3.5) to (3.7) of Boyd et al. ([2011](#)):

$$\begin{aligned} \mathbf{U}^{(k+1)} &= \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{U} - \mathbf{X}\|_F^2 + \frac{\rho}{2} \|\mathbf{D}\mathbf{U} - \mathbf{V}^{(k)} + \mathbf{Z}^{(k)}\|_F^2 \\ \mathbf{V}^{(k+1)} &= \arg \min_{\mathbf{V} \in \mathbb{R}^{|\mathcal{E}| \times p}} \lambda P(\mathbf{V}; \mathbf{w}, q) + \frac{\rho}{2} \|\mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V} + \mathbf{Z}^{(k)}\|_F^2 \\ \mathbf{Z}^{(k+1)} &= \mathbf{Z}^{(k)} + \mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V}^{(k+1)} \end{aligned}$$

where the dual variable is denoted by  $\mathbf{Z}$ . The analytical solution to the first subproblem is given by:

$$\mathbf{U}^{(k+1)} = (\mathbf{I} + \rho \mathbf{D}^T \mathbf{D})^{-1} (\mathbf{X} + \rho \mathbf{D}^T (\mathbf{V}^{(k)} - \mathbf{Z}^{(k)}))$$

This update is the most expensive step in the ADMM, though it can be significantly sped up by pre-calculating caching the Cholesky factorization of  $\mathbf{I} + \rho \mathbf{D}^T \mathbf{D}$  and using it at each  $\mathbf{U}$ -update:

$$\mathbf{U}^{(k+1)} = \mathbf{L}^{-T} \mathbf{L}^{-1} (\mathbf{X} + \rho \mathbf{D}^T (\mathbf{V}^{(k)} - \mathbf{Z}^{(k)})) \quad \text{where} \quad \mathbf{L}\mathbf{L}^T = \mathbf{I} + \rho \mathbf{D}^T \mathbf{D}$$

To solve the second subproblem, we note that it can be written as a proximal operator:

$$\begin{aligned} \arg \min_{\mathbf{V} \in \mathbb{R}^{|\mathcal{E}| \times n}} \lambda P(\mathbf{V}; \mathbf{w}, q) + \frac{\rho}{2} \left\| \mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V} + \mathbf{Z}^{(k)} \right\|_F^2 &= \arg \min_{\mathbf{V} \in \mathbb{R}^{|\mathcal{E}| \times n}} \frac{\lambda}{\rho} P(\mathbf{V}; \mathbf{w}, q) + \frac{\rho}{2} \left\| \mathbf{V} - \left( \mathbf{D}\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)} \right) \right\|_F^2 \\ &= \text{prox}_{\frac{\lambda}{\rho} P(\cdot; \mathbf{w}, q)} \left( \mathbf{D}\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)} \right) \end{aligned}$$

We note that, due to the row-wise structure of  $P$ , this proximal operator can be computed separately across the rows of its argument. In the cases  $q = 1$  or  $q = 2$ , the proximal operator reduces to element-wise ( $q = 1$ ) or group ( $q = 2$ ) soft-thresholding row  $l$  at the level  $w_l \lambda / \rho$ . If  $q = \infty$ , Moreau’s decomposition (Moreau, 1962) can be combined with the the efficient projection onto the  $\ell_1$  ball developed by Duchi et al. (2008) to evaluate the prox in  $\mathcal{O}(p \log p)$  steps. For other values of  $q$ , an iterative algorithm must be used.

Several stopping criteria for the ADMM have been proposed in the literature. We have found a simple stopping rule based on the change in  $\mathbf{U}$  being small sufficient in all cases. While some authors report speed-ups due to varying the ADMM relaxation parameter  $\rho$ , we have found that fixing  $\rho$  and re-using the Cholesky factor  $\mathbf{L}$  to be more efficient. Combining these steps, we obtain Algorithm A1.

## A.2 Algorithmic Regularization for Convex Clustering

In this section, we give a the full version of the CARP algorithm presented in Algorithm 2. CARP can be obtained from the standard ADMM for convex clustering (Algorithm A1) by replacing the inner ADMM loop with a single iteration. This modification gives CARP (Algorithm A2). As with Algorithm A1, we prefer to use a matrix formulation, instead of a fully vectorized formulation, to simplify the implementation and to take advantage of high-performance numerical linear algebra libraries.

---

**Algorithm A1** Warm-Started ADMM for the Convex Clustering Problem (1)

---

1. Input:
    - Data Matrix:  $\mathbf{X} \in \mathbb{R}^{n \times p}$
    - Weighted Directed Edge Set:  $\mathcal{E} = \{(e_l, w_l)\}$
    - Relaxation Parameter:  $\rho \in \mathbb{R}_{>0}$
    - Initial Regularization Parameter  $\epsilon$  and Multiplicative Step-Size  $t$
  2. Precompute:
    - Difference Matrix:  $\mathbf{D} \in \mathbb{R}^{|\mathcal{E}| \times n}$  where  $D_{ij}$  is 1 if edge  $i$  begins at node  $j$ ,  $-1$  if edge  $i$  ends at node  $j$ , and 0 otherwise
    - Cholesky Factor:  $\mathbf{L} = \text{chol}(\mathbf{I} + \rho \mathbf{D}^T \mathbf{D}) \in \mathbb{R}^{n \times n}$
  3. Initialize:
    - $\mathbf{U}^{(0)} = \mathbf{X}$
    - $\mathbf{V}^{(0)} = \mathbf{Z}^{(0)} = \mathbf{D}\mathbf{X}$
    - $l = 0, \lambda_0 = \epsilon, k = 0,$
  4. Repeat until  $\|\mathbf{V}^{(k)}\| = 0$ :
    - Repeat until convergence:
      - (i)  $\mathbf{U}^{(k+1)} = \mathbf{L}^{-T} \mathbf{L}^{-1} (\mathbf{X} + \rho \mathbf{D}^T (\mathbf{V}^{(k)} - \mathbf{Z}^{(k)}))$
      - (ii)  $\mathbf{V}^{(k+1)} = \text{prox}_{\lambda_l / \rho P(\cdot; \mathbf{w}, q)} (\mathbf{D}\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)})$
      - (iii)  $\mathbf{Z}^{(k+1)} = \mathbf{Z}^{(k)} + \mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V}^{(k+1)}$
      - (iv)  $k := k + 1$
    - Store  $\hat{\mathbf{U}}_{\lambda_l} = \mathbf{U}^{(k)}$
    - Update Regularization Parameter  $l := l + 1; \lambda_l := \lambda_{l-1} * t$
  5. Return  $\{\hat{\mathbf{U}}_{\lambda_i}\}_{i=0}^{l-1}$  as the regularization path
-

---

**Algorithm A2** CARP: Convex Clustering via Algorithmic Regularization Paths

---

1. Input:
    - Data Matrix:  $\mathbf{X} \in \mathbb{R}^{n \times p}$
    - Weighted Directed Edge Set:  $\mathcal{E} = \{(e_l, w_l)\}$
    - Relaxation Parameter:  $\rho \in \mathbb{R}_{>0}$
    - Initial Regularization Parameter  $\epsilon$  and Multiplicative Step-Size  $t$
  2. Precompute:
    - Difference Matrix:  $\mathbf{D} \in \mathbb{R}^{|\mathcal{E}| \times n}$  where  $D_{ij}$  is 1 if edge  $i$  begins at node  $j$ ,  $-1$  if edge  $i$  ends at node  $j$ , and 0 otherwise
    - Cholesky Factor:  $\mathbf{L} = \text{chol}(\mathbf{I} + \rho \mathbf{D}^T \mathbf{D}) \in \mathbb{R}^{n \times n}$
  3. Initialize:
    - $\mathbf{U}^{(0)} = \mathbf{X}$
    - $\mathbf{V}^{(0)} = \mathbf{Z}^{(0)} = \mathbf{D}\mathbf{X}$
    - $k = 0, \gamma^{(0)} = \epsilon$
  4. Repeat until  $\|\mathbf{V}^{(k)}\| = 0$ :
    - (i)  $\mathbf{U}^{(k+1)} = \mathbf{L}^{-T} \mathbf{L}^{-1} (\mathbf{X} + \rho \mathbf{D}^T (\mathbf{V}^{(k)} - \mathbf{Z}^{(k)}))$
    - (ii)  $\mathbf{V}^{(k+1)} = \text{prox}_{\gamma^{(k)}/\rho P(\cdot; \mathbf{w}, q)} (\mathbf{D}\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)})$
    - (iii)  $\mathbf{Z}^{(k+1)} = \mathbf{Z}^{(k)} + \mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V}^{(k+1)}$
    - (iv)  $k := k + 1, \gamma^{(k)} = \gamma^{(k-1)} * t$
  5. Return  $\{\mathbf{U}^{(k)}\}_{i=0}^k$  as the CARP algorithmic regularization path
-

### A.3 AMA for Convex Clustering

In addition to the AMA, Chi and Lange (2015) also show that the convex clustering problem (1) can be efficiently solved using the Alternating Minimization Algorithm (AMA) of Tseng (1991). In our notation, the AMA becomes

$$\begin{aligned} \mathbf{U}^{(k+1)} &= \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{U} - \mathbf{X}\|_F^2 + \langle \mathbf{Z}^{(k)}, \mathbf{D}\mathbf{U} - \mathbf{V}^{(k)} \rangle \\ \mathbf{V}^{(k+1)} &= \arg \min_{\mathbf{V} \in \mathbb{R}^{|\mathcal{E}| \times n}} \lambda P(\mathbf{V}; \mathbf{w}, q) + \langle \mathbf{Z}^{(k)}, \mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V}^{(k)} \rangle + \frac{\rho}{2} \|\mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V}\|_F^2 \\ \mathbf{Z}^{(k+1)} &= \mathbf{Z}^{(k)} + \rho(\mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V}^{(k+1)}) \end{aligned}$$

(Note that we use the unscaled updates for  $\mathbf{V}$ ,  $\mathbf{Z}$  here as the AMA uses different values of the relaxation parameter in the  $\mathbf{U}$  and  $\mathbf{V}$  updates. In particular, this means that the dual variables  $\mathbf{Z}$  from the ADMM are not the same as those from the AMA.) Simplifying these updates as before, the AMA becomes:

$$\begin{aligned} \mathbf{U}^{(k+1)} &= \mathbf{X} - \mathbf{D}^T \mathbf{Z}^{(k)} \\ \mathbf{V}^{(k+1)} &= \text{prox}_{\frac{\lambda}{\rho} P(\cdot; \mathbf{w}, q)} (\mathbf{D}\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)} / \rho) \\ \mathbf{Z}^{(k+1)} &= \mathbf{Z}^{(k)} + \rho(\mathbf{D}\mathbf{U}^{(k+1)} - \mathbf{V}^{(k+1)}) \end{aligned}$$

Chi and Lange (2015) note that a clever application of Moreau’s decomposition (Moreau, 1962) allows the  $\mathbf{V}$ -updates to be elided and for the AMA to be simplified into a two-step scheme. The  $\mathbf{V}^{(k)}$  iterates are key to dendrogram reconstruction, however, so such a simplification could not be used in an AMA-based CARP variant.

In our experiments, this elision is necessary for the AMA to outperform the ADMM and so, without it, the AMA does not appear to be a promising basis for an algorithmic regularization scheme. In general, the ADMM appears to converge more rapidly *per iteration* than the AMA, while the simplified AMA has much faster updates, allowing better overall

computational performance in a standard optimization scheme. Since **CARP** performs only a single iteration per regularization level, however, the faster per iteration convergence of the ADMM is more important to us than the faster calculation of the AMA.

Finally, Chi and Lange (2015) also discuss the use of accelerated variants of the ADMM and AMA (Goldstein et al., 2014) to improve convergence. Because **CARP** uses only a single iteration for each regularization level, it is not amenable to acceleration.

## B Proof of Theorem 1

In this section we prove Theorem 1 on the Hausdorff convergence of **CARP** to the convex clustering regularization path. We begin with 3 technical lemmas which may be of independent interest: Lemma 1 provides a convergence rate for the optimization step embedded within a **CARP** iteration; Lemma 2 establishes a form of Lipschitz continuity for convex clustering regularization paths; Lemma 3 provides a global bound for the approximation error induced by **CARP** at any iteration. In one step, our results are stated and proven for **CARP** with an  $\ell_2$ -fusion penalty, but can be easily extended to other  $\ell_q$ -fusion penalties.

