

# Learning to Cluster Using High Order Graphical Models with Latent Variables (Supplemental Material)

Nikos Komodakis  
University of Crete, Computer Science Department  
<http://www.csd.uoc.gr/~komod>

## Abstract

This document provides technical proofs for all theorems in the main paper.

## 1. Proofs

**Lemma 1** Let  $\hat{\mathbf{x}}^{k,p}, \hat{\mathbf{x}}^{k,C}$  be binary minimizers of the energy functions  $\bar{E}_p^k, \bar{E}_C^k$ . Define  $f_{pq}^k \equiv f_{pq}(\mathbf{y}^k)$ ,  $\hat{X}_q^k \equiv \hat{x}_{qq}^{k,C} + \sum_p \hat{x}_{pq}^{k,p}, \forall q \in C$ . Update (30) then reduces to

$$\begin{bmatrix} \mathbf{w} \\ \lambda_{pq}^k \\ \lambda_{Cq}^k \end{bmatrix} = s_t \begin{bmatrix} \tau \nabla J(\mathbf{w}) + \sum_k \delta_{\mathbf{w}}^k \\ \frac{\hat{X}_q^k}{|S^k|+1} - \hat{x}_{qq}^{k,p} \\ \frac{\hat{X}_q^k}{|S^k|+1} - \hat{x}_{qq}^{k,C} \end{bmatrix}, \quad (36)$$

where  $\delta_{\mathbf{w}}^k = \sum_{p,q} x_{pq}^k f_{pq}^k - \sum_{p \neq q} \hat{x}_{pq}^{k,p} f_{pq}^k - \frac{\sum_q \hat{X}_q^k f_{qq}^k}{|S^k|+1}$ .

**Note:** If  $J(\mathbf{w})$  is non-differentiable (e.g., if  $J(\mathbf{w}) = \|\mathbf{w}\|_1$ ) then  $\nabla J(\mathbf{w})$  should refer to a subgradient of  $J(\cdot)$  at  $\mathbf{w}$ .

*Proof.* Update (30) requires computing a subgradient of the objective function (28) with respect to  $\mathbf{w}, \lambda^k$  (for a fixed  $\mathbf{x}^k$ ). To this end, we need to compute the corresponding subgradient for each of the terms  $\bar{\mathcal{L}}_{\bar{E}_p^k}(\mathbf{x}^k; \mathbf{w}, \lambda^k)$  and  $\bar{\mathcal{L}}_{\bar{E}_C^k}(\mathbf{x}^k; \mathbf{w}, \lambda^k)$  that are included in function (28). By definition (21) it holds that<sup>1</sup>

$$\bar{\mathcal{L}}_{\bar{E}_p^k}(\mathbf{x}^k; \mathbf{w}, \lambda^k) = \bar{E}_p^k(\mathbf{x}^k; \mathbf{w}, \lambda^k) - \min_{\mathbf{x}} \bar{E}_p^k(\mathbf{x}; \mathbf{w}, \lambda^k) \quad (37)$$

$$= \bar{E}_p^k(\mathbf{x}^k; \mathbf{w}, \lambda^k) + \max_{\mathbf{x}} (-\bar{E}_p^k(\mathbf{x}; \mathbf{w}, \lambda^k)) \quad (38)$$

A subgradient for a pointwise maximum function  $g(\mathbf{w}, \lambda^k) = \max_{\mathbf{x}} g_{\mathbf{x}}(\mathbf{w}, \lambda^k)$ , where each  $g_{\mathbf{x}}(\cdot, \cdot)$  is convex and differentiable, is given by  $\nabla g_{\hat{\mathbf{x}}}(\mathbf{w}, \lambda^k)$  for any  $\hat{\mathbf{x}}$  that satisfies  $g(\mathbf{w}, \lambda^k) = g_{\hat{\mathbf{x}}}(\mathbf{w}, \lambda^k)$ , i.e.,  $\max_{\mathbf{x}} g_{\mathbf{x}}(\mathbf{w}, \lambda^k) = g_{\hat{\mathbf{x}}}(\mathbf{w}, \lambda^k)$ . Since function  $-\bar{E}_p^k(\mathbf{x}; \mathbf{w}, \lambda^k)$  is linear (and hence both convex and differentiable) with respect to  $\mathbf{w}, \lambda^k$ , a subgradient of function  $\bar{\mathcal{L}}_{\bar{E}_p^k}(\mathbf{x}^k; \mathbf{w}, \lambda^k)$  (with respect to  $\mathbf{w}, \lambda^k$ ) will thus equal

