

Structured Prediction of Large Output Domain via Dual Decomposition and Factorwise Oracles

Ian En-Hsu Yen (Enxu Yan)

Tuesday 15th November, 2016

Abstract

Many applications of machine learning involve structured outputs with large domains, where learning of a structured predictor is prohibitive due to repetitive calls to an expensive inference oracle. In this work, we show that by decomposing training of a Structural Support Vector Machine (SVM) into a series of multiclass SVM problems connected through messages, one can replace an expensive structured oracle with Factorwise Maximization Oracles (FMOs) that allow efficient implementation of complexity sublinear to the factor domain. A Greedy Direction Method of Multiplier (GDMM) algorithm is then proposed to exploit the sparsity of messages while guarantees convergence to ϵ sub-optimality after $O(\log(1/\epsilon))$ passes of FMOs over every factor. We conduct experiments on chain-structured and fully-connected problems of large output domains, where the proposed approach is orders-of-magnitude faster than current state-of-the-art algorithms for training Structural SVMs.

Keywords: Structural SVM, Dual Decomposition, Structured Prediction, Frank-Wolfe Method, Greedy Coordinate Descent, Augmented Lagrangian.

DAP Committee members:

Pradeep Ravikumar (Machine Learning Department);
Matt Gormley (Machine Learning Department).

1 Introduction

Structured prediction has become prevalent with wide applications in Natural Language Processing (NLP), Computer Vision, and Bioinformatics to name a few, where one is interested in outputs of strong interdependence. Although many dependency structures yield intractable inference problems, approximation techniques such as convex relaxations with theoretical guarantees [7, 10, 14] have been developed. However, solving the relaxed problems (LP, QP, SDP, etc.) is computationally expensive for factor graphs of large output domain and results in prohibitive training time when embedded into an learning algorithm relying on inference oracles [6, 9]. For instance, many applications in NLP such as Machine Translation [3], Speech Recognition [21], and Semantic Parsing [1] have output domains as large as the size of vocabulary, for which the prediction of even a single sentence takes considerable time.

One approach to avoid inference during training is by introducing a loss function conditioned on the given labels of neighboring output variables [15]. However, it also introduces more variance to the estimation of model and could degrade testing performance significantly. Another thread of research aims to formulate parameter learning and output inference as a joint optimization problem that avoids treating inference as a subroutine [11, 12]. In this approach, the structured hinge loss is reformulated via dual decomposition, so both messages between factors and model parameters are treated as first-class variables. The new formulation, however, does not yield computational advantage due to the

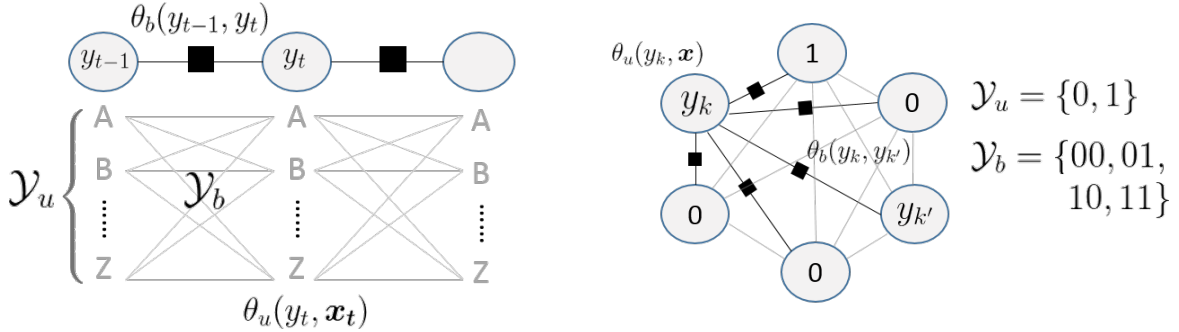


Figure 1: (left) Factors with large output domains in Sequence Labeling. (right) Large number of factors in a Correlated Multilabel Prediction problem. Circles denote variables and black boxes denote factors. (\mathcal{Y}_u : domain of unigram factor. \mathcal{Y}_b : domain of bigram factor.)

constraints entangling the two types of variables. In particular, [11] employs a hybrid method (DLPW) that alternately optimizes model parameters and messages, but the algorithm is not significantly faster than directly performing stochastic gradient on the structured hinge loss. More recently, [12] proposes an approximate objective for structural SVMs that leads to an algorithm considerably faster than DLPW on problems requiring expensive inference. However, the approximate objective requires a trade-off between efficiency and approximation quality, yielding an $O(1/\epsilon^2)$ overall iteration complexity for achieving ϵ sub-optimality.

The contribution of this work is twofold. First, we propose a Greedy Direction Method of Multiplier (GDMM) algorithm that decomposes the training of a structural SVM into factorwise multiclass SVMs connected through sparse messages confined to the active labels. The algorithm guarantees an $O(\log(1/\epsilon))$ iteration complexity for achieving an ϵ sub-optimality and each iteration requires only one pass of *Factorwise Maximization Oracles (FMOs)* over every factor. Second, we show that the FMO can be realized in time sublinear to the cardinality of factor domains, hence is considerably more efficient than a structured maximization oracle when it comes to large output domain. For problems consisting of numerous binary variables, we further give realization of a joint FMO that has complexity sublinear to the number of factors. We conduct experiments on both chain-structured problems that allow exact inference and fully-connected problems that rely on Linear Program relaxations, where we show the proposed approach is orders-of-magnitude faster than current state-of-the-art training algorithms for Structured SVMs.