**Lemma 1** (Q-Linear Error Decrease). *At each iteration  $k$ , the **CARP** approximation error decreases by a factor  $c < 1$  not depending on  $t$  or  $\epsilon$ . That is,*

$$\|\mathbf{U}^{(k)} - \hat{\mathbf{U}}_{\gamma^{(k)}}\| + \|\mathbf{Z}^{(k)} - \hat{\mathbf{Z}}_{\gamma^{(k)}}\| < c \left[ \|\mathbf{U}^{(k-1)} - \hat{\mathbf{U}}_{\gamma^{(k)}}\| + \|\mathbf{Z}^{(k-1)} - \hat{\mathbf{Z}}_{\gamma^{(k)}}\| \right]$$

for some  $c$  strictly less than 1.

*Proof.* By construction, each **CARP** step is a single iteration of the ADMM for the convex clustering problem (Algorithm A1) initialized at  $(\mathbf{U}^{(k-1)}, \mathbf{V}^{(k-1)}, \mathbf{Z}^{(k-1)})$ . Hence it suffices to analyze the convergence of the ADMM for the convex clustering problem and to establish linear convergence.

The convex clustering problem (1) is strongly convex due to squared Frobenius norm term. Linear convergence of the standard ADMM for strongly convex problems was first shown by Lions and Mercier (1979) and has since been refined by several other authors including Shi et al. (2014), Nishihara et al. (2015), Deng and Yin (2016), and Yang and Han (2016).

In vectorized form, with  $\mathbf{x} = \text{vec}(\mathbf{X})$ ,  $\mathbf{u} = \text{vec}(\mathbf{U})$ , and  $\mathbf{v} = \text{vec}(\mathbf{V})$ , the convex clustering problem (1) can be expressed as:

$$\arg \min_{\mathbf{u} \in \mathbb{R}^{np}, \mathbf{v} \in \mathbb{R}^{|\mathcal{E}|p}} \frac{\|\mathbf{x} - \mathbf{u}\|_2^2}{2} + \lambda \|\mathbf{v}\|_{\text{vec}(\ell_q)} \quad \text{subject to} \quad (\mathbf{I} \otimes \mathbf{D})\mathbf{u} = \mathbf{v}$$

where  $\|\cdot\|_{\text{vec}(\ell_q)}$  is an appropriately vectorized version of the row-wise  $\ell_q$  norm (a standard  $\ell_1$  norm in the case  $q = 1$  and a mixed  $\ell_q/\ell_1$  norm otherwise) and we have omitted the fusion weights for brevity.

In the notation of Hong and Luo (2017), the constraint matrix for the convex clustering problem is given by  $\mathbf{E} = \begin{pmatrix} \mathbf{I} \otimes \mathbf{D} & -\mathbf{I} \end{pmatrix}$ , for appropriately sized identity matrices, which is clearly row-independent, yielding linear convergence of the primal and dual variables at a rate  $c_\lambda < 1$  which may depend on  $\lambda$ . (We do not need to verify their additional technical assumptions as we are only using a two-block ADMM instead of the more general multi-block ADMM which is the focus of their paper.) Taking  $c = \sup_{\lambda \leq \lambda_{\max}} c_\lambda$ , we observe that the CARP iterates are uniformly Q-linearly convergent at a rate  $c$ .  $\square$

*Remark.* Recently, Deng and Yin (2016) gave a readable and precise analysis of the linear convergence of the ADMM, including estimates of the convergence rate  $c$ . The specific proof technique they employ does not strictly apply to the convex clustering problem 1, however, as the  $\mathbf{D}$  matrix is rank-deficient (excluding their Scenario 1) and the norm used for the fusion penalty is non-differentiable at the origin (excluding their Scenario 3). If an estimate of the convergence rate is required, the analysis of Deng and Yin (2016) can be applied to the convex clustering problem 1 by re-parameterizing the problem to address the rank-deficiency of  $\mathbf{D}$ . In particular, if the redundant rows of  $\mathbf{D}$  are combined (eliminating the nullspace

of  $\mathbf{D}$ ), the resulting matrix will be full-row rank, allowing Scenario 1 and Case 2 of Deng and Yin (2016) to be applied. This re-parameterization results in different split and dual variables ( $\mathbf{V}$  and  $\mathbf{Z}$ , corresponding to the  $\mathbf{D}$  matrix), however, so we do not pursue that approach here. The primal variable,  $\mathbf{U}$ , remains unchanged under this re-parameterization.

**Lemma 2** (Lipschitz Continuity of Solution Paths).  $(\hat{\mathbf{U}}_\lambda, \hat{\mathbf{Z}}_\lambda)$  is  $L$ -Lipschitz with respect to  $\lambda$ . That is,

$$\|\hat{\mathbf{U}}_{\lambda_1} - \hat{\mathbf{U}}_{\lambda_2}\| + \|\hat{\mathbf{Z}}_{\lambda_1} - \hat{\mathbf{Z}}_{\lambda_2}\| \leq L * |\lambda_1 - \lambda_2|$$

for some  $L > 0$ .

We note that this not the only form of Lipschitz continuity commonly considered for regularized estimation problems. In particular, Lipschitz continuity of the *solution with respect to the data* is a key element of various consistency results, while Lipschitz continuity of the *objective function with respect to the parameters* is a key assumption used to prove convergence of many optimization schemes.

*Proof.* It suffices to prove Lipschitz continuity of  $\hat{\mathbf{U}}_\lambda$  and  $\hat{\mathbf{Z}}_\lambda$  separately and then take the sum of their Lipschitz moduli as the joint Lipschitz modulus.

We first show that  $\hat{\mathbf{U}}_\lambda$  is Lipschitz. In vectorized form, the convex clustering problem is

$$\hat{\mathbf{u}}_\lambda = \arg \min_{\mathbf{u} \in \mathbb{R}^{np}} \frac{1}{2} \|\mathbf{u} - \mathbf{x}\|_2^2 + \lambda f_q(\tilde{\mathbf{D}}\mathbf{u})$$

where  $\mathbf{u} = \text{vec}(\mathbf{U})$ ,  $\mathbf{x} = \text{vec}(\mathbf{X})$ ,  $f_q$  is a convex function, and  $\tilde{\mathbf{D}} = \mathbf{I} \otimes \mathbf{D}$  is a fixed matrix (cf. Tan and Witten, 2015).

The KKT conditions give

$$0 \in \mathbf{u}_\lambda - \mathbf{x} + \lambda \tilde{\mathbf{D}}^T \partial f_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda)$$

where  $\partial f_q(\cdot)$  is the subdifferential of  $f_q$ . Since  $f_q$  is convex, it is differentiable almost every-

where (Rockafellar, 1970, Theorem 25.5), so the following holds for almost all  $\mathbf{u}_\lambda$ :

$$0 = \mathbf{u}_\lambda - \mathbf{x} + \lambda \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda)$$

Differentiating with respect to  $\lambda$ , we obtain (c.f. Rosset and Zhu, 2007)

$$\begin{aligned} 0 &= \mathbf{u}_\lambda - \mathbf{x} + \lambda \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \\ \frac{\partial}{\partial \lambda} [0] &= \frac{\partial}{\partial \lambda} \left[ \mathbf{u}_\lambda - \mathbf{x} + \lambda \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \right] \\ 0 &= \frac{\partial \mathbf{u}_\lambda}{\partial \lambda} - 0 + \lambda \frac{\partial}{\partial \lambda} \left[ \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \right] + \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \\ 0 &= \frac{\partial \mathbf{u}_\lambda}{\partial \lambda} + \lambda \tilde{\mathbf{D}}^T f''_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \tilde{\mathbf{D}} \frac{\partial \mathbf{u}_\lambda}{\partial \lambda} + \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \\ \implies \frac{\partial \mathbf{u}}{\partial \lambda} &= -[\mathbf{I} + \lambda \tilde{\mathbf{D}}^T f''_q(\tilde{\mathbf{D}}\mathbf{u}) \tilde{\mathbf{D}}]^{-1} \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}). \end{aligned}$$

Note that  $\mathbf{u}_\lambda$  depends on  $\lambda$  so the chain rule must be used here. From here, we note

$$\left\| \frac{\partial \mathbf{u}_\lambda}{\partial \lambda} \right\|_\infty = \left\| -[\mathbf{I} + \lambda \tilde{\mathbf{D}}^T f''_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \tilde{\mathbf{D}}]^{-1} \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \right\|_\infty \leq \| -[\mathbf{I} + 0]^{-1} \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \|_\infty = \| \tilde{\mathbf{D}}^T f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda) \|_\infty.$$

For the convex clustering problem, we recall that  $f_q(\cdot)$  is a norm and hence has bounded gradient; hence  $f'_q(\tilde{\mathbf{D}}\mathbf{u}_\lambda)$  is bounded so the gradient of the regularization path is bounded and exists almost everywhere. This implies that the regularization path is *piecewise* Lipschitz. Since the solution path is constant for  $\lambda \geq \lambda_{\max}$  and is continuous (Chi and Lange, 2015, Proposition 2.1), the solution path is globally Lipschitz with a Lipschitz modulus equal to the maximum of the piecewise Lipschitz moduli.

A similar argument shows Lipschitz continuity of  $\hat{\mathbf{Z}}_\lambda$  or one can use the relationships between  $\hat{\mathbf{U}}_\lambda$  and  $\hat{\mathbf{Z}}_\lambda$  discussed in Section 2.1 of Tan and Witten (2015).

□

**Lemma 3** (Global Error Bound). *The following error bound holds for all  $k$ :*

$$\|\mathbf{U}^{(k)} - \hat{\mathbf{U}}_{\gamma^{(k)}}\| + \|\mathbf{Z}^{(k)} - \hat{\mathbf{Z}}_{\gamma^{(k)}}\| \leq c^k L\epsilon + L(t-1)\epsilon t^k \sum_{i=1}^{k-1} \left(\frac{c}{t}\right)^i$$

*Proof.* Throughout, we let

$$\hat{\mathbf{Y}}_\lambda = \begin{pmatrix} \hat{\mathbf{U}}_\lambda \\ \hat{\mathbf{Z}}_\lambda \end{pmatrix} \text{ and } \hat{\mathbf{Y}}^{(k)} = \begin{pmatrix} \mathbf{U}^{(k)} \\ \mathbf{Z}^{(k)} \end{pmatrix}.$$

Our proof proceeds by induction on  $k$ . First note that, at initialization:

$$\|\mathbf{Y}^{(0)} - \hat{\mathbf{Y}}_\epsilon\| \leq L\epsilon$$

by Lemma 2.

Next, at  $k = 1$ , we note that

$$\|\mathbf{Y}^{(1)} - \hat{\mathbf{Y}}_{t\epsilon}\| \leq c\|\mathbf{Y}^{(0)} - \hat{\mathbf{Y}}_{t\epsilon}\|$$

by Lemma 1. We now use the triangle inequality to split the right hand side:

$$\|\mathbf{Y}^{(0)} - \hat{\mathbf{Y}}_{t\epsilon}\| \leq \underbrace{\|\mathbf{Y}^{(0)} - \hat{\mathbf{Y}}_\epsilon\|}_{\text{RHS-1}} + \underbrace{\|\hat{\mathbf{Y}}_\epsilon - \hat{\mathbf{Y}}_{t\epsilon}\|}_{\text{RHS-2}}$$

From above, we have RHS-1  $\leq L\epsilon$ . Using Lemma 2, RHS-2 can be bounded by

$$\|\hat{\mathbf{Y}}_\epsilon - \hat{\mathbf{Y}}_{t\epsilon}\| \leq L|t\epsilon - \epsilon| = L(t-1)\epsilon.$$

Putting these together, we get

$$\|\mathbf{Y}^{(1)} - \hat{\mathbf{Y}}_{t\epsilon}\| \leq c[\text{RHS-1} + \text{RHS-2}] \leq c[L\epsilon + L(t-1)\epsilon] = cLte$$

Repeating this argument for  $k = 2$ , we see

$$\begin{aligned}
\|\mathbf{Y}^{(2)} - \hat{\mathbf{Y}}_{t^2\epsilon}\| &\leq c\|\mathbf{Y}^{(1)} - \hat{\mathbf{Y}}_{t^2\epsilon}\| \\
&\leq c\left[\|\mathbf{Y}^{(1)} - \hat{\mathbf{Y}}_{t\epsilon}\| + \|\hat{\mathbf{Y}}_{t\epsilon} - \hat{\mathbf{Y}}_{t^2\epsilon}\|\right] \\
&\leq c[cLt\epsilon + L|t^2\epsilon - t\epsilon|] \\
&= c^2Lt\epsilon + cL(t-1)\epsilon * t \\
&= c^2Lt\epsilon + L\epsilon(t-1)t^2 * \left(\frac{c}{t}\right) \\
&= c^2Lt\epsilon + L\epsilon(t-1)t^2 * \sum_{i=1}^{k-1} \left(\frac{c}{t}\right)^i
\end{aligned}$$