$$\nabla \bar{E}_p^k(\mathbf{x}^k; \mathbf{w}, \lambda^k) - \nabla \bar{E}_p^k(\hat{\mathbf{x}}^{k,p}; \mathbf{w}, \lambda^k), \quad (39)$$

where  $\hat{\mathbf{x}}^{k,p}$  denotes a binary minimizer of function  $\bar{E}_p^k(\cdot; \mathbf{w}, \lambda^k)$ . Therefore, based on (39) and the fact that  $d_{pq}^k = \mathbf{w}^T f_{pq}^k$ , a

<sup>1</sup>Note that both here and in the main paper all vectors of CRF variables  $\mathbf{x}$  are *always* assumed to be *integral*. Therefore, in order to reduce notational clutter we often omit stating this integrality constraint when using such vectors (e.g., we simply write  $\min_{\mathbf{x}}$  instead of  $\min_{\{\mathbf{x}: \mathbf{x} \text{ has integral components}\}}$ ).

subgradient of  $\bar{\mathcal{L}}_{\bar{E}_p^k}$  will have components  $\delta \mathbf{w}^{k,p}$ ,  $\{\delta \lambda_q^{k,p}\}_q$  (corresponding to variables  $\mathbf{w}$ ,  $\{\lambda_{pq}^k\}_q$  respectively) given by

$$\delta \mathbf{w}^{k,p} = \sum_{q:q \neq p} x_{pq}^k f_{pq}^k + \sum_q \frac{x_{qq}^k f_{qq}^k}{|S^k|+1} - \left( \sum_{q:q \neq p} \hat{x}_{pq}^{k,p} f_{pq}^k + \sum_q \frac{\hat{x}_{qq}^{k,p} f_{qq}^k}{|S^k|+1} \right) \quad (40)$$

$$\delta \lambda_q^{k,p} = x_{qq}^k - \hat{x}_{qq}^{k,p} . \quad (41)$$

Similarly, we can prove that a subgradient of function  $\bar{\mathcal{L}}_{\bar{E}_C^k}(\mathbf{x}^k; \mathbf{w}, \boldsymbol{\lambda}^k)$  will have components  $\delta \mathbf{w}^{k,C}$ ,  $\{\delta \lambda_q^{k,C}\}_{q \in C}$  (corresponding to variables  $\mathbf{w}$ ,  $\{\lambda_{Cq}^k\}_{q \in C}$  respectively) given by

$$\delta \mathbf{w}^{k,C} = \sum_{q \in C} \frac{x_{qq}^k f_{qq}^k}{|S^k|+1} - \sum_{q \in C} \frac{\hat{x}_{qq}^{k,C} f_{qq}^k}{|S^k|+1} \quad (42)$$

$$\delta \lambda_q^{k,C} = x_{qq}^k - \hat{x}_{qq}^{k,C} , \forall q \in C \quad (43)$$

where  $\hat{\mathbf{x}}^{k,C}$  denotes a binary minimizer of function  $\bar{E}_C^k(\cdot; \mathbf{w}, \boldsymbol{\lambda}^k)$ .

Therefore, a total subgradient of the objective function (28) will have components  $\delta \mathbf{w}$ ,  $\delta \lambda_q^{k,p}$ ,  $\delta \lambda_q^{k,C}$  (corresponding to variables  $\mathbf{w}$ ,  $\lambda_{pq}^k$ ,  $\lambda_{Cq}^k$  respectively), where

$$\delta \mathbf{w} = \tau \nabla J(\mathbf{w}) + \sum_k \left( \sum_{p \in S^k} \delta \mathbf{w}^{k,p} + \sum_{C \in \mathcal{C}^k} \delta \mathbf{w}^{k,C} \right) \stackrel{(40),(42)}{=} \tau \nabla J(\mathbf{w}) + \sum_k \delta \mathbf{w}^k . \quad (44)$$