2 Problem Formulation

Structured prediction aims to predict a set of outputs $\mathbf{y} \in \mathcal{Y}(\mathbf{x})$ from their interdependency and inputs $\mathbf{x} \in \mathcal{X}$. Given a feature map $\phi(\mathbf{x}, \mathbf{y}) : \mathcal{X} \times \mathcal{Y}(\mathbf{x}) \rightarrow \mathbb{R}^d$ that extracts relevant information from (\mathbf{x}, \mathbf{y}) , a linear classifier with parameters \mathbf{w} can be defined as $h(\mathbf{x}; \mathbf{w}) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) \rangle$, where we estimate the parameters \mathbf{w} from a training set $\mathcal{D} = \{(\mathbf{x}_i, \bar{\mathbf{y}}_i)\}_{i=1}^n$ by solving a regularized *Empirical Risk Minimization (ERM)* problem

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n L(\mathbf{w}; \mathbf{x}_i, \bar{\mathbf{y}}_i). \quad (1)$$

In case of a Structural SVM [19, 20], we consider the structured hinge loss

$$L(\mathbf{w}; \mathbf{x}, \bar{\mathbf{y}}) = \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) - \phi(\mathbf{x}, \bar{\mathbf{y}}) \rangle + \delta(\mathbf{y}, \bar{\mathbf{y}}), \quad (2)$$

where $\delta(\mathbf{y}, \bar{\mathbf{y}}_i)$ is a task-dependent error function, for which the Hamming distance $\delta_H(\mathbf{y}, \bar{\mathbf{y}}_i)$ is commonly used. Since the size of domain $|\mathcal{Y}(\mathbf{x})|$ typically grows exponentially with the number of output variables, the tractability of problem (1) lies in the decomposition of the responses $\langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) \rangle$ into several factors, each involving only a few outputs. The factor decomposition can be represented as a bipartite graph $G(\mathcal{F}, \mathcal{V}, \mathcal{E})$ between factors \mathcal{F} and variables \mathcal{V} , where an edge $(f, j) \in \mathcal{E}$ exists if the factor f involves the variable j . Typically, a set of factor templates \mathcal{T} exists so that factors of the same template $F \in \mathcal{T}$ share the same feature map $\phi_F(\cdot)$ and parameter vector \mathbf{w}_F . Then the response on input-output pair (\mathbf{x}, \mathbf{y}) is given by

$$\langle \mathbf{w}, \phi(\mathbf{x}, \mathbf{y}) \rangle = \sum_{F \in \mathcal{T}} \sum_{f \in F(\mathbf{x})} \langle \mathbf{w}_F, \phi_F(\mathbf{x}_f, \mathbf{y}_f) \rangle, \quad (3)$$

where $F(\mathbf{x})$ denotes the set of factors on \mathbf{x} that share a template F , and \mathbf{y}_f denotes output variables relevant to factor f of domain $\mathcal{Y}_f = \mathcal{Y}_F$. We will use $\mathcal{F}(\mathbf{x})$ to denote the union of factors of different templates $\{F(\mathbf{x})\}_{F \in \mathcal{T}}$. Figure 1 shows two examples that both have two factor templates (i.e. unigram and bigram) for which the responses have decomposition $\sum_{f \in u(\mathbf{x})} \langle \mathbf{w}_u, \phi_u(\mathbf{x}_f, \mathbf{y}_f) \rangle + \sum_{f \in b(\mathbf{x})} \langle \mathbf{w}_b, \phi_b(\mathbf{y}_f) \rangle$. Unfortunately, even with such decomposition, the maximization in (2) is still computationally expensive. First, most of graph structures do not allow exact maximization, so in practice one would minimize an upper bound of the original loss (2) obtained from relaxation [10, 18]. Second, even for the relaxed loss or a tree-structured graph that allows polynomial-time maximization, its complexity is at least linear to the cardinality of factor domain $|\mathcal{Y}_f|$ times the number of factors $|\mathcal{F}|$. This results in a prohibitive computational cost for problems with large output domain. As in Figure 1, one example has a factor domain $|\mathcal{Y}_b|$ which grows quadratically with the size of output domain; the other has the number of factors $|\mathcal{F}|$ which grows quadratically with the number of outputs. A key observation of this paper is, in contrast to the structural maximization (2) that requires larger extent of exploration on locally suboptimal assignments in order to achieve global optimality, the *Factorwise Maximization Oracle (FMO)*

$$\mathbf{y}_f^* := \underset{\mathbf{y}_f}{\operatorname{argmax}} \langle \mathbf{w}_F, \phi(\mathbf{x}_f, \mathbf{y}_f) \rangle \quad (4)$$

can be realized in a more efficient way by maintaining data structures on the factor parameters \mathbf{w}_F . In the next section, we develop globally-convergent algorithms that rely only on FMO, and provide realizations of *message-augmented FMO* with cost sublinear to the size of factor domain or to the number of factors.