We use this as a base case for our inductive proof and prove the general case:

$$\begin{aligned}
\|\mathbf{Y}^{(k)} - \hat{\mathbf{Y}}_{t^k\epsilon}\| &\leq c\|\mathbf{Y}^{(k-1)} - \hat{\mathbf{Y}}_{t^k\epsilon}\| \\
&\leq c\left[\|\mathbf{Y}^{(k-1)} - \hat{\mathbf{Y}}_{t^{k-1}\epsilon}\| + \|\hat{\mathbf{Y}}_{t^{k-1}\epsilon} - \hat{\mathbf{Y}}_{t^k\epsilon}\|\right] \\
&\leq c\left[c^{k-1}Lt\epsilon + L\epsilon(t-1)t^{k-1} \sum_{i=1}^{k-2} \left(\frac{c}{t}\right)^i + L|t^k\epsilon - t^{k-1}\epsilon|\right] \\
&= c^kLt\epsilon + cL\epsilon(t-1)t^{k-1} \sum_{i=1}^{k-2} \left(\frac{c}{t}\right)^i + cL\epsilon(t^k - t^{k-1}) \\
&= c^kLt\epsilon + L\epsilon(t-1)t^k \left[\frac{c}{t} \sum_{i=1}^{k-2} \left(\frac{c}{t}\right)^i + \frac{c}{t}\right] \\
&= c^kLt\epsilon + L\epsilon(t-1)t^k \left[\sum_{i=2}^{k-1} \left(\frac{c}{t}\right)^i + \frac{c}{t}\right] \\
&= c^kLt\epsilon + L\epsilon(t-1)t^k \sum_{i=1}^{k-1} \left(\frac{c}{t}\right)^i
\end{aligned}$$

Expanding the definitions of  $\mathbf{Y}^{(k)}$ ,  $\hat{\mathbf{Y}}_\lambda$ , we get the desired result.

□

With these results, we are now ready to prove Theorem 1:

**Theorem 1.** As  $(t, \epsilon) \rightarrow (1, 0)$ , where  $t$  is the multiplicative step-size update and  $\epsilon$  is the initial regularization level, the primal and dual CARP paths converge to the primal and dual convex clustering paths in the Hausdorff metric: that is,

$$d_H(\{\mathbf{U}^{(k)}\}, \{\hat{\mathbf{U}}_\lambda\}) \equiv \max \left\{ \sup_\lambda \inf_k \left\| \mathbf{U}^{(k)} - \hat{\mathbf{U}}_\lambda \right\|, \sup_k \inf_\lambda \left\| \mathbf{U}^{(k)} - \hat{\mathbf{U}}_\lambda \right\| \right\} \xrightarrow{(t, \epsilon) \rightarrow (1, 0)} 0$$

$$d_H(\{\mathbf{Z}^{(k)}\}, \{\hat{\mathbf{Z}}_\lambda\}) \equiv \max \left\{ \sup_\lambda \inf_k \left\| \mathbf{Z}^{(k)} - \hat{\mathbf{Z}}_\lambda \right\|, \sup_k \inf_\lambda \left\| \mathbf{Z}^{(k)} - \hat{\mathbf{Z}}_\lambda \right\| \right\} \xrightarrow{(t, \epsilon) \rightarrow (1, 0)} 0$$

where  $\mathbf{U}^{(k)}, \mathbf{Z}^{(k)}$  are the values of the  $k^{\text{th}}$  CARP iterate and  $\hat{\mathbf{U}}_\lambda, \hat{\mathbf{Z}}_\lambda$  are the exact solutions to the convex clustering problem (1) and its dual at  $\lambda$ .

*Proof of Theorem 1.* It suffices to show that  $\{\mathbf{\Upsilon}^{(k)}\}, \{\hat{\mathbf{\Upsilon}}_\lambda\}$  converge in the Hausdorff metric to show that the primal and dual paths converge separately. We break our proof into three steps:

- i.  $\sup_\lambda \inf_k \left\| \mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_\lambda \right\| \rightarrow 0$ ;
- ii.  $\epsilon t^{k^*}$  remains bounded as  $t, \epsilon$  decrease and  $k^*$  increases, where  $k^*$  is the iteration at which CARP terminates; and
- iii.  $\sup_k \inf_\lambda \left\| \mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_\lambda \right\| \rightarrow 0$ .

Together, these give the desired result.

**Step i.** We first show that

$$\sup_\lambda \inf_k \left\| \mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_\lambda \right\|$$

tends to zero. We begin by fixing temporarily  $\lambda$  and bounding

$$\inf_k \left\| \mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_\lambda \right\|$$

The infimum over all  $k$  is less than the distance at any particular  $k$ , so it suffices to choose a value of  $k$  which gives convergence to 0. Let  $\tilde{k}$  be the value of  $k$  which gives the closest

value of  $\gamma^{(k)}$  to  $\lambda$  along the CARP path; and let  $\tilde{\lambda} = \gamma^{(\tilde{k})} = \epsilon t^{\tilde{k}}$ . That is,

$$\tilde{k} = \arg \min_k |\gamma^{(k)} - \lambda| \quad \text{and} \quad \tilde{\lambda} = \gamma^{(\tilde{k})}$$

Then

$$\inf_k \left\| \mathbf{Y}^{(k)} - \hat{\mathbf{Y}}_\lambda \right\| \leq \left\| \mathbf{Y}^{(\tilde{k})} - \hat{\mathbf{Y}}_\lambda \right\| \leq \underbrace{\left\| \mathbf{Y}^{(\tilde{k})} - \hat{\mathbf{Y}}_{\tilde{\lambda}} \right\|}_{\text{RHS-1}} + \underbrace{\left\| \hat{\mathbf{Y}}_{\tilde{\lambda}} - \hat{\mathbf{Y}}_\lambda \right\|}_{\text{RHS-2}}$$

Using Lemma 2, we can bound RHS-2 as

$$\text{RHS-2} \leq L|\tilde{\lambda} - \lambda| \leq L|\gamma^{(\tilde{k}+1)} - \gamma^{(\tilde{k}-1)}| = L * \epsilon t^{\tilde{k}-1} * [t^2 - 1] = L * \epsilon t^{\tilde{k}-1} * [t^2 - 1] \leq L * \lambda_{\max} * [t^2 - 1]$$

Using Lemma 3, we can bound RHS-1 as

$$\begin{aligned} \text{RHS-1} &= \left\| \mathbf{Y}^{(\tilde{k})} - \hat{\mathbf{Y}}_{\tilde{\lambda}} \right\| \\ &= \left\| \mathbf{Y}^{(\tilde{k})} - \hat{\mathbf{Y}}_{\gamma^{(\tilde{k})}} \right\| \\ &= c^{\tilde{k}} L \epsilon + L(t-1) * \epsilon t^{\tilde{k}} \sum_{i=1}^{k-1} \left( \frac{c}{1+t} \right)^i \leq c^{\tilde{k}} L \epsilon + L(t-1) * \epsilon t^{\tilde{k}} * C \end{aligned} \quad (\text{A1})$$

where  $C = \sum_{i=1}^{\infty} \left( \frac{c}{1+t} \right)^i$  is large but finite. Since  $c < 1$  and  $\tilde{\lambda} = \epsilon t^{\tilde{k}} \leq \lambda_{\max}$ , we can replace the  $k$ -dependent quantities to get

$$\text{RHS-1} = \left\| \mathbf{Y}^{(\tilde{k})} - \hat{\mathbf{Y}}^{\tilde{\lambda}} \right\| \leq L \epsilon + C * L(t-1) * \lambda_{\max}$$

Putting these together, we have

$$\inf_k \left\| \mathbf{Y}^{(k)} - \hat{\mathbf{Y}}_\lambda \right\| \leq \text{RHS-1} + \text{RHS-2} \leq L \epsilon + C * L(t-1) * \lambda_{\max} + L * \lambda_{\max} * [t^2 - 1]$$

Since the right-hand side doesn't depend on  $\lambda$ , we have

$$\sup_{\lambda} \inf_k \|\Upsilon^{(k)} - \hat{\Upsilon}_{\lambda}\| \leq L\epsilon + C * L(t-1) * \lambda_{\max} + L * \lambda_{\max} * [t^2 - 1]$$

As  $(t, \epsilon) \rightarrow (1, 0)$ , we have that the right-hand side converges to zero and hence

$$\sup_{\lambda} \inf_k \|\Upsilon^{(k)} - \hat{\Upsilon}_{\lambda}\| \rightarrow 0$$

as desired.

**Step ii.** Before showing the other half of Hausdorff convergence, we pause to prove an intermediate result: As  $(t, \epsilon) \rightarrow (1, 0)$ ,  $\epsilon t^{k^*}$  remains bounded, where  $k^* = k^*(t, \epsilon)$  is the iteration at which CARP halts. For this step, we specialize to the  $\ell_2$ -case for concreteness, though our results are easily generalized.

CARP terminates when  $\|\mathbf{V}^{(k+1)}\|_{\infty, q} = \max_{i,j} \|\mathbf{U}_i^{(k+1)} - \mathbf{U}_j^{(k+1)}\|_q = 0$ ; that is, CARP terminates when all of the pairwise differences have gone to zero and the data has been fused into a single cluster.

Note that the update

$$\mathbf{V}_i^{(k+1)} = \left[ 1 - \frac{w_i \gamma^{(k)}}{\|(\mathbf{DU}^{(k+1)} + \mathbf{Z}^{(k)})_i\|_2} \right] (\mathbf{DU}^{(k+1)} + \mathbf{Z}^{(k)})_i$$

will set  $\mathbf{V}_i^{(k+1)}$  to zero when

$$\|(\mathbf{DU}^{(k+1)} + \mathbf{Z}^{(k)})_i\|_2 < w_i \gamma^{(k)}$$

Letting  $(j, k)$  be the endpoints of edge  $i$ , we find

$$\begin{aligned}
\|(\mathbf{D}\mathbf{U}^{(k+1)} + \mathbf{Z}^{(k)})_i\|_2 &= \|\mathbf{U}_{j\cdot}^{(k+1)} - \mathbf{U}_{k\cdot}^{(k+1)} + \mathbf{Z}_{i\cdot}^{(k)}\|_2 \\
&= \left\| (\mathbf{U}_{j\cdot}^{(k+1)} - \bar{\mathbf{x}}) - (\mathbf{U}_{k\cdot}^{(k+1)} - \bar{\mathbf{x}}) + \mathbf{Z}_{i\cdot}^{(k)} \right\|_2 \\
&\leq \left\| \mathbf{U}_{j\cdot}^{(k+1)} - \bar{\mathbf{x}} \right\|_2 + \left\| \mathbf{U}_{k\cdot}^{(k+1)} - \bar{\mathbf{x}} \right\|_2 + \left\| \mathbf{Z}_{i\cdot}^{(k)} \right\|_2
\end{aligned}$$

where  $\bar{\mathbf{x}}$  is the column-wise mean of  $\mathbf{X}$ .

Our strategy will be to show that this quantity is less than  $w_i\gamma^{(k)}$  for some  $k > k^*$  small enough that  $\epsilon t^k$  remains bounded and hence  $\epsilon t^{k^*}$  remains bounded. Let

$$\tilde{k} = \left\lceil \frac{\log(\lambda_{\max}/\epsilon)}{\log(t)} \right\rceil = \lceil \log_t(\lambda_{\max}/\epsilon) \rceil$$

be the first value of  $k$  such that  $\gamma^{(k)} > \lambda_{\max}$ , *i.e.*, the value of  $\lambda$  such that all of the pairwise differences have gone to zero and the data has been fused into a single cluster in the regularization path ( $\tilde{k}$  is to the regularization path as  $k^*$  is to the CARP path).