Furthermore, projection onto the set  $\boldsymbol{\Lambda}^k = \{\boldsymbol{\lambda}^k : \sum_{p \in S^k} \lambda_{pq}^k + \lambda_{Cq}^k = 0, \forall C \in \mathcal{C}^k, q \in C\}$  simply requires to first subtract the average  $\frac{\sum_{p \in S^k} \delta \lambda_q^{k,p} + \delta \lambda_q^{k,C}}{|S^k|+1} \stackrel{(41),(43)}{=} x_{qq}^k - \frac{X_q^k}{|S^k|+1}$  from each of the elements  $\{\delta \lambda_q^{k,p}\}_p$ ,  $\delta \lambda_q^{k,C}$  before applying the updates  $\mathbf{w} \leftarrow s_t \delta \mathbf{w}$ ,  $\lambda_{pq}^k \leftarrow s_t \delta \lambda_q^{k,p}$ ,  $\lambda_{Cq}^k \leftarrow s_t \delta \lambda_q^{k,C}$  (where  $s_t$  is the multiplier used during the  $t$ -th iteration). This is easily seen to lead to updates (36), which concludes the proof of the lemma.  $\square$

**Lemma 2** Let  $[a]_+ \equiv \max(a, 0)$ ,  $[a]_- \equiv \min(a, 0)$ .

1. For fixed  $p$ , let  $\theta_q^k \equiv \frac{\bar{u}_{qq}^k(1)}{|S^k|+1} + \lambda_{pq}^k, \forall q$  and let us define  $\bar{\theta}_q^k \equiv \bar{u}_{pq}^k(1) + [\theta_q^k]_+, \forall q \neq p$  and  $\bar{\theta}_p^k = \theta_p^k$ . A minimizer  $\hat{\mathbf{x}}$  of  $\bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  can be computed as follows:

$$\forall q \neq p, \hat{x}_{qq} \leftarrow [\theta_q^k < 0] \quad (45)$$

$$\forall q, \hat{x}_{pq} \leftarrow [q = \bar{q}], \text{ where } \bar{q} = \arg \min_q \bar{\theta}_q^k \quad (46)$$

2. For fixed  $C \in \mathcal{C}^k$ , let  $\theta_q^k \equiv \frac{\bar{u}_{qq}^k(1)}{|S^k|+1} + \lambda_{Cq}^k, \forall q \in C$ . A minimizer  $\hat{\mathbf{x}}$  of  $\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  is given by

$$\forall q \in C, \hat{x}_{qq} = \begin{cases} [\theta_q^k < \alpha], & \text{if } 2\alpha + \sum_{q' \in C} [\theta_{q'}^k - \alpha]_- < 0 \\ 0, & \text{otherwise} \end{cases} \quad (47)$$

*Proof.* 1. It holds that

$$\bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) = \sum_{q:q \neq p} \bar{u}_{pq}^k(x_{pq}) + \sum_q \left( \frac{\bar{u}_{qq}^k(x_{qq})}{|S^k|+1} + \lambda_{pq}^k x_{qq} \right) + \sum_q \bar{\phi}_{pq}(x_{pq}, x_{qq}) + \bar{\phi}_p(\mathbf{x}_p) - \beta \quad (48)$$

$$= \sum_{q:q \neq p} \bar{u}_{pq}^k(1) x_{pq} + \sum_q \theta_q^k x_{qq} + \sum_q \bar{\phi}_{pq}(x_{pq}, x_{qq}) + \bar{\phi}_p(\mathbf{x}_p) - \beta \quad (49)$$

$$= \sum_{q:q \neq p} \bar{u}_{pq}^k(1) x_{pq} + \sum_q (\theta_q^k x_{qq} + \bar{\phi}_{pq}(x_{pq}, x_{qq})) + \bar{\phi}_p(\mathbf{x}_p) - \beta , \quad (50)$$

where  $\bar{\phi}_{pq}(x_{pq}, x_{qq}) = \delta(x_{pq} \leq x_{qq})$ ,  $\bar{\phi}_p(\mathbf{x}_p) = \delta\left(\sum_q x_{pq} = 1\right)$  and  $\delta(\cdot)$  equals 0 if the expression in parenthesis is satisfied and  $\infty$  otherwise.