3 A Dual-Decomposed Approach to Learning

We consider an upper bound of the loss (2) based on a Linear Program (LP) relaxation that is tight in case of a tree-structured graph and leads to a tractable approximation for general factor graphs [11, 18]:

$$L^{LP}(\mathbf{w}; \mathbf{x}, \bar{\mathbf{y}}) = \max_{(\mathbf{q}, \mathbf{p}) \in \mathcal{M}_L} \sum_{f \in \mathcal{F}(\mathbf{x})} \langle \boldsymbol{\theta}_f(\mathbf{w}), \mathbf{q}_f \rangle \quad (5)$$

where $\boldsymbol{\theta}_f(\mathbf{w}) := (\langle \mathbf{w}_F, \phi_F(\mathbf{x}_f, \mathbf{y}_f) - \phi_F(\mathbf{x}_f, \bar{\mathbf{y}}_f) \rangle + \delta_f(\mathbf{y}_f, \bar{\mathbf{y}}_f))_{\mathbf{y}_f \in \mathcal{Y}_f}$. \mathcal{M}_L is a polytope that constrains \mathbf{q}_f in a $|\mathcal{Y}_f|$ -dimensional simplex $\Delta^{|\mathcal{Y}_f|}$ and also enforces local consistency:

$$\mathcal{M}_L := \left\{ \begin{array}{l} \mathbf{q} = (\mathbf{q}_f)_{f \in \mathcal{F}(\mathbf{x})} \mid \mathbf{q}_f \in \Delta^{|\mathcal{Y}_f|}, \quad \forall f \in F(\mathbf{x}), \forall F \in \mathcal{T} \\ \mathbf{p} = (\mathbf{p}_j)_{j \in \mathcal{V}(\mathbf{x})} \mid M_{jf} \mathbf{q}_f = \mathbf{p}_j, \quad \forall (j, f) \in \mathcal{E}(\mathbf{x}) \end{array} \right\},$$

where M_{jf} is a $|\mathcal{Y}_j|$ by $|\mathcal{Y}_f|$ matrix that has $M_{jf}(y_j, \mathbf{y}_f) = 1$ if y_j is consistent with \mathbf{y}_f (i.e. $y_j = [\mathbf{y}_f]_j$) and $M_{jf}(y_j, \mathbf{y}_f) = 0$ otherwise. For a tree-structured graph $G(\mathcal{F}, \mathcal{V}, \mathcal{E})$, the LP relaxation is tight and

thus loss (5) is equivalent to (2). For a general factor graph, (5) is an upper bound on the original loss (2). It is observed that parameters \mathbf{w} learned from the upper bound (5) tend to tightening the LP relaxation and thus in practice lead to tight LP in the testing phase [10]. Instead of solving LP (5) as a subroutine, a recent attempt formulates (1) as a problem that optimizes (\mathbf{p}, \mathbf{q}) and \mathbf{w} jointly via dual decomposition [11, 12]. We denote $\boldsymbol{\lambda}_{jf}$ as dual variables associated with constraint $M_{jf}\mathbf{q}_f = \mathbf{p}_j$, and $\boldsymbol{\lambda}_f := (\boldsymbol{\lambda}_{jf})_{j \in \mathcal{N}(f)}$ where $\mathcal{N}(f) = \{j \mid (j, f) \in \mathcal{E}\}$. We have

$$L^{LP}(\mathbf{w}; \mathbf{x}, \bar{\mathbf{y}}) = \max_{\mathbf{q}, \mathbf{p}} \min_{\boldsymbol{\lambda}} \sum_{f \in \mathcal{F}(\mathbf{x})} \langle \boldsymbol{\theta}_f(\mathbf{w}), \mathbf{q}_f \rangle + \sum_{j \in \mathcal{N}(f)} \langle \boldsymbol{\lambda}_{jf}, M_{jf}\mathbf{q}_f - \mathbf{p}_j \rangle \quad (6)$$

$$= \min_{\boldsymbol{\lambda} \in \Lambda} \sum_{f \in \mathcal{F}(\mathbf{x})} \max_{\mathbf{q}_f \in \Delta^{|\mathcal{Y}_f|}} (\boldsymbol{\theta}_f(\mathbf{w}) + \sum_{j \in \mathcal{N}(f)} M_{jf}^T \boldsymbol{\lambda}_{jf})^T \mathbf{q}_f \quad (7)$$

$$= \min_{\boldsymbol{\lambda} \in \Lambda} \sum_{f \in \mathcal{F}(\mathbf{x})} \left(\max_{\mathbf{y}_f \in \mathcal{Y}_f} \theta_f(\mathbf{y}_f; \mathbf{w}) + \sum_{j \in \mathcal{N}(f)} \lambda_{jf}([\mathbf{y}_f]_j) \right) = \min_{\boldsymbol{\lambda} \in \Lambda} \sum_{f \in \mathcal{F}(\mathbf{x})} L_f(\mathbf{w}; \mathbf{x}_f, \bar{\mathbf{y}}_f, \boldsymbol{\lambda}_f) \quad (8)$$

where (7) follows the strong duality, and the domain $\Lambda = \left\{ \boldsymbol{\lambda} \mid \sum_{(j,f) \in \mathcal{E}(\mathbf{x})} \boldsymbol{\lambda}_{jf} = \mathbf{0}, \forall j \in \mathcal{V}(\mathbf{x}) \right\}$ follows the maximization w.r.t. \mathbf{p} in (6). The result (8) is a loss function $L_f(\cdot)$ that penalizes the response of each factor separately given $\boldsymbol{\lambda}_f$. The ERM problem (1) can then be expressed as

$$\min_{\mathbf{w}, \boldsymbol{\lambda} \in \Lambda} \sum_{F \in \mathcal{T}} \left(\frac{1}{2} \|\mathbf{w}_F\|^2 + C \sum_{f \in F} L_f(\mathbf{w}_F; \mathbf{x}_f, \bar{\mathbf{y}}_f, \boldsymbol{\lambda}_f) \right), \quad (9)$$

where $F = \bigcup_{i=1}^N F(\mathbf{x}_i)$ and $\mathcal{F} = \bigcup_{F \in \mathcal{T}} F$. The formulation (9) has an insightful interpretation: each factor template F learns a multiclass SVM given by parameters \mathbf{w}_F from factors $f \in F$, while each factor is augmented with messages $\boldsymbol{\lambda}_f$ passed from all variables related to f .