Using the bound from Equation (A1), we have that

$$\|\mathbf{U}_{j\cdot}^{(k+1)} - \bar{\mathbf{x}}\| = \|\mathbf{u}_l^{(k+1)} - (\hat{\mathbf{U}}_{\lambda_{\max}})_{j\cdot}\| < L\epsilon + C * L(t-1) * \lambda_{\max}$$

so

$$\|\mathbf{U}_{j\cdot}^{(k+1)} - \bar{\mathbf{x}}\| + \|\mathbf{U}_{k\cdot}^{(k+1)} - \bar{\mathbf{x}}\| < 2(L\epsilon + C * L(t-1) * \lambda_{\max})$$

Bounding  $\|\mathbf{Z}_{i\cdot}^{(k)}\|_2$  is more subtle, but a rough bound can be obtained again using Equation (A1) to obtain:

$$\|\mathbf{Z}_{i\cdot}^{(k)} - (\hat{\mathbf{Z}}_{\lambda_{\max}})_{i\cdot}\|_2 < L\epsilon + C * L(t-1) * \lambda_{\max}$$

so

$$\|\mathbf{Z}_{i\cdot}^{(k)}\|_2 \leq \|(\hat{\mathbf{Z}}_{\lambda_{\max}})_{i\cdot}\|_2 + L\epsilon + C * L(t-1) * \lambda_{\max}$$

Putting these together, we obtain

$$\|\mathbf{u}_l^{(k+1)} - \mathbf{u}_m^{(k+1)} - \mathbf{z}_{l,m}^{(k)}\|_2 \leq \|(\hat{\mathbf{z}}_{\lambda_{\max}})_{l,m}\|_2 + 3(L\epsilon + C * L(t-1) * \lambda_{\max})$$

To stop, we require that

$$\|(\hat{\mathbf{Z}}_{\lambda_{\max}})_i\|_2 + 3(L\epsilon + C * L(t-1) * \lambda_{\max}) < w_{l,m} \underbrace{\epsilon t^k}_{\gamma^{(k)}}$$

which occurs when

$$k > \log_t \frac{\|(\hat{\mathbf{Z}}_{\lambda_{\max}})_i\|_2 + 3(L\epsilon + C * L(t-1) * \lambda_{\max})}{w_{l,m}\epsilon}$$

Taking the max over all  $(l, m)$ -pairs we find

$$k^* \leq \max_{l,m} \log_t \frac{\|(\hat{\mathbf{z}}_{\lambda_{\max}})_{l,m}\|_2 + 3(L\epsilon + C * L(t-1) * \lambda_{\max})}{w_{l,m}\epsilon}$$

Hence it suffices to note

$$\epsilon t^{k^*} \leq \epsilon t^{\max_i \log_t \frac{\|(\hat{\mathbf{Z}}_{\lambda_{\max}})_i\|_2 + 3(L\epsilon + C * L(t-1) * \lambda_{\max})}{w_i \epsilon}} \leq \max_i \frac{\|(\hat{\mathbf{Z}}_{\lambda_{\max}})_i\|_2 + 3(L\epsilon + C * L(t-1) * \lambda_{\max})}{w_i}$$

which clearly remains bounded as  $t, \epsilon \rightarrow (1, 0)$ .

**Step iii.** With this result in hand, the proof is similar to the first half. Again, we can invoke Lemma 3 to find that

$$\inf_{\lambda} \|\Upsilon^{(k)} - \hat{\Upsilon}_{\lambda}\| \leq \|\Upsilon^{(k)} - \hat{\Upsilon}_{\epsilon t^k}\| \leq c^k L\epsilon + CL * (t-1) * \epsilon t^k$$

With the result from above,  $\epsilon t^k$  remains bounded above by some  $B < \infty$ , so

$$\sup_k \inf_{\lambda} \|\mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_{\lambda}\| = \sup_{1 \leq k \leq k^*} \inf_{\lambda} \|\mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_{\lambda}\| \leq \sup_{1 \leq k \leq k^*} c^k L \epsilon + CL * (t-1) * \epsilon t^k \leq L \epsilon + CL * (t-1) * B$$

As  $(t, \epsilon) \rightarrow (1, 0)$ , the right hand-side goes to zero so

$$\sup_k \inf_{\lambda} \|\mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_{\lambda}\| \rightarrow 0$$

Combining this with step i, we have

$$d_H(\{\mathbf{\Upsilon}^{(k)}\}, \{\hat{\mathbf{\Upsilon}}_{\lambda}\}) = \max \left\{ \sup_{\lambda} \inf_k \|\mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_{\lambda}\|, \sup_k \inf_{\lambda} \|\mathbf{\Upsilon}^{(k)} - \hat{\mathbf{\Upsilon}}_{\lambda}\| \right\} \xrightarrow{(t, \epsilon) \rightarrow (1, 0)} 0$$

as desired. □

## C CBASS: Algorithmic Regularization Paths for Convex Bi-Clustering

Having explored the computational, theoretical, and practical advantages of CARP, we now turn to the closely related problem of bi-clustering. *Bi-clustering* refers to the simultaneous clustering of rows and columns. Building on the convex clustering formulation (1), Chi et al. (2017) propose the following convex formulation of bi-clustering:

$$\hat{\mathbf{U}}_{\lambda} = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{U} - \mathbf{X}\|_F^2 + \lambda \left( \sum_{\substack{i,j=1 \\ i \neq j}}^n w_{ij} \|\mathbf{U}_{i \cdot} - \mathbf{U}_{j \cdot}\|_q + \sum_{\substack{k,l=1 \\ k \neq l}}^p \tilde{w}_{kl} \|\mathbf{U}_{\cdot k} - \mathbf{U}_{\cdot l}\|_q \right). \quad (\text{A2})$$

The second penalty term induces row fusions, similarly to how the first term induces column fusions. The resulting  $\hat{\mathbf{U}}_{\lambda}$  has a ‘checkerboard’ pattern where groups of rows and columns

are clustered together. Note that for bi-clustering the centroids are *scalars* instead of vectors as in the clustering case.

Despite their relatively similar appearances, the convex bi-clustering problem (A2) is significantly more complicated than the convex clustering problem (1) and cannot be directly solved directly using an operator splitting method. Chi et al. (2017) propose the use of the *Dykstra-Like Proximal Algorithm* (DLPA) of Bauschke and Combettes (2008) to solve the convex bi-clustering problem (A2) and refer to the resulting algorithm as **COBRA** (**C**onvex **B**i-Cluste**R**ing **A**lgorithm). COBRA works by alternating solving row- and column-wise convex clustering problems until convergence. As with convex clustering, calculating the bi-clustering solution path with sufficient accuracy to accurately reconstruct both row and column dendrograms poses significant computational burden, which is exacerbated by COBRA’s requirement to evaluate several convex clustering subproblems for each value of  $\lambda$ . While CARP could be used to solve each subproblem quickly, we would still have to run CARP many times, incurring a non-trivial total cost.

Instead, we apply the technique of algorithmic regularization to COBRA directly: we take only a single DLPA step and, within that step, we take only a single ADMM step for each of the row- and column-subproblems. We refer to the resulting algorithm as **CBASS**—**C**onvex **B**i-clustering via **A**lgorithmic **R**egularization with **S**mall **S**teps. Details of the CBASS algorithm are given in Algorithm A6 below. Our `clustRviz` software implements CBASS with and without a back-tracking step to ensure exact recovery of both the row- and column-dendrograms.

## C.1 Algorithms for Convex Bi-Clustering

The DLPA can be used to solve problems of the form

$$\text{prox}_{(f+g)(\cdot)}(\mathbf{r}) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{r}\|_2^2 + f(\mathbf{x}) + g(\mathbf{x})$$

where  $f$  and  $g$  are proximable but  $f + g$  is not using the following iterative algorithm (see also Algorithm 10.18 in Combettes and Pesquet (2011)):

---

**Algorithm A3** DLPA: Dykstra-Like Proximal Algorithm

---

1. Initialize:  $\mathbf{x}^{(0)} = \mathbf{r}$ ,  $\mathbf{p}^{(0)} = \mathbf{q}^{(0)} = 0$ ,  $k = 0$
  2. Repeat until convergence:
    - $\mathbf{y} = \text{prox}_f(\mathbf{x}^{(k)} + \mathbf{p}^{(k)})$
    - $\mathbf{p}^{(k+1)} = \mathbf{p}^{(k)} + \mathbf{x}^{(k)} - \mathbf{y}$
    - $\mathbf{x}^{(k+1)} = \text{prox}_g(\mathbf{y}^{(k+1)} + \mathbf{q}^{(k)})$
    - $\mathbf{q}^{(k+1)} = \mathbf{q}^{(k)} + \mathbf{y} - \mathbf{x}^{(k+1)}$
    - $k := k + 1$
  3. Return  $\mathbf{x}^{(k)}$
- 

To apply Algorithm A3 to convex bi-clustering (A2), we note that the problem can be rewritten as:

$$\arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{U} - \mathbf{X}\|_F^2 + \lambda \underbrace{\left( \sum_{(e_l, w_l) \in \mathcal{E}_{\text{row}}} w_l \|(\mathbf{D}_{\text{row}} \mathbf{U})_{\cdot l}\|_q \right)}_{f(\mathbf{U}) = P_{\text{row}}(\mathbf{U}; \mathbf{w}_{\text{row}}, q)} + \lambda \underbrace{\left( \sum_{(e_l, w_l) \in \mathcal{E}_{\text{col}}} w_l \|(\mathbf{U} \mathbf{D}_{\text{col}})_{\cdot l}\|_q \right)}_{g(\mathbf{U}) = P_{\text{col}}(\mathbf{U}; \mathbf{w}_{\text{col}}, q)}$$

and apply the DLPA with  $f(\mathbf{U}) = P_{\text{row}}(\mathbf{U}; \mathbf{w}_{\text{row}}, q)$  and  $g(\mathbf{U}) = P_{\text{col}}(\mathbf{U}; \mathbf{w}_{\text{col}}, q)$ . We note that  $\text{prox}_f$  is a standard convex clustering problem and can be evaluated using the ADMM or AMA approaches described above. To evaluate  $\text{prox}_g$ , we note that  $\|(\mathbf{U} \mathbf{D}_{\text{col}})_{\cdot l}\|_q = \|(\mathbf{D}_{\text{col}}^T \mathbf{U}^T)_{\cdot l}\|_q$  so we simply need to perform standard convex clustering on transposed data. The DLPA then becomes:

---

**Algorithm A4** DLPA for Convex Bi-Clustering

---

1. Initialize:  $\mathbf{U}^{(0)} = \mathbf{X}$ ,  $\mathbf{P}^{(0)} = \mathbf{Q}^{(0)} = 0$ ,  $k = 0$
  2. Repeat until convergence:
    - $\mathbf{T} = \text{Convex-Clustering}(\mathbf{U}^{(k)} + \mathbf{P}^{(k)}; \mathcal{E}_{\text{row}})$
    - $\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} + \mathbf{X}^{(k)} - \mathbf{T}$
    - $\mathbf{U}^{(k+1)} = \text{Convex-Clustering}((\mathbf{Q}^{(k)} + \mathbf{T})^T; \mathcal{E}_{\text{col}})^T$
    - $\mathbf{Q}^{(k+1)} = \mathbf{Q}^{(k)} + \mathbf{T} - \mathbf{U}^{(k+1)}$
    - $k := k + 1$
  3. Return  $\mathbf{U}^{(k)}$
- 

Expanding the Convex-Clustering steps with the ADMM from Algorithm A1 yields Algorithm

A5. To obtain CBASS from Algorithm A5, we replace the inner row- and column-subproblem loops with a single iteration of the convex clustering ADMM. Additionally, we do not reset the auxiliary  $\mathbf{U}, \mathbf{P}, \mathbf{Q}$  variables, instead carrying forward their values from each CBASS iteration to the next. These two modifications yield CBASS (Algorithm A6).