Due to the term  $\theta_q^k x_{qq}$ , it is easy to see that if we set  $x_{qq} = 1$  for any  $q \neq p$  then the value of the function  $\bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  will decrease if and only if it holds  $\theta_q^k < 0$ . Therefore, to minimize  $\bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  we must set

$$\hat{x}_{qq} = [\theta_q^k < 0], \forall q \neq p. \quad (51)$$

Furthermore, the fact that the components of an optimal solution  $\hat{\mathbf{x}}$  must belong to  $\{0, 1\}$  in conjunction with the form of the potential  $\bar{\phi}_p(\mathbf{x}_p) = \delta\left(\sum_q x_{pq} = 1\right)$  impose the constraint that we must set equal to 1 exactly one of the variables in the set  $\{\hat{x}_{pq}\}_q$ . If we set variable  $\hat{x}_{pq}$  (with  $q \neq p$ ) equal to 1 then the cost we must pay is  $\bar{u}_{pq}^k(1)$ , due to the term  $\bar{u}_{pq}^k(1)\hat{x}_{pq}$ , plus  $[\theta_q^k]_+$ , due to the term  $\theta_q^k \hat{x}_{qq} + \bar{\phi}_{pq}(\hat{x}_{pq}, \hat{x}_{qq})$  that requires also setting  $\hat{x}_{qq} = 1$  (note that we are paying  $[\theta_q^k]_+$  and not  $\theta_q^k$  because for  $q \neq p$  if  $\theta_q^k < 0$  then  $\hat{x}_{qq}$  is set to 1 anyway due to (51) and thus no extra cost is paid in this case). On the other hand, if we set  $\hat{x}_{pp} = 1$  then the cost we must pay is  $\theta_p^k$  due to the term  $\theta_p^k \hat{x}_{pp}$ . Therefore, for any  $q$ , the cost we pay if we choose to set  $\hat{x}_{pq} = 1$  is given by  $\bar{\theta}_q^k$ . As a result, we should set  $\hat{x}_{pq} = [q = \bar{q}]$ , where  $\bar{q} = \arg \min_q \bar{\theta}_q^k$ .

2. Energy  $\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  can be expressed as

$$\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) = \sum_{q \in C} \left( \frac{\bar{u}_{qq}^k(x_{qq})}{|S^k| + 1} + \lambda_{Cq}^k x_{qq} \right) + \bar{\phi}_C(\mathbf{x}_C) \quad (52)$$

$$= \sum_{q \in C} \theta_q^k x_{qq} + \bar{\phi}_C(\mathbf{x}_C) \quad (53)$$

$$= \sum_{q \in C} \theta_q^k x_{qq} - \alpha \left| 1 - \sum_{q \in C} x_{qq} \right|. \quad (54)$$

We will consider two cases:

(a) The minimizer of function  $\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  is given by  $\hat{\mathbf{x}} = \mathbf{0}$  (i.e., none of the binary variables  $\{\hat{x}_{qq}\}_{q \in C}$  is equal to 1). In this case the minimum of function  $\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  must equal

$$\text{OPT}_1 = -\alpha. \quad (55)$$

(b) The minimizer of function  $\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  is given by  $\hat{\mathbf{x}} \neq \mathbf{0}$ . In this case at least one of the binary variables  $\{\hat{x}_{qq}\}_{q \in C}$  will equal 1 and so  $\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  can be written as

$$\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) = \sum_{q \in C} \theta_q^k x_{qq} - \alpha \left| 1 - \sum_{q \in C} x_{qq} \right| \quad (56)$$

$$= \sum_{q \in C} \theta_q^k x_{qq} - \alpha \left( \sum_{q \in C} x_{qq} - 1 \right) \quad (57)$$

$$= \sum_{q \in C} (\theta_q^k - \alpha) x_{qq} + \alpha. \quad (58)$$

Therefore, the minimizer  $\hat{\mathbf{x}}$  will be given by

$$\hat{x}_{qq} = [\theta_q^k < \alpha] \quad (59)$$

and so the optimum value of  $\bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k)$  will equal