Despite the insightful interpretation, formulation (9) does not yield computational advantage directly. In particular, the non-smooth loss $L_f(\cdot)$ entangles parameters \mathbf{w} and messages $\boldsymbol{\lambda}$, which leads to a difficult optimization problem. Previous attempts to solve (9) either have slow convergence [11] or rely on an approximation objective [12]. In the next section, we propose a *Greedy Direction Method of Multiplier (GDMM)* algorithm for solving (9), which achieves ϵ sub-optimality in $O(\log(1/\epsilon))$ iterations while requires only one pass of FMOs for each iteration.

3.1 Greedy Direction Method of Multiplier

Let $\alpha_f(\mathbf{y}_f)$ be dual variables for the factor responses $z_f(\mathbf{y}_f) = \langle \mathbf{w}, \phi(\mathbf{x}_f, \mathbf{y}_f) \rangle$ and $\{\boldsymbol{\alpha}_j\}_{j \in \mathcal{V}}$ be that for constraints in Λ . The dual problem of (9) can be expressed as ¹

$$\begin{aligned} \min_{\boldsymbol{\alpha}_f \in \Delta^{|\mathcal{Y}_f|}} \quad & G(\boldsymbol{\alpha}) := \frac{1}{2} \sum_{F \in \mathcal{T}} \|\mathbf{w}_F(\boldsymbol{\alpha})\|^2 - \sum_{j \in \mathcal{V}} \boldsymbol{\delta}_j^T \boldsymbol{\alpha}_j \\ \text{s.t.} \quad & M_{jf} \boldsymbol{\alpha}_f = \boldsymbol{\alpha}_j, \quad j \in \mathcal{N}(f), f \in \mathcal{F}. \\ & \mathbf{w}_F(\boldsymbol{\alpha}) = \sum_{f \in F} \Phi_f^T \boldsymbol{\alpha}_f \end{aligned} \quad (10)$$

where $\boldsymbol{\alpha}_f$ lie in the shifted simplex

$$\Delta^{|\mathcal{Y}_f|} := \left\{ \boldsymbol{\alpha}_f \mid \alpha_f(\bar{\mathbf{y}}_f) \leq C, \alpha_f(\mathbf{y}_f) \leq 0, \forall \mathbf{y}_f \neq \bar{\mathbf{y}}_f, \sum_{\mathbf{y}_f \in \mathcal{Y}_f} \alpha_f(\mathbf{y}_f) = 0. \right\}. \quad (11)$$

Algorithm 1 Greedy Direction Method of Multiplier

0. Initialize $t = 0$, $\boldsymbol{\alpha}^0 = \mathbf{0}$, $\boldsymbol{\lambda}^0 = \mathbf{0}$ and $\mathcal{A}^0 = \mathcal{A}^{init}$.
for $t = 0, 1, \dots$ **do**
 1. Compute $(\boldsymbol{\alpha}^{t+1}, \mathcal{A}^{t+1})$ via one pass of Algorithm 2, 3, or 4.
 2. $\boldsymbol{\lambda}_{jf}^{t+1} = \boldsymbol{\lambda}_{jf}^t + \eta \left(M_{jf} \boldsymbol{\alpha}_f^{t+1} - \boldsymbol{\alpha}_j^{t+1} \right)$, $j \in \mathcal{N}(f)$, $\forall f \in \mathcal{F}$.
end for
-

Problem (10) can be interpreted as a summation of the dual objectives of $|\mathcal{T}|$ multiclass SVMs (each per factor template), connected with consistency constraints. To minimize (10) one factor at a time, we adopt a *Greedy Direction Method of Multiplier (GDMM)* algorithm that alternates between minimizing the *Augmented Lagrangian* function

$$\min_{\boldsymbol{\alpha}_f \in \Delta^{|\mathcal{Y}_f|}} \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda}^t) := G(\boldsymbol{\alpha}) + \frac{\rho}{2} \sum_{j \in \mathcal{N}(f), f \in \mathcal{F}} \|\mathbf{m}_{jf}(\boldsymbol{\alpha}, \boldsymbol{\lambda}^t)\|^2 - \|\boldsymbol{\lambda}_{jf}^t\|^2 \quad (12)$$

and updating the Lagrangian Multipliers (of consistency constraints)

$$\boldsymbol{\lambda}_{jf}^{t+1} = \boldsymbol{\lambda}_{jf}^t + \eta (M_{jf} \boldsymbol{\alpha}_f - \boldsymbol{\alpha}_j). \quad \forall j \in \mathcal{N}(f), f \in \mathcal{F}, \quad (13)$$

where $\mathbf{m}_{jf}(\boldsymbol{\alpha}, \boldsymbol{\lambda}^t) = M_{jf} \boldsymbol{\alpha}_f - \boldsymbol{\alpha}_j + \boldsymbol{\lambda}_{jf}^t$ plays the role of messages between $|\mathcal{T}|$ multiclass problems, and η is a constant step size. The procedure is outlined in Algorithm 1. The minimization (12) is conducted in an approximate and greedy fashion, in the aim of involving as few dual variables as possible. We discuss two greedy algorithms that suit two different cases in the following.