---

**Algorithm A5** Warm-Started DLPA + ADMM for the Convex Bi-Clustering Problem (A2)

---

1. Input:
    - Data Matrix:  $\mathbf{X} \in \mathbb{R}^{n \times p}$
    - Weighted Directed Edge Sets:  $\mathcal{E}_{\text{row}} = \{(e_l, w_l)\}$ ,  $\mathcal{E}_{\text{col}} = \{(e_l, w_l)\}$
    - Relaxation Parameter:  $\rho \in \mathbb{R}_{>0}$
    - Initial Regularization Parameter  $\epsilon$  and Multiplicative Step-Size  $t$
  2. Precompute:
    - Difference Matrices:  $\mathbf{D}_{\text{row}} \in \mathbb{R}^{|\mathcal{E}_{\text{row}}| \times n}$  and  $\mathbf{D}_{\text{col}} \in \mathbb{R}^{p \times |\mathcal{E}_{\text{col}}|}$
    - Cholesky Factors:  $\mathbf{L}_{\text{row}} = \text{chol}(\mathbf{I} + \rho \mathbf{D}_{\text{row}}^T \mathbf{D}_{\text{row}}) \in \mathbb{R}^{n \times n}$  and  $\mathbf{L}_{\text{col}} = \text{chol}(\mathbf{I} + \rho \mathbf{D}_{\text{col}} \mathbf{D}_{\text{col}}^T) \in \mathbb{R}^{p \times p}$
  3. Initialize:
    - $\mathbf{U}^{(0)} = \mathbf{X}$
    - $\mathbf{V}_{\text{row}}^{(0)} = \mathbf{Z}_{\text{row}}^{(0)} = \mathbf{D}_{\text{row}} \mathbf{X}$
    - $\mathbf{V}_{\text{col}}^{(0)} = \mathbf{Z}_{\text{col}}^{(0)} = (\mathbf{X} \mathbf{D}_{\text{col}})^T = \mathbf{D}_{\text{col}}^T \mathbf{X}^T$
    - $\mathbf{P}^{(0)} = \mathbf{Q}^{(0)} = \mathbf{0}$
    - $l = 0$ ,  $\lambda_0 = \epsilon$ ,  $k = 0$ ,
  4. Repeat until  $\|\mathbf{V}_{\text{row}}^{(k)}\| = \|\mathbf{V}_{\text{col}}^{(k)}\| = 0$ :
    - Repeat Until Convergence:
      - Row Sub-Problem – Repeat Until Convergence:
        - (i)  $\mathbf{T} = \mathbf{L}_{\text{row}}^{-T} \mathbf{L}_{\text{row}}^{-1} \left( \mathbf{U}^{(k)} + \mathbf{P}^{(k)} + \rho \mathbf{D}_{\text{row}}^T (\mathbf{V}_{\text{row}}^{(k)} - \mathbf{Z}_{\text{row}}^{(k)}) \right)$
        - (ii)  $\mathbf{V}_{\text{row}}^{(k+1)} = \text{prox}_{\lambda_l / \rho P(\cdot; \mathbf{w}_{\text{row}}, q)} \left( \mathbf{D}_{\text{row}} \mathbf{T} + \mathbf{Z}_{\text{row}}^{(k)} \right)$
        - (iii)  $\mathbf{Z}_{\text{row}}^{(k+1)} = \mathbf{Z}_{\text{row}}^{(k)} + \mathbf{D} \mathbf{T} - \mathbf{V}_{\text{row}}^{(k+1)}$
      - $\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} + \mathbf{U}^{(k-1)} - \mathbf{T}$
      - Column Sub-Problem – Repeat Until Convergence:
        - (i)  $\mathbf{S} = \mathbf{L}_{\text{col}}^{-T} \mathbf{L}_{\text{col}}^{-1} \left( (\mathbf{T} + \mathbf{Q}^{(k)})^T + \rho \mathbf{D}_{\text{col}} (\mathbf{V}_{\text{col}}^{(k)} - \mathbf{Z}_{\text{col}}^{(k)}) \right)$
        - (ii)  $\mathbf{V}_{\text{col}}^{(k+1)} = \text{prox}_{\lambda_l / \rho P(\cdot; \mathbf{w}_{\text{col}}, q)} \left( \mathbf{D}_{\text{col}}^T \mathbf{S} + \mathbf{Z}_{\text{col}}^{(k)} \right)$
        - (iii)  $\mathbf{Z}_{\text{col}}^{(k+1)} = \mathbf{Z}_{\text{col}}^{(k)} + \mathbf{D}_{\text{col}}^T \mathbf{S} - \mathbf{V}_{\text{col}}^{(k+1)}$
      - $\mathbf{U}^{(k+1)} = \mathbf{S}^T$
      - $\mathbf{Q}^{(k+1)} = \mathbf{Q}^{(k)} + \mathbf{T} - \mathbf{U}^{(k+1)}$
      - $k := k + 1$
    - Store  $\hat{\mathbf{U}}_{\lambda_l} = \mathbf{U}^{(k)}$
    - Reset Auxiliary Variables:  $\mathbf{U}^{(k+1)} = \mathbf{X}$ ,  $\mathbf{P}^{(k+1)} = \mathbf{Q}^{(k+1)} = \mathbf{0}$
    - Update Regularization Parameter  $\lambda_l := \lambda_{l-1} * t$ ,  $l := l + 1$
  5. Return  $\left\{ \hat{\mathbf{U}}_{\lambda_i} \right\}_{i=0}^{l-1}$  as the regularization path
-

---

**Algorithm A6** CBASS: Convex Bi-Clustering via Algorithmic regularization with Small Steps

---

1. Input:
    - Data Matrix:  $\mathbf{X} \in \mathbb{R}^{n \times p}$
    - Weighted Directed Edge Sets:  $\mathcal{E}_{\text{row}} = \{(e_l, w_l)\}$ ,  $\mathcal{E}_{\text{col}} = \{(e_l, w_l)\}$
    - Relaxation Parameter:  $\rho \in \mathbb{R}_{>0}$
    - Initial Regularization Parameter  $\epsilon$  and Multiplicative Step-Size  $t$
  2. Precompute:
    - Difference Matrices:  $\mathbf{D}_{\text{row}} \in \mathbb{R}^{|\mathcal{E}_{\text{row}}| \times n}$  and  $\mathbf{D}_{\text{col}} \in \mathbb{R}^{p \times |\mathcal{E}_{\text{col}}|}$
    - Cholesky Factors:  $\mathbf{L}_{\text{row}} = \text{chol}(\mathbf{I} + \rho \mathbf{D}_{\text{row}}^T \mathbf{D}_{\text{row}}) \in \mathbb{R}^{n \times n}$  and  $\mathbf{L}_{\text{col}} = \text{chol}(\mathbf{I} + \rho \mathbf{D}_{\text{col}} \mathbf{D}_{\text{col}}^T) \in \mathbb{R}^{p \times p}$
  3. Initialize:
    - $\mathbf{U}^{(0)} = \mathbf{X}$
    - $\mathbf{V}_{\text{row}}^{(0)} = \mathbf{Z}_{\text{row}}^{(0)} = \mathbf{D}_{\text{row}} \mathbf{X}$
    - $\mathbf{V}_{\text{col}}^{(0)} = \mathbf{Z}_{\text{col}}^{(0)} = (\mathbf{X} \mathbf{D}_{\text{col}})^T = \mathbf{D}_{\text{col}}^T \mathbf{X}^T$
    - $\mathbf{P}^{(0)} = \mathbf{Q}^{(0)} = \mathbf{0}$
    - $k = 0$ ,  $\gamma^{(0)} = \epsilon$
  4. Repeat until  $\|\mathbf{V}_{\text{row}}^{(k)}\| = \|\mathbf{V}_{\text{col}}^{(k)}\| = 0$ :
    - Row Updates:
      - (i)  $\mathbf{T} = \mathbf{L}_{\text{row}}^{-T} \mathbf{L}_{\text{row}}^{-1} \left( \mathbf{U}^{(k)} + \mathbf{P}^{(k)} + \rho \mathbf{D}_{\text{row}}^T (\mathbf{V}_{\text{row}}^{(k)} - \mathbf{Z}_{\text{row}}^{(k)}) \right)$
      - (ii)  $\mathbf{V}^{(k+1)} = \text{prox}_{\gamma^{(k)}/\rho P(\cdot; \mathbf{w}_{\text{row}}, q)} \left( \mathbf{D}_{\text{row}} \mathbf{T} + \mathbf{Z}^{(k)} \right)$
      - (iii)  $\mathbf{Z}_{\text{row}}^{(k+1)} = \mathbf{Z}_{\text{row}}^{(k)} + \mathbf{D}_{\text{row}} \mathbf{T} - \mathbf{V}_{\text{row}}^{(k+1)}$
    - $\mathbf{P}^{(k+1)} = \mathbf{P}^{(k)} + \mathbf{U}^{(k-1)} - \mathbf{T}$
    - Column Updates:
      - (i)  $\mathbf{S} = \mathbf{L}_{\text{col}}^{-T} \mathbf{L}_{\text{col}}^{-1} \left( (\mathbf{T} + \mathbf{Q}^{(k)})^T + \rho \mathbf{D}_{\text{col}} (\mathbf{V}_{\text{col}}^{(k)} - \mathbf{Z}_{\text{col}}^{(k)}) \right)$
      - (ii)  $\mathbf{V}_{\text{col}}^{(k+1)} = \text{prox}_{\gamma^{(k)}/\rho P(\cdot; \mathbf{w}_{\text{col}}, q)} \left( \mathbf{D}_{\text{col}}^T \mathbf{S} + \mathbf{Z}_{\text{col}}^{(k)} \right)$
      - (iii)  $\mathbf{Z}_{\text{col}}^{(k+1)} = \mathbf{Z}_{\text{col}}^{(k)} + \mathbf{D}_{\text{col}}^T \mathbf{S} - \mathbf{V}_{\text{col}}^{(k+1)}$
    - $\mathbf{U}^{(k+1)} = \mathbf{S}^T$
    - $\mathbf{Q}^{(k+1)} = \mathbf{Q}^{(k)} + \mathbf{T} - \mathbf{U}^{(k+1)}$
    - $k := k + 1$ ,  $\gamma^{(k)} = \gamma^{(k-1)} * t$
  5. Return  $\{\mathbf{U}^{(k)}\}_{i=0}^k$  as the CBASS algorithmic regularization path
-

## C.2 Visualizations for Convex Bi-Clustering

While it is possible to construct row- and column-wise CBASS analogues of the CARP dendrogram and path plots discussed above, the primary visualization associated with bi-clustering is the *cluster heatmap*, which combines a heatmap visualization of the raw data with independent row- and column-dendrograms (Wilkinson and Friendly, 2009). We modify the standard cluster heatmap by creating dendrograms using the fusions identified by CBASS. As Chi et al. (2017) argue, the joint estimation of dendrograms provided by convex bi-clustering often produces better results than independent dendrogram construction.

We applied CBASS to the Presidents data and show the resulting cluster heatmap in Figure A1. A close examination reveals several interesting patterns. This data clearly exhibits a bi-clustered structure, with certain words being strongly associated with certain groups of presidents. Examining the two clear bi-clusters on the left, we see that words such as “billion,” “soviet,” and “technology” are frequently used by modern presidents and rarely used by pre-modern presidents. Conversely, we see that words which may be considered somewhat antiquated, such as “vessel” or “shall,” are associated with pre-modern presidents. For data with less clear structure, the interpretability of the cluster heatmap can sometimes be increased by plotting the smoothed estimates  $\mathbf{U}^{(k)}$  rather than the raw data.

In simulation studies, CBASS appears to converge to the exact regularization path as  $t \rightarrow 1$ . While this is consistent with both our theory and observations for CARP, we leave the theoretical analysis of CBASS to future work. As far as we know, a rate of convergence has not been established for the DLPA in the optimization literature, without which the techniques used to prove Theorem 1 cannot be applied to CBASS.



## D Additional Comparisons

Figure A2 compares the accuracy of CARP, CBASS, hierarchical clustering, and  $K$ -means clustering on the TCGA and Authors data sets discussed in Section 4. While certain forms of hierarchical clustering perform well on this data, CARP achieves superior performance without requiring the user to select a distance or linkage.

Figure A3 compares the performance of CARP, hierarchical clustering with Euclidean distance and Ward’s, complete, and single linkage, and  $K$ -means on data simulated from a Gaussian mixture model. The cluster centroids were equally spaced on a 2-dimensional subspace and  $n = 54$  observations were generated from a Gaussian distribution with unit variance centered at the cluster centroid. Each of the clustering methods exhibit similar behaviors, with improved performance as the inter-cluster distance increases and decreased performance with higher ambient dimensionality or more clusters. Because these data were generated from isotropic Gaussians, all methods except single linkage hierarchical clustering perform well.

Figure A4 compares the performance of the same methods on non-convex clusters. In particular, we consider a version of the “half-moons” example proposed by Hocking et al. (2011). (See also Figure A5.) The data were generated on a two-dimensional subspace with  $n = 50$  observations from each cluster and Gaussian noise orthogonal to the signal subspace were added. Not surprisingly, the performance of all methods degrades as the degree of noise and the ambient dimensionality are increased. Despite this, we see that CARP and single-linkage hierarchical clustering clearly outperform other methods, with CARP being more robust to the presence of noise.