$$\text{OPT}_2 = \sum_{q \in C} [\theta_q^k - \alpha]_- + \alpha. \quad (60)$$

To conclude the proof, it suffices to notice that the second case will hold true if and only if

$$\text{OPT}_2 < \text{OPT}_1 \Leftrightarrow \sum_{q \in C} [\theta_q^k - \alpha]_- + \alpha < -\alpha \Leftrightarrow \sum_{q \in C} [\theta_q^k - \alpha]_- + 2\alpha < 0. \quad (61)$$

□

**Lemma 3: Minimizations (27) and (28) in the main paper are equivalent.**

*Proof.* It holds that

$$\min_{\{\mathbf{x}^k \in \mathcal{X}(C^k)\}, \mathbf{w}} \tau J(\mathbf{w}) + \sum_k (\bar{E}^k(\mathbf{x}^k; \mathbf{w}) - \mathcal{R}^k(\mathbf{w})) \quad (62)$$

$$\stackrel{(26)}{=} \min_{\{\mathbf{x}^k \in \mathcal{X}(C^k)\}, \mathbf{w}} \tau J(\mathbf{w}) + \sum_k \left( \bar{E}^k(\mathbf{x}^k; \mathbf{w}) - \max_{\boldsymbol{\lambda}^k \in \Lambda^k} \left( \sum_p \min_{\mathbf{x}} \bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) + \sum_C \min_{\mathbf{x}} \bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) \right) \right) \quad (63)$$

$$= \min_{\{\mathbf{x}^k \in \mathcal{X}(C^k)\}, \mathbf{w}, \{\boldsymbol{\lambda}^k \in \Lambda^k\}} \tau J(\mathbf{w}) + \sum_k \left( \bar{E}^k(\mathbf{x}^k; \mathbf{w}) - \sum_p \min_{\mathbf{x}} \bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) - \sum_C \min_{\mathbf{x}} \bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) \right) \quad (64)$$

$$\stackrel{(25)}{=} \min_{\{\mathbf{x}^k \in \mathcal{X}(C^k)\}, \mathbf{w}, \{\boldsymbol{\lambda}^k \in \Lambda^k\}} \tau J(\mathbf{w}) + \sum_k \left( \sum_p \bar{E}_p^k(\mathbf{x}^k; \mathbf{w}, \boldsymbol{\lambda}^k) + \sum_C \bar{E}_C^k(\mathbf{x}^k; \mathbf{w}, \boldsymbol{\lambda}^k) - \sum_p \min_{\mathbf{x}} \bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) - \sum_C \min_{\mathbf{x}} \bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) \right) \quad (65)$$

$$= \min_{\{\mathbf{x}^k \in \mathcal{X}(C^k)\}, \mathbf{w}, \{\boldsymbol{\lambda}^k \in \Lambda^k\}} \tau J(\mathbf{w}) + \sum_k \sum_p \left( \bar{E}_p^k(\mathbf{x}^k; \mathbf{w}, \boldsymbol{\lambda}^k) - \min_{\mathbf{x}} \bar{E}_p^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) \right) + \sum_k \sum_C \left( \bar{E}_C^k(\mathbf{x}^k; \mathbf{w}, \boldsymbol{\lambda}^k) - \min_{\mathbf{x}} \bar{E}_C^k(\mathbf{x}; \mathbf{w}, \boldsymbol{\lambda}^k) \right) \quad (66)$$

$$= \min_{\{\mathbf{x}^k \in \mathcal{X}(C^k)\}, \mathbf{w}, \{\boldsymbol{\lambda}^k \in \Lambda^k\}} \tau J(\mathbf{w}) + \sum_k \sum_p \bar{\mathcal{L}}_{\bar{E}_p^k}(\mathbf{x}^k; \mathbf{w}, \boldsymbol{\lambda}^k) + \sum_k \sum_C \bar{\mathcal{L}}_{\bar{E}_C^k}(\mathbf{x}^k; \mathbf{w}, \boldsymbol{\lambda}^k), \quad (67)$$

which concludes the proof. □