Factor of Large Domain For problems with large factor domains, we minimize (12) via a variant of *Frank-Wolfe* algorithm with *away steps* (AFW) [8], outlined in Algorithm 2. The AFW algorithm maintains the iterate $\boldsymbol{\alpha}^t$ as a linear combination of bases constructed during iterates

$$\boldsymbol{\alpha}^t = \sum_{\mathbf{v} \in \mathcal{A}^t} c_{\mathbf{v}}^t \mathbf{v}, \quad \mathcal{A}^t := \{\mathbf{v} \mid c_{\mathbf{v}}^t \neq 0\} \quad (14)$$

where \mathcal{A}^t maintains an active set of bases of non-zero coefficients. Each iteration of AFW finds a direction $\mathbf{v}^+ := (\mathbf{v}_f^+)_{f \in \mathcal{F}}$ leading to the most descent amount according to the current gradient, subject to the simplex constraints:

$$\mathbf{v}_f^+ := \underset{\mathbf{v}_f \in \Delta^{|\mathcal{Y}_f|}}{\operatorname{argmin}} \langle \nabla_{\boldsymbol{\alpha}_f} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t), \mathbf{v}_f \rangle = C(\mathbf{e}_{\bar{\mathbf{y}}_f} - \mathbf{e}_{\mathbf{y}_f^*}), \quad \forall f \in \mathcal{F} \quad (15)$$

where $\mathbf{y}_f^* := \operatorname{argmax}_{\mathbf{y}_f \in \mathcal{Y}_f \setminus \{\bar{\mathbf{y}}_f\}} \langle \nabla_{\boldsymbol{\alpha}_f} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t), \mathbf{e}_{\mathbf{y}_f} \rangle$ is the non-ground-truth labeling of factor f of highest response. In addition, AFW finds the *away direction*

$$\mathbf{v}^- := \underset{\mathbf{v} \in \mathcal{A}^t}{\operatorname{argmax}} \langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t), \mathbf{v} \rangle, \quad (16)$$

which corresponds to the basis that leads to the most descent amount when being removed. Then the update is determined by

$$\boldsymbol{\alpha}^{t+1} := \begin{cases} \boldsymbol{\alpha}^t + \gamma_F \mathbf{d}_F, & \langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}, \mathbf{d}_F \rangle < \langle \nabla_{\boldsymbol{\alpha}} \mathcal{L}, \mathbf{d}_A \rangle \\ \boldsymbol{\alpha}^t + \gamma_A \mathbf{d}_A, & \text{otherwise.} \end{cases} \quad (17)$$

¹ $\boldsymbol{\alpha}_j$ is also dual variables for responses on unigram factors. We define $\mathcal{U} := \mathcal{V}$ and $\boldsymbol{\alpha}_f := \boldsymbol{\alpha}_j$, $\forall f \in \mathcal{U}$.

Algorithm 2 Away-step Frank-Wolfe (AFW)

repeat

1. Find a greedy direction \mathbf{v}^+ satisfying (15).
2. Find an away direction \mathbf{v}^- satisfying (16).
3. Compute $\boldsymbol{\alpha}^{t+1}$ according to (17).
4. Maintain active set \mathcal{A}^t by (14).
5. Maintain $\mathbf{w}_F(\boldsymbol{\alpha})$ according to (10).

until A non-drop step is performed.

Algorithm 3 Block-Greedy Coordinate Descent (BGCD)

for $i \in [n]$ **do**

1. Find f^* satisfying (18) for i -th sample.
2. $\mathcal{A}_i^{s+1} = \mathcal{A}_i^s \cup \{f^*\}$.

for $f \in \mathcal{A}_i$ **do**

- 3.1 Update $\boldsymbol{\alpha}_f$ according to (19).
- 3.2 Maintain $\mathbf{w}_F(\boldsymbol{\alpha})$ according to (10).

end for

end for

where we choose between two descent directions $\mathbf{d}_F := \mathbf{v}^+ - \boldsymbol{\alpha}^t$ and $\mathbf{d}_A := \boldsymbol{\alpha}^t - \mathbf{v}^-$. The step size of each direction $\gamma_F := \arg \min_{\gamma \in [0,1]} \mathcal{L}(\boldsymbol{\alpha}^t + \gamma \mathbf{d}_F)$ and $\gamma_A := \arg \min_{\gamma \in [0, c_{v^-}]} \mathcal{L}(\boldsymbol{\alpha}^t + \gamma \mathbf{d}_A)$ can be computed exactly due to the quadratic nature of (12).

A step is called *drop step* if a step size $\gamma^* = c_{v^-}$ is chosen, which leads to the removal of a basis \mathbf{v}^- from the active set, and therefore the total number of drop steps can be bounded by half of the number of iterations t . Since a drop step could lead to insufficient descent, Algorithm 2 stops only if a *non-drop step* is performed. Note Algorithm 2 requires only a factorwise greedy search (15) instead of a structural maximization (2). In section 3.2 we show how the factorwise search can be implemented much more efficiently than structural ones. All the other steps (2-5) in Algorithm 2 can be computed in $O(|\mathcal{A}_f| \text{nnz}(\boldsymbol{\phi}_f))$, where $|\mathcal{A}_f|$ is the number of active states in factor f , which can be much smaller than $|\mathcal{Y}_f|$ when output domain is large.