Comparing these two simulations, we see that only convex clustering (CARP) is able to consistently perform well on both the convex and non-convex simulated data without requiring the user to select a distance metric or linkage. This is in large part due to the sparse weighting

scheme used in the `clustRviz` package, which is able to flexibly and robustly adapt to the observed data distribution. Our findings should be contrasted with those of Tan and Witten (2015) who focus only on the case of uniform weights and show that, without informative weights, convex clustering performs similarly to single linkage convex clustering.

## E Back-Tracking, Post-Processing, and Dendrogram Construction

The `CARP-VIZ` variant of our `CARP` algorithm implements a back-tracking scheme in order to improve dendrogram recovery. Because a relatively large value of  $\lambda$  is typically required for any fusions to occur in convex clustering (1), `CARP-VIZ` begin with a large step-size (by default,  $t = 1.1$ ) and performs standard `CARP` iterations until the first fusion is identified (*i.e.*, a row of  $\mathbf{V}^{(k)}$  is set to zero). After the first fusion is identified, `CARP-VIZ` switches to a smaller step-size (by default,  $t = 1.01$ ) for the remainder of the algorithm. At each iteration, `CARP-VIZ` counts the number of fusions that occur. If more than one fusion occurs, instead of proceeding, `CARP-VIZ` attempts to determine which fusion occurred first. It does so using a back-tracking scheme, similar to those used in optimization methods. `CARP-VIZ` discards the iteration with multiple fusions, halves the step-size, and performs another iteration. If this half-step iteration has only one fusion, `CARP-VIZ` accepts it and continues as before. Otherwise, `CARP-VIZ` again halves the step-size and repeats this process until the correct order of fusions is identified (or a limit on the number of back-tracking steps is hit). Once the first fusion is identified, `CARP-VIZ` resets  $t$  and continues. `CBASS-VIZ` uses essentially the same scheme, though it checks for both row and column fusions. We have found that, because it only uses a small step-size at “interesting” parts of the solution space, this back-tracking scheme typically produces more accurate dendrogram recovery at less expense than running standard `CARP` with a very small step-size.



### Clustering Accuracy: Gaussian Mixture Model

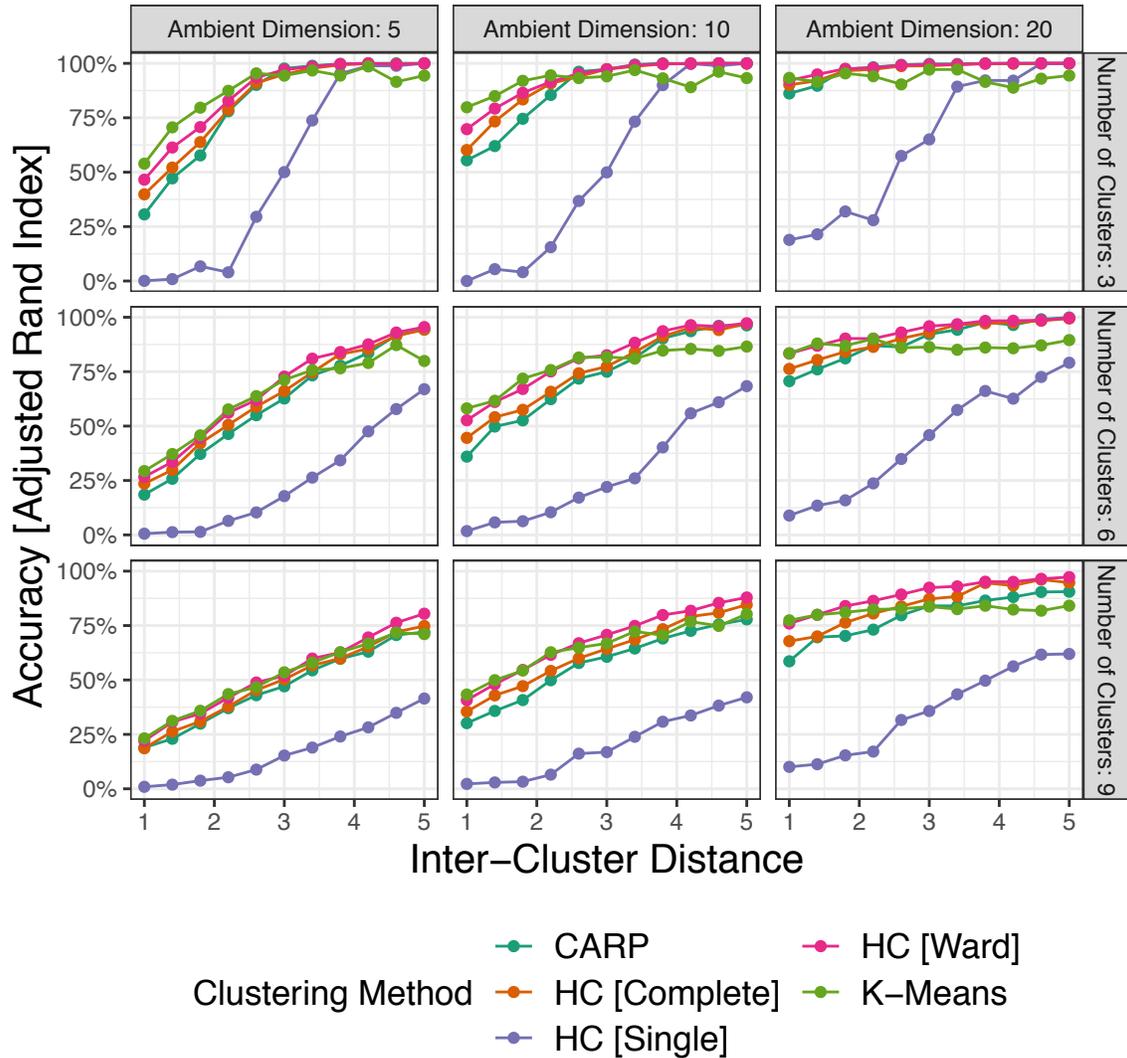


Figure A3: Accuracy of convex clustering, hierarchical clustering with Euclidean distance and several linkages, and  $K$ -means clustering (teal) on data simulated from a Gaussian mixture model, as measured by the Adjusted Rand (Hubert and Arabie, 1985) index. Because these clusters are spherical, with sufficient inter-cluster separation all methods except hierarchical clustering with single linkage perform well.

### Clustering Accuracy: Half Moons [Non Convex Clusters]

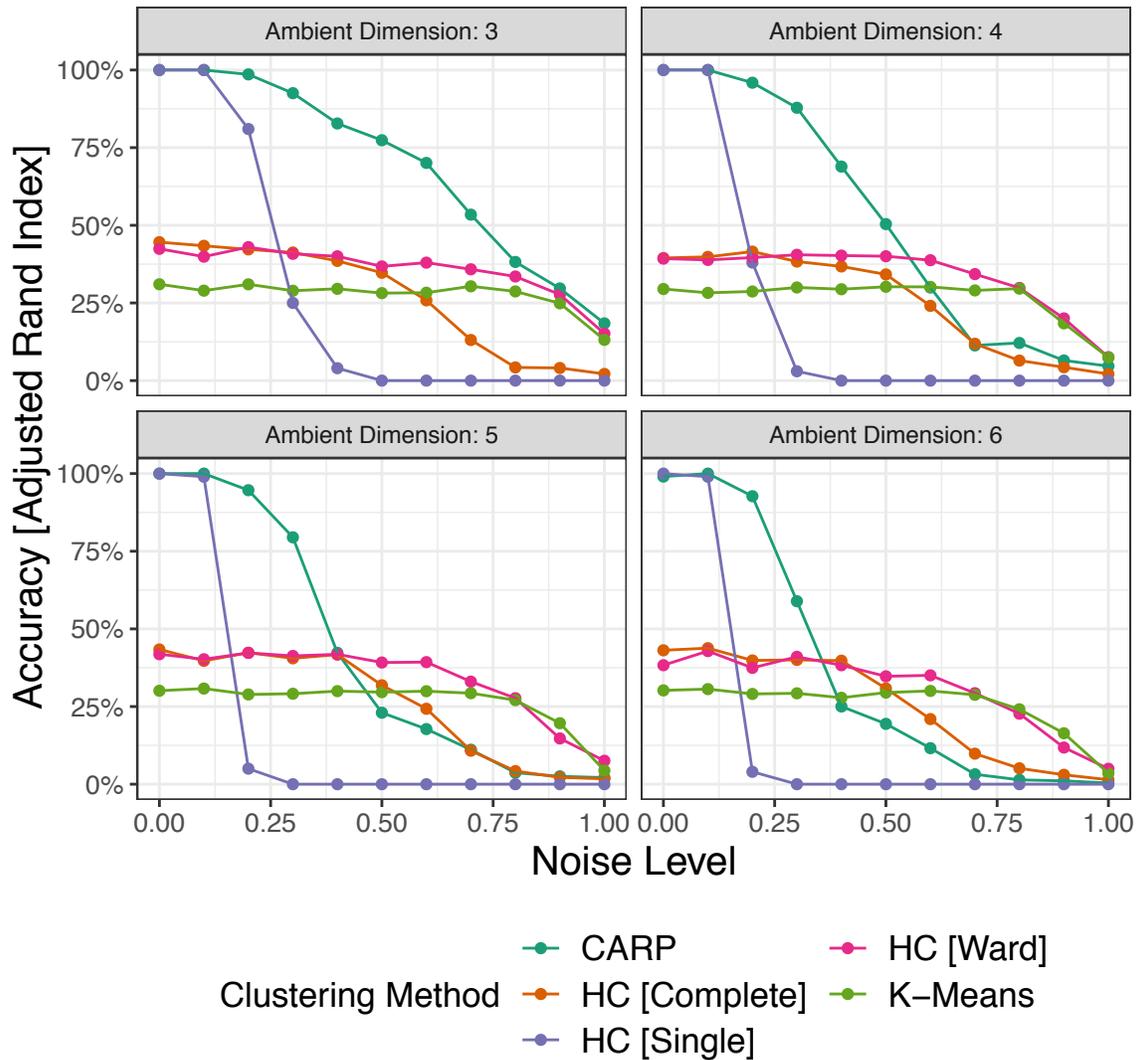


Figure A4: Accuracy of convex clustering, hierarchical clustering with Euclidean distance and several linkages, and  $K$ -means clustering (teal) on data simulated from a the two-circles and two-half-moons model, as measured by the Adjusted Rand (Hubert and Arabie, 1985) index. Note that only CARP and single linkage hierarchical clustering are able to adapt to the non-convex cluster shapes.

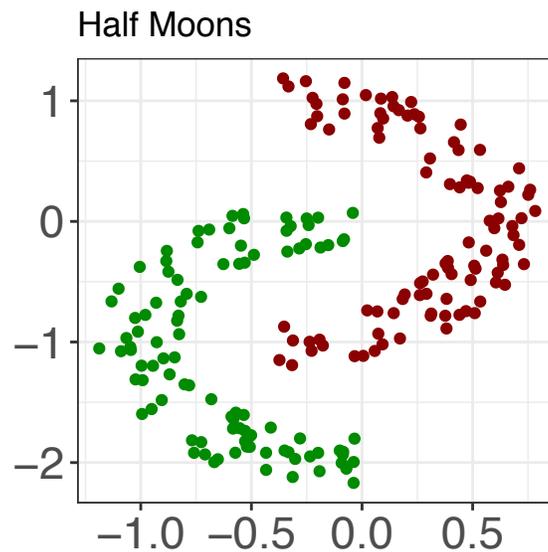


Figure A5: A sample realization of the interlocking half moons test data of Hocking et al. (2011) used for Figure A4. Data were generated from these clusters along a random two-dimensional subspace and Gaussian noise orthogonal to the signal subspace was added to increase the difficulty of the clustering problem. CARP and single-linkage hierarchical clustering are able to exactly recover the true clustering in the noiseless case; the performance of hierarchical clustering quickly degrades as more noise is added, however, while CARP is more robust.