In practice, a *Block-Coordinate Frank-Wolfe (BCFW)* method has much faster convergence than *Frank-Wolfe* method (Algorithm 2) [9, 13], but proving linear convergence for *BCFW* is also much more difficult [13], which prohibits its use in our analysis. In our implementation, however, we adopt the *BCFW* version since it turns out to be much more efficient. We include a detailed description on the *BCFW* version in Appendix-A (Algorithm 4).

Large Number of Factors Many structured prediction problems, such as alignment, segmentation, and multilabel prediction (Fig. 1, right), comprise binary variables and large number of factors with small domains, for which Algorithm 2 does not yield any computational advantage. For this type of problem, we minimize (12) via one pass of *Block-Greedy Coordinate Descent (BGCD)* (Algorithm 3) instead. Let Q_{\max} be an upper bound on the eigenvalue of Hessian matrix of each block $\nabla_{\boldsymbol{\alpha}_f}^2 \mathcal{L}(\boldsymbol{\alpha})$. For binary variables of pairwise factor, we have $Q_{\max} = 4(\max_{f \in F} \|\boldsymbol{\phi}_f\|^2 + 1)$. Each iteration of BGCD finds a factor that leads to the most progress

$$f^* := \underset{f \in \mathcal{F}(\mathbf{x}_i)}{\operatorname{argmin}} \left(\min_{\boldsymbol{\alpha}_f + \mathbf{d} \in \Delta^{|\mathcal{Y}_f|}} \langle \nabla_{\boldsymbol{\alpha}_f} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t), \mathbf{d} \rangle + \frac{Q_{\max}}{2} \|\mathbf{d}\|^2 \right). \quad (18)$$

for each instance \mathbf{x}_i , adds them into the set of active factors \mathcal{A}_i , and performs updates by solving block subproblems

$$\mathbf{d}_f^* = \underset{\boldsymbol{\alpha}_f + \mathbf{d} \in \Delta^{|\mathcal{Y}_f|}}{\operatorname{argmin}} \quad \langle \nabla_{\boldsymbol{\alpha}_f} \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t), \mathbf{d} \rangle + \frac{Q_{\max}}{2} \|\mathbf{d}\|^2 \quad (19)$$

for each factor $f \in \mathcal{A}_i$. Note $|\mathcal{A}_i|$ is bounded by the number of GDMM iterations and it converges to a constant much smaller than $|\mathcal{F}(\mathbf{x}_i)|$ in practice. We address in the next section how a joint FMO can be performed to compute (18) in time sublinear to $|\mathcal{F}(\mathbf{x}_i)|$ in the binary-variable case.

3.2 Greedy Search via Factorwise Maximization Oracle (FMO)

The main difference between the FMO and structural maximization oracle (2) is that the former involves only simple operations such as inner products or table look-ups for which one can easily come up with data structures or approximation schemes to lower the complexity. In this section, we present two approaches to realize sublinear-time FMOs for two types of factors widely used in practice. We will describe in terms of pairwise factors, but the approach can be naturally generalized to factors involving more variables.

Indicator Factor Factors $\theta_f(\mathbf{x}_f, \mathbf{y}_f)$ of the form

$$\langle \mathbf{w}_F, \phi_F(\mathbf{x}_f, \mathbf{y}_f) \rangle = v(\mathbf{x}_f, \mathbf{y}_f) \quad (20)$$

are widely used in practice. It subsumes the bigram factor $v(y_i, y_j)$ that is prevalent in sequence, grid, and network labeling problems, and also factors that map an input-output pair (x, y) directly to a score $v(x, y)$. For this type of factor, one can maintain ordered multimaps for each factor template F , which support ordered visits of $\{v(\mathbf{x}, (y_i, y_j))\}_{(y_i, y_j) \in \mathcal{Y}_f}$, $\{v(\mathbf{x}, (y_i, y_j))\}_{y_j \in \mathcal{Y}_j}$ and $\{v(\mathbf{x}, (y_i, y_j))\}_{y_i \in \mathcal{Y}_i}$. Then to find \mathbf{y}_f that maximizes (26), we compare the maximizers in 4 cases: (i) $(y_i, y_j) : m_{if}(y_i) = m_{jf}(y_j) = 0$, (ii) $(y_i, y_j) : m_{if}(y_i) = 0$, (iii) $(y_i, y_j) : m_{jf}(y_j) = 0$, (iv) $(y_i, y_j) : m_{jf}(y_j) \neq 0, m_{if}(y_i) \neq 0$. The maximization requires $O(|\mathcal{A}_i||\mathcal{A}_j|)$ in cases (ii)-(iv) and $O(\max(|\mathcal{A}_i||\mathcal{Y}_j|, |\mathcal{Y}_i||\mathcal{A}_j|))$ in case (i) (see details in Appendix C-1). However, in practice we observe an $O(1)$ cost for case (i) and the bottleneck is actually case (iv), which requires $O(|\mathcal{A}_i||\mathcal{A}_j|)$.

Note the ordered multimaps need maintenance whenever the vector $\mathbf{w}_F(\boldsymbol{\alpha})$ is changed. Fortunately, since the indicator factor has $v(\mathbf{y}_f, \mathbf{x}) = \sum_{f \in F, \mathbf{x}_f = \mathbf{x}} \alpha_f(\mathbf{y}_f)$, each update (25) leads to at most $|\mathcal{A}_f|$ changed elements, which gives a maintenance cost bounded by $O(|\mathcal{A}_f| \log(|\mathcal{Y}_F|))$. On the other hand, the space complexity is bounded by $O(|\mathcal{Y}_F||\mathcal{X}_F|)$ since the map is shared among factors.