Once `CARP` or `CARP-VIZ` terminate, `clustRviz` performs an additional post-processing step to isolate individual fusions. `clustRviz` reviews the fusions at each iteration and, if an iteration has multiple fusions, linearly interpolates between  $\mathbf{U}^{(k)}$  and  $\mathbf{U}^{(k+1)}$  to determine the approximate regularization level at which each fusion occurred. The interpolated iterate is only approximate, but is necessary for dendrogram construction. We note that no interpolation is typically needed for `CARP-VIZ` results, due to the back-tracking step used to isolate individual fusions, but, by default, `clustRviz` post-processes both `CARP` and `CARP-VIZ` output. The same post-processing scheme is applied separately to the row and column fusions from `CBASS`.

Once post-processing is performed, a dendrogram is constructed from the interpolated iterates. The dendrogram construction proceeds in the opposite order as hierarchical clustering: we begin with the fully fused data and decrease  $\gamma^{(k)}$ , noting the order in which centroids were fused. (We use the reverse ordering so that, in the rare case where the path contains fissions, the *final* fusion is reflected in the resulting dendrogram.) The dendrogram height associated with each fusion is the  $\gamma^{(k)}$  at which that fusion is first observed. Finally, we check whether fusions are more uniformly distributed on the  $\gamma^{(k)}$  scale or the  $\log(\gamma^{(k)})$  and adjust the dendrogram height accordingly to provide less cluttered visualizations.

Since the weight selection, post-processing, and dendrogram reconstruction steps could potentially be applied to any convex clustering algorithm, they are omitted from all timing results shown in this paper.

## F Additional Related Work

Following its original introduction by Pelckmans et al. (2005) and popularization by Hocking et al. (2011) and Lindsten et al. (2011), convex clustering has been the subject of much methodological and theoretical research. In this section, we review some of this related

work which, while not directly relevant to the computational or visualization strategies we propose, may be of interest to readers interested in convex clustering.

The convex clustering problem can be generalized as

$$\hat{\mathbf{U}}_\lambda = \arg \min_{\mathbf{U} \in \mathbb{R}^{n \times p}} \frac{1}{2} \|\mathbf{X} - \mathbf{U}\|_F^2 + \lambda \sum_{\substack{i,j=1 \\ i < j}}^n w_{ij} p(\mathbf{U}_{i\cdot} - \mathbf{U}_{j\cdot})$$

where  $p(\cdot)$  is any sparsity-inducing function. The choice of an  $\ell_q$ -norm ( $p(\cdot) = \|\cdot\|_q$ ) gives standard convex clustering as considered in this paper. Pan et al. (2013), Marchetti and Zhou (2014), Wu et al. (2016), and Shah and Koltun (2017) have all considered the use of non-convex choices of  $p(\cdot)$ , typically using the popular SCAD or MCP penalty functions to reduce bias and improve estimation performance (Fan and Li, 2001; Zhang, 2010). We do not consider non-convex  $p(\cdot)$  in this paper, though the computational techniques and visualizations we propose could be adapted to non-convex penalties in a relatively straightforward manner.

Restricting our attention to standard convex clustering ( $p(\cdot) = \|\cdot\|_q$ ), several useful methodological extensions have been proposed in the literature. For example, Wang et al. (2016) augment the convex clustering problem (1) with an additional sparse component to add robustness to outliers, similar to the robust PCA formulation of Candès et al. (2011), while Wang et al. (2018) propose a variant which incorporates feature selection into the clustering objective using an  $\ell_1$  penalty (Tibshirani, 1996). As discussed in Section C, Chi et al. (2017) extend convex clustering to the bi-clustering setting, where rows and columns are simultaneously clustered. Building on this work, Chi et al. (2018) extend bi-clustering to general co-clustering of  $k$ -order tensors, where they note several surprising theoretical advantages. The recent paper by Park et al. (2018+) extends convex clustering to *histogram-valued* data by replacing the Euclidean distance with an appropriate metric on the space of histograms.

The squared Frobenius loss function of the convex clustering problem may be interpreted as an isotropic Gaussian likelihood, suggesting another avenue for generalization. Sui et al. (2018) replace the Frobenius loss with a squared Mahalanobis distance to improve performance on non-spherical clusters. If the metric (inverse covariance matrix) is known, simple variants on the techniques used in this paper may be used; if the metric must be estimated from the data, the resulting problem is bi-convex and an alternating minimization scheme must be used, only guaranteeing convergence to a stationary point.

The use of a convex formulation allows the sophisticated tools of modern high-dimensional statistics to be brought to bear (Bühlmann and Geer, 2011; Hastie et al., 2015). In addition to the work of Tan and Witten (2015) proving a form prediction consistency and of Radchenko and Mukherjee (2017) proving asymptotic dendrogram recovery, Zhu et al. (2014) give sufficient conditions for exact cluster recovery in the two-cluster case. The results of Zhu et al. (2014) were later extended by Panahi et al. (2017) and by Sun et al. (2018) to the more general multi-cluster case.

In addition to the general purpose operator-splitting algorithms proposed by Chi and Lange (2015), specialized algorithms have been proposed for convex clustering in the “large  $n$ ” (many observations) setting. Panahi et al. (2017) propose a stochastic incremental algorithm based on the framework of Bertsekas (2011), while Sun et al. (2018) propose a semi-smooth Newton algorithm based on the framework of Li et al. (2016). Chen et al. (2015) propose a proximal distance-based algorithm (Lange and Keys, 2014) and provide a GPU-based implementation. Recently, Ho et al. (2019) proposed a generalized dual gradient ascent algorithm with linear convergence, though their approach only works for the  $q = 1$  case; their approach is likely amenable to algorithmic regularization schemes similar to those we have propose for the ADMM.

The special case of convex clustering in  $\mathbb{R}$  has been studied under various names, including *total variation denoising* (Rudin et al., 1992), the *edge lasso* (Sharpnack et al., 2012) and the

*graph-fused lasso* (Hoeffling, 2010), or as a special case of the generalized lasso (Tibshirani and Taylor, 2011). When the underlying graph is a chain graph, convex clustering simplifies to the well-studied *fused lasso* problem (Tibshirani et al., 2005; Rinaldo, 2009; Johnson, 2013).

Convex clustering has not yet seen significant adoption outside of the statistics and machine learning communities, though Chen et al. (2015) discuss applications to human genomics. Nagorski and Allen (2018) propose an alternative weighting scheme based on genetic distances which they use to perform genomic region segmentation.

## Additional References

- Bauschke, Heinz H. and Patrick L. Combettes (2008). “A Dykstra-Like Algorithm for Two Monotone Operators”. *Pacific Journal of Optimization* 4.3, pp. 383–391. URL: <http://www.ybook.co.jp/online2/oppjo/vol4/p383.html>.
- Bertsekas, Dimitri P. (2011). “Incremental proximal methods for large scale convex optimization”. *Mathematical Programming, Series B* 129, pp. 163–195. DOI: [10.1007/s10107-011-0472-0](https://doi.org/10.1007/s10107-011-0472-0).
- Boyd, Stephen, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein (2011). “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers”. *Foundations and Trends<sup>®</sup> in Machine Learning* 3.1, pp. 1–122. DOI: [10.1561/22000000016](https://doi.org/10.1561/22000000016).
- Bühlmann, Peter and Sara van de Geer (2011). *Statistics for High-Dimensional Data: Methods, Theory, and Applications*. Springer Series in Statistics. Springer Verlag. ISBN: 978-3-642-20191-2. DOI: [10.1007/978-3-642-20192-9](https://doi.org/10.1007/978-3-642-20192-9).
- Candès, Emmanuel J., Xiaodong Li, Yi Ma, and John Wright (2011). “Robust Principal Component Analysis?” *Journal of the ACM* 58.3, 11:1–11:37. DOI: [10.1145/1970392.1970395](https://doi.org/10.1145/1970392.1970395).

- Chen, Gary K., Eric C. Chi, John Michael O. Ranola, and Kenneth Lange (2015). “Convex Clustering: An Attractive Alternative to Hierarchical Clustering”. *PLOS Computational Biology* 11.5, e1004228. DOI: [10.1371/journal.pcbi.1004228](https://doi.org/10.1371/journal.pcbi.1004228).
- Chi, Eric C., Genevera I. Allen, and Richard G. Baraniuk (2017). “Convex Biclustering”. *Biometrics* 73.1, pp. 10–19. DOI: [10.1111/biom.12540](https://doi.org/10.1111/biom.12540).
- Chi, Eric C., Brian R. Gaines, Will Wei Sun, Hua Zhou, and Jian Yang (2018). “Provable Convex Co-Clustering of Tensors”. *ArXiv Pre-Print 1803.06518*. URL: <https://arxiv.org/abs/1803.06518>.
- Chi, Eric C. and Kenneth Lange (2015). “Splitting Methods for Convex Clustering”. *Journal of Computational and Graphical Statistics* 24.4, pp. 994–1013. DOI: [10.1080/10618600.2014.948181](https://doi.org/10.1080/10618600.2014.948181).
- Combettes, Patrick L. and Jean-Cristophe Pesquet (2011). “Proximal Splitting Methods in Signal Processing”. *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*. Ed. by Heinz H. Bauschke, Regina S. Burachik, Patrick L. Combettes, Veit Elser, D. Russell Luke, and Henry Wolkowicz. Springer. Chap. 10, pp. 185–212. DOI: [10.1007/978-1-4419-9569-8\\_10](https://doi.org/10.1007/978-1-4419-9569-8_10).
- Deng, Wei and Wotao Yin (2016). “On the Global and Linear Convergence of the Generalized Alternating Direction Method of Multipliers”. *Journal of Scientific Computing* 66.3, pp. 889–916. DOI: [10.1007/s10915-015-0048-x](https://doi.org/10.1007/s10915-015-0048-x).
- Duchi, John, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra (2008). “Efficient Projections onto the  $\ell_1$ -Ball for Learning in High Dimensions”. *ICML 2008: Proceedings of the 25<sup>th</sup> International Conference on Machine Learning*. Ed. by Andrew McCallum and Sam Roweis. Helsinki, Finland: Omnipress, pp. 272–279. DOI: [10.1145/1390156.1390191](https://doi.org/10.1145/1390156.1390191).
- Fan, Jianqing and Runze Li (2001). “Variable Selection via Nonconcave Penalized Likelihood and its Oracle Properties”. *Journal of the American Statistical Association* 96.456, pp. 1348–1360. DOI: [10.1198/016214501753382273](https://doi.org/10.1198/016214501753382273).

- Gabay, Daniel and Bertrand Mercier (1976). “A Dual Algorithm for the Solution of Nonlinear Variational Problems via Finite Element Approximation”. *Computers & Mathematics with Applications* 2.1, pp. 17–40. DOI: [10.1016/0898-1221\(76\)90003-1](https://doi.org/10.1016/0898-1221(76)90003-1).
- Glowinski, R. and A. Marroco (1975). “Sur l’Approximation, par Éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de Dirichlet non linéaires”. *Revue française d’automatique, informatique, recherche opérationnelle. Analyse numérique* 9.R2, pp. 41–76.
- Goldstein, Tom, Brendan O’Donoghue, Simon Setzer, and Richard Baraniuk (2014). “Fast Alternating Direction Optimization Methods”. *SIAM Journal on Imaging Sciences* 7.3, pp. 1588–1623. DOI: [10.1137/120896219](https://doi.org/10.1137/120896219).
- Hastie, Trevor, Robert Tibshirani, and Martin Wainwright (2015). *Statistical Learning with Sparsity: The Lasso and Generalizations*. Monographs on Statistics and Applied Probability 143. Chapman and Hall/CRC. ISBN: 978-1-498-71216-3.
- Ho, Nhat, Tianyi Lin, and Michael I. Jordan (2019). “Global Error Bounds and Linear Convergence for Gradient-Based Algorithms for Trend Filtering and  $\ell_1$ -Convex Clustering”. *ArXiv Pre-Print 1904.07462*. URL: <https://arxiv.org/abs/1904.07462>.
- Hocking, Toby Dylan, Armand Joulin, Francis Bach, and Jean-Philippe Vert (2011). “Clustertpath: An Algorithm for Clustering using Convex Fusion Penalties”. *ICML 2011: Proceedings of the 28<sup>th</sup> International Conference on Machine Learning*. Ed. by Lise Getoor and Tobias Scheffer. Bellevue, Washington, USA: ACM, pp. 745–752. ISBN: 978-1-4503-0619-5. URL: [http://www.icml-2011.org/papers/419\\_icmlpaper.pdf](http://www.icml-2011.org/papers/419_icmlpaper.pdf).
- Hoefling, Holger (2010). “A Path Algorithm for the Fused Lasso Signal Approximator”. *Journal of Computational and Graphical Statistics* 19.4, pp. 984–1006. DOI: [10.1198/jcgs.2010.09208](https://doi.org/10.1198/jcgs.2010.09208).
- Hong, Mingyi and Zhi-Quan Luo (2017). “On the linear convergence of the alternating direction method of multipliers”. *Mathematical Programming, Series A* 162.1-2, pp. 165–1699. DOI: [10.1007/s10107-016-1034-2](https://doi.org/10.1007/s10107-016-1034-2).

- Hubert, Lawrence and Phipps Arabie (1985). “Comparing Partitions”. *Journal of Classification* 2.1, pp. 193–218. DOI: [10.1007/BF01908075](https://doi.org/10.1007/BF01908075).
- Johnson, Nicholas A. (2013). “A Dynamic Programming Algorithm for the Fused Lasso and  $L_0$ -Segmentation”. *Journal of Computational and Graphical Statistics* 22.2, pp. 246–260. DOI: [10.1080/10618600.2012.681238](https://doi.org/10.1080/10618600.2012.681238).
- Lange, Kenneth and Kevin L. Keys (2014). “The Proximal Distance Algorithm”. *ICM 2014: Proceedings of the International Congress of Mathematicians*. Ed. by Sun Young Jang, Young Rock Kim, Dae-Woong Lee, and Ikkwon Yie. Vol. 4. Seoul, Korea, pp. 95–116. URL: <http://www.icm2014.org/en/vod/proceedings.html>.
- Li, Xudong, Defeng Sun, and Kim-Chuan Toh (2016). “A Highly Efficient Semismooth Newton Augmented Lagrangian Method for Solving Lasso Problems”. *SIAM Journal on Optimization* 28.1, pp. 433–458. DOI: [10.1137/16M1097572](https://doi.org/10.1137/16M1097572).
- Lindsten, Fredrik, Henrik Ohlsson, and Lennart Ljung (2011). “Clustering using sum-of-norms regularization: With application to particle filter output computation”. *SSP 2011: Proceedings of the 2011 IEEE Statistical Signal Processing Workshop*. Ed. by Petar M. Djuric. Nice, France: Curran Associates, Inc., pp. 201–204. DOI: [10.1109/SSP.2011.5967659](https://doi.org/10.1109/SSP.2011.5967659).
- Lions, P.L. and B. Mercier (1979). “Splitting Algorithms for the Sum of Two Nonlinear Operators”. *SIAM Journal on Numerical Analysis* 16.6, pp. 964–979. DOI: [10.1137/0716071](https://doi.org/10.1137/0716071).
- Marchetti, Yuliya and Qing Zhou (2014). “Solution Path Clustering with Adaptive Concave Penalty”. *Electronic Journal of Statistics* 8.1, pp. 1569–1603. DOI: [10.1214/14-EJS934](https://doi.org/10.1214/14-EJS934).
- Moreau, Jean Jacques (1962). “Décomposition Orthogonale d’un Espace Hilbertien Selon Deux Cônes Mutuellement Polaires”. *Comptes Rendus Hebdomadaires des Séances de l’Académie des Sciences* 255.1, pp. 238–240.
- Nagorski, John and Genevera I. Allen (2018). “Genomic Region Detection via Spatial Convex Clustering”. *PLoS One* 13.9, e0203007. DOI: [10.1371/journal.pone.0203007](https://doi.org/10.1371/journal.pone.0203007).

- Nishihara, Robert, Laurent Lessard, Ben Recht, Andrew Packard, and Michael Jordan (2015). “A General Analysis of the Convergence of ADMM”. *ICML:2015: Proceedings of the 32<sup>nd</sup> International Conference on Machine Learning*. Ed. by Francis Bach and David Blei. Lille, France: PMLR, pp. 343–352. URL: <http://proceedings.mlr.press/v37/nishihara15.html>.
- Pan, Wei, Xiaotong Shen, and Binghui Liu (2013). “Cluster Analysis: Unsupervised Learning via Supervised Learning with a Non-convex Penalty”. *Journal of Machine Learning Research* 14, pp. 1865–1889. URL: <http://www.jmlr.org/papers/v14/pan13a.html>.
- Panahi, Ashkan, Devdatt Dubhashi, Fredrik D. Johansson, and Chiranjib Bhattacharyya (2017). “Clustering by Sum of Norms: Stochastic Incremental Algorithm, Convergence, and Cluster Recovery”. *ICML:2017: Proceedings of the 34<sup>th</sup> International Conference on Machine Learning*. Ed. by Doina Precup and Yee Whye Teh. Vol. 70. Sydney, Australia: PMLR, pp. 2769–2777. URL: <http://proceedings.mlr.press/v70/panahi17a.html>.
- Park, Cheolwoo, Hosik Choi, Chris Delcher, Yanning Wang, and Young Joo Yoon (2018+). “Convex Clustering Analysis for Histogram-Valued Data”. *Biometrics* To appear. DOI: [10.1111/biom.13004](https://doi.org/10.1111/biom.13004).
- Pelckmans, Kristiaan, Joseph de Brabanter, Bart de Moor, and Johan Suykens (2005). “Convex Clustering Shrinkage”. *PASCAL Workshop on Statistics and Optimization of Clustering*.
- Radchenko, Peter and Gourab Mukherjee (2017). “Convex Clustering via  $\ell_1$  Fusion Penalization”. *Journal of the Royal Statistical Society, Series B: Statistical Methodology* 79.5, pp. 1527–1546. DOI: [10.1111/rssb.12226](https://doi.org/10.1111/rssb.12226).
- Rand, William M. (1971). “Objective Criteria for the Evaluation of Clustering Methods”. *Journal of the American Statistical Association* 66.366, pp. 846–850. DOI: [10.1080/01621459.1971.10482356](https://doi.org/10.1080/01621459.1971.10482356).
- Rinaldo, Alessandro (2009). “Properties and Refinements of the Fused Lasso”. *Annals of Statistics* 37.5B, pp. 2922–2952. DOI: [10.1214/08-AOS665](https://doi.org/10.1214/08-AOS665).

- Rockafellar, R. Tyrrell (1970). *Convex Analysis*. Princeton University Press. ISBN: 978-0-691-01586-6.
- Rosset, Saharon and Ji Zhu (2007). “Piecewise Linear Regularized Solution Paths”. *Annals of Statistics* 35.3, pp. 1012–1030. DOI: [10.1214/009053606000001370](https://doi.org/10.1214/009053606000001370).
- Rudin, Leonid I., Stanley Osher, and Emad Fatemi (1992). “Nonlinear Total Variation Based Noise Removal Algorithms”. *Physica D: Nonlinear Phenomena* 60.1-4, pp. 259–268. DOI: [10.1016/0167-2789\(92\)90242-F](https://doi.org/10.1016/0167-2789(92)90242-F).
- Shah, Sohil Atul and Vladlen Koltun (2017). “Robust continuous clustering”. *Proceedings of the National Academy of Sciences of the United States* 114.37, pp. 9814–9819. DOI: [10.1073/pnas.1700770114](https://doi.org/10.1073/pnas.1700770114).
- Sharpnack, James, Aarti Singh, and Alessandro Rinaldo (2012). “Sparsistency of the Edge Lasso over Graphs”. *AISTATS 2012: Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*. Ed. by Neil D. Lawrence and Mark Girolami. Vol. 22. La Palma, Canary Islands: PMLR, pp. 1028–1036. URL: <http://proceedings.mlr.press/v22/sharpnack12.html>.
- Shi, Wei, Qing Ling, Kun Yuan, Gang Wu, and Wotao Yin (2014). “On the Linear Convergence of the ADMM in Decentralized Consensus Optimization”. *IEEE Transactions on Signal Processing* 62.7, pp. 1750–1761. DOI: [10.1109/TSP.2014.2304432](https://doi.org/10.1109/TSP.2014.2304432).
- Sui, Xiaopeng Lucia, Li Xu, Xiaoning Qian, and Tie Liu (2018). “Convex Clustering with Metric Learning”. *Pattern Recognition* 81, pp. 575–584. DOI: [10.1016/j.patcog.2018.04.019](https://doi.org/10.1016/j.patcog.2018.04.019).
- Sun, Defeng, Kim-Chuan Toh, and Yancheng Yuan (2018). “Convex Clustering: Model, Theoretical Guarantee and Efficient Algorithm”. *ArXiv Pre-Print 1810.02677*. URL: <https://arxiv.org/abs/1810.02677>.
- Tan, Kean Ming and Daniela Witten (2015). “Statistical Properties of Convex Clustering”. *Electronic Journal of Statistics* 9.2, pp. 2324–2347. DOI: [10.1214/15-EJS1074](https://doi.org/10.1214/15-EJS1074).

- Tibshirani, Robert (1996). “Regression Shrinkage and Selection via the Lasso”. *Journal of the Royal Statistical Society, Series B (Methodological)* 58.1, pp. 267–288. DOI: [10.1111/j.2517-6161.1996.tb02080.x](https://doi.org/10.1111/j.2517-6161.1996.tb02080.x).
- Tibshirani, Robert, Michael Saunders, Saharon Rosset, Ji Zhu, and Keith Knight (2005). “Sparsity and Smoothness via the Fused Lasso”. *Journal of the Royal Statistical Society, Series B: Statistical Methodology* 67.1, pp. 91–108. DOI: [10.1111/j.1467-9868.2005.00490.x](https://doi.org/10.1111/j.1467-9868.2005.00490.x).
- Tibshirani, Ryan J. and Jonathan Taylor (2011). “The Solution Path of the Generalized Lasso”. *Annals of Statistics* 39.3, pp. 1335–1371. DOI: [10.1214/11-AOS878](https://doi.org/10.1214/11-AOS878).
- Tseng, Paul (1991). “Applications of a Splitting Algorithm to Decomposition in Convex Programming and Variational Inequalities”. *SIAM Journal on Control and Optimization* 29.1, pp. 119–138. DOI: [10.1137/0329006](https://doi.org/10.1137/0329006).
- Wang, Binhuan, Yilong Zhang, Will Wei Sun, and Yixin Fang (2018). “Sparse Convex Clustering”. *Journal of Computational and Graphical Statistics* 27.2, pp. 393–403. DOI: [10.1080/10618600.2017.1377081](https://doi.org/10.1080/10618600.2017.1377081).
- Wang, Qi, Pinghua Gong, Shiyu Chang, Thomas S. Huang, and Jiayu Zhou (2016). “Robust Convex Clustering Analysis”. *ICDM 2016: Proceedings of the 16<sup>th</sup> IEEE International Conference on Data Mining*. Ed. by Francesco Bonchi, Josep Domingo-Ferrer, Ricardo Baeza-Yates, Zhi-Hua Zhou, and Xindong Wu. Barcelona, Spain, pp. 1263–1268. DOI: [10.1109/ICDM.2016.0170](https://doi.org/10.1109/ICDM.2016.0170).
- Wilkinson, Leland and Michael Friendly (2009). “The History of the Cluster Heat Map”. *The American Statistician* 63.2, pp. 179–184. DOI: [10.1198/tas.2009.0033](https://doi.org/10.1198/tas.2009.0033).
- Wu, Chong, Sunghoon Kwon, Xiaotong Shen, and Wei Pan (2016). “A New Algorithm and Theory for Penalized Regression-based Clustering”. *Journal of Machine Learning Research* 17.188, pp. 1–25. URL: <http://jmlr.org/papers/v17/15-553.html>.

- Yang, Wei Hong and Deren Han (2016). “Linear Convergence of the Alternating Direction Method of Multipliers for a Class of Convex Optimization Problems”. *SIAM Journal on Numerical Analysis* 54.2, pp. 625–640. DOI: [10.1137/140974237](https://doi.org/10.1137/140974237).
- Zhang, Cun-Hui (2010). “Nearly Unbiased Variable Selection under Minimax Concave Penalty”. *Annals of Statistics* 38.2, pp. 894–942. DOI: [10.1214/09-AOS729](https://doi.org/10.1214/09-AOS729).
- Zhu, Changbo, Huan Xu, Chenlei Leng, and Shuicheng Yan (2014). “Convex Optimization Procedure for Clustering: Theoretical Revisit”. *NIPS 2014: Advances in Neural Information Processing Systems 27*. Ed. by Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Killian Q. Weinberger. Montréal, Canada: Curran Associates, Inc., pp. 1619–1627. URL: <https://papers.nips.cc/paper/5307-convex-optimization-procedure-for-clustering-theoretical-revisit>.