Binary-Variable Interaction Factor Many problems consider pairwise-interaction factor between binary variables, where the factor domain is small but the number of factors is large. For this type of problem, there is typically an rare outcome $\mathbf{y}_f^A \in \mathcal{Y}_f$. We call factors exhibiting such outcome as *active factors* and the score of a labeling is determined by the score of the active factors (inactive factors give score 0). For example, in the problem of *multilabel prediction with pairwise interactions* (Fig. 1, right), an active unigram factor has outcome $\mathbf{y}_f^A = 1$ and an active bigram factor has $\mathbf{y}_f^A = (1, 1)$, and each sample typically has only few outputs with value 1.

For this type of problem, we show that the gradient magnitude w.r.t. $\boldsymbol{\alpha}_f$ for a bigram factor f can be determined by the gradient w.r.t. $\alpha_f(\mathbf{y}_f^A)$ when one of its incoming message m_{jf} or m_{if} is 0 (see details in Appendix C-2). Therefore, we can find the greedy factor (18) by maintaining an ordered multimap for the scores of outcome \mathbf{y}_f^A in each factor $\{v(\mathbf{y}_f^A, \mathbf{x}_f)\}_{f \in F}$. The resulting complexity for finding a factor that maximizes (18) is then reduced from $O(|\mathcal{Y}_i||\mathcal{Y}_j|)$ to $O(|\mathcal{A}_i||\mathcal{A}_j|)$, where the latter is for comparison among factors that have both messages m_{if} and m_{jf} being non-zero.

Inner-Product Factor We consider another widely-used type of factor of the form

$$\theta_f(\mathbf{x}_f, \mathbf{y}_f) = \langle \mathbf{w}_F, \phi_F(\mathbf{x}_f, \mathbf{y}_f) \rangle = \langle \mathbf{w}_F(\mathbf{y}_f), \phi_F(\mathbf{x}_f) \rangle$$

where all labels $\mathbf{y}_f \in \mathcal{Y}_f$ share the same feature mapping $\phi_F(\mathbf{x}_f)$ but with different parameters $\mathbf{w}_F(\mathbf{y}_f)$. We propose a simple sampling approximation method with a performance guarantee for the convergence of GDM. Note although one can apply similar sampling schemes to the structural maximization oracle (2), it is hard to guarantee the quality of approximation. The sampling method

divides \mathcal{Y}_f into ν mutually exclusive subsets $\mathcal{Y}_f = \bigcup_{k=1}^{\nu} \mathcal{Y}_f^{(k)}$, and realizes an approximate FMO by first sampling k uniformly from $[\nu]$ and returning

$$\hat{\mathbf{y}}_f \in \arg \max_{\mathbf{y}_f \in \mathcal{Y}_f^{(k)}} \langle \mathbf{w}_F(\mathbf{y}_f), \phi_F(\mathbf{x}_f) \rangle. \quad (21)$$

Note there is at least $1/\nu$ probability that $\hat{\mathbf{y}}_f \in \arg \max_{\mathbf{y}_f \in \mathcal{Y}_f} \langle \mathbf{w}_F(\mathbf{y}_f), \phi_F(\mathbf{x}_f) \rangle$ since at least one partition $\mathcal{Y}_f^{(k)}$ contains a label of the highest score. In section 3.3, we show that this approximate FMO still ensures convergence with a rate scaled by $1/\nu$. In practice, since the set of active labels is not changing frequently during training, once an active label \mathbf{y}_f is sampled, it will be kept in the active set \mathcal{A}_f till the end of the algorithm and thus results in a convergence rate similar to that of an exact FMO. Note for problems of binary variables with large number of inner-product factors, the sampling technique applies similarly by simply partitioning factors as $\mathcal{F}_i = \bigcup_{k=1}^{\nu} \mathcal{F}_i^{(k)}$ and searching active factors only within one randomly chosen partition at a time.

3.3 Analysis

We show the iteration complexity of the GDMM algorithm with an $1/\nu$ -approximated FMO given in section 3.2. The convergence guarantee for exact FMOs can be obtained by setting $\nu = 1$. The analysis leverages recent analysis on the global linear convergence of Frank-Wolfe variants [8] for function of the form (12) with a polyhedral domain, and also the analysis in [5] for Augmented Lagrangian based method. This type of greedy Augmented Lagrangian Method was also analyzed previously under different context [22, 23, 24]. Let $d(\boldsymbol{\lambda}) = \min_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda})$ be the dual objective of (12), and let

$$\Delta_d^t := d^* - d(\boldsymbol{\lambda}^t), \quad \Delta_p^t := \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - d(\boldsymbol{\lambda}^t) \quad (22)$$

be the dual and primal suboptimality of problem (10) respectively. We have the following theorems.

Theorem 1 (Convergence of GDMM with AFW). *The iterates $\{(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t)\}_{t=1}^{\infty}$ produced by Algorithm 1 with step 1 performed by Algorithm 2 has*

$$E[\Delta_p^t + \Delta_d^t] \leq \epsilon \text{ for } t \geq \omega \log\left(\frac{1}{\epsilon}\right) \quad (23)$$

for any $0 < \eta \leq \frac{\rho}{4+16(1+\nu)mQ/\mu_{\mathcal{M}}}$ with $\omega = \max\left\{2\left(1 + 4\frac{mQ(1+\nu)}{\mu_{\mathcal{M}}}\right), \frac{\tau}{\eta}\right\}$, where $\mu_{\mathcal{M}}$ is the generalized geometric strong convexity constant of (12), Q is the Lipschitz-continuous constant for the gradient of objective (12), and $\tau > 0$ is a constant depending on optimal solution set.

Theorem 2 (Convergence of GDMM with BGCD). *The iterates $\{(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t)\}_{t=1}^{\infty}$ produced by Algorithm 1 with step 1 performed by Algorithm 3 has*

$$E[\Delta_p^t + \Delta_d^t] \leq \epsilon \text{ for } t \geq \omega_1 \log\left(\frac{1}{\epsilon}\right) \quad (24)$$

for any $0 < \eta \leq \frac{\rho}{4(1+Q_{\max\nu/\mu_1})}$ with $\omega_1 = \max\left\{2\left(1 + \frac{Q_{\max\nu}}{\mu_1}\right), \frac{\tau}{\eta}\right\}$, where μ_1 is the generalized strong convexity constant of objective (12) and $Q_{\max} = \max_{f \in \mathcal{F}} Q_f$ is the factorwise Lipschitz-continuous constant on the gradient.

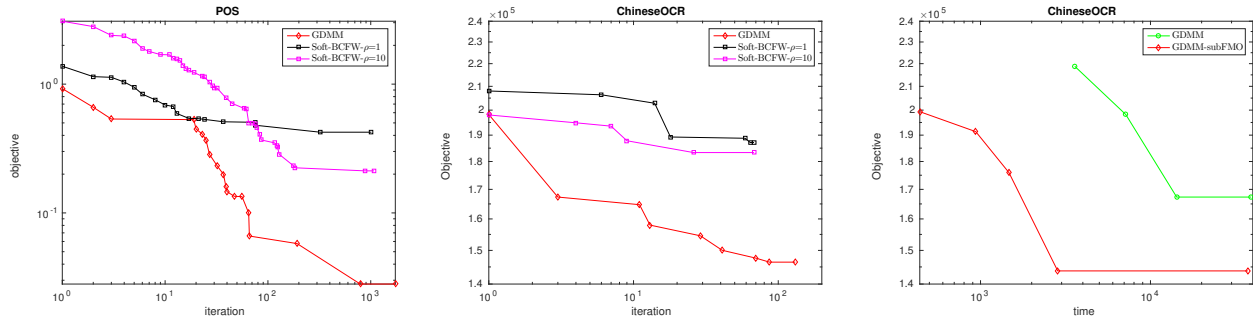


Figure 2: (left) Compare two FMO-based algorithms (GDMM, Soft-BCFW) in number of iterations. (right) Improvement in training time given by sublinear-time FMO.

4 Experiment

4.1 Data

Our experiments are conducted on 4 public datasets: *POS*, *ChineseOCR*, *RCV1-regions*, and *EUR-Lex* (directory codes). For sequence labeling we experiment on *POS* and *ChineseOCR*. The *POS* dataset is a subset of Penn treebank² that contains 3,808 sentences, 196,223 words, and 45 POS labels. The HIT-MW³ *ChineseOCR* dataset is a hand-written Chinese character dataset from [17]. The dataset has 12,064 hand-written sentences, and a total of 174,074 characters. The vocabulary (label) size is 3,039. For the Correlated Multilabel Prediction problems, we experiment on two benchmark datasets *RCV1-regions*⁴ and *EUR-Lex* (directory codes)⁵. The *RCV1-regions* dataset has 228 labels, 23,149 training instances and 47,236 features. Note that a smaller version of *RCV1* with only 30 labels and 6000 instances is used in [11, 12]. *EUR-Lex* (directory codes) has 410 directory codes as labels with a sample size of 19,348.

4.2 Results

In this section, we compare with existing approaches on Sequence Labeling and Multi-label prediction with pairwise interaction. The algorithms in comparison are: (i) *BCFW*: a Block-Coordinate Frank-Wolfe method based on structural oracle [9], which outperforms other competitors such as Cutting-Plane, FW, and online-EG methods in [9]. (ii) *SSG*: an implementation of the Stochastic Subgradient method [16]. (iii) *Soft-BCFW*: Algorithm proposed in ([12]), which avoids structural oracle by minimizing an approximate objective, where a parameter ρ controls the precision of the approximation. We tuned the parameter and chose two of the best on the figure. For BCFW and SSG, we adapted the MATLAB implementation provided by authors of [9] into C++, which is an order of magnitude faster. All other implementations are also in C++. The results are compared in terms of primal objective (achieved by w) and test accuracy.

We first compare GDMM (without subFMO) with Soft-BCFW in Figure 2. Due to the approximation (controlled by ρ), Soft-BCFW can converge to a suboptimal primal objective value. While the gap decreases as ρ increases, its convergence becomes also slower. GDMM, on the other hand, enjoys a faster convergence. The sublinear-time implementation of FMO also reduces the training time by an order of magnitude on the ChineseOCR data set, as showed in Figure 2 (right). More general experiments are showed in Figure 3. When the size of output domain is small (*POS* dataset),

²<https://catalog.ldc.upenn.edu/LDC99T42>

³<https://sites.google.com/site/hitmwd/>

⁴www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/multilabel.html

⁵mulan.sourceforge.net/datasets-mlc.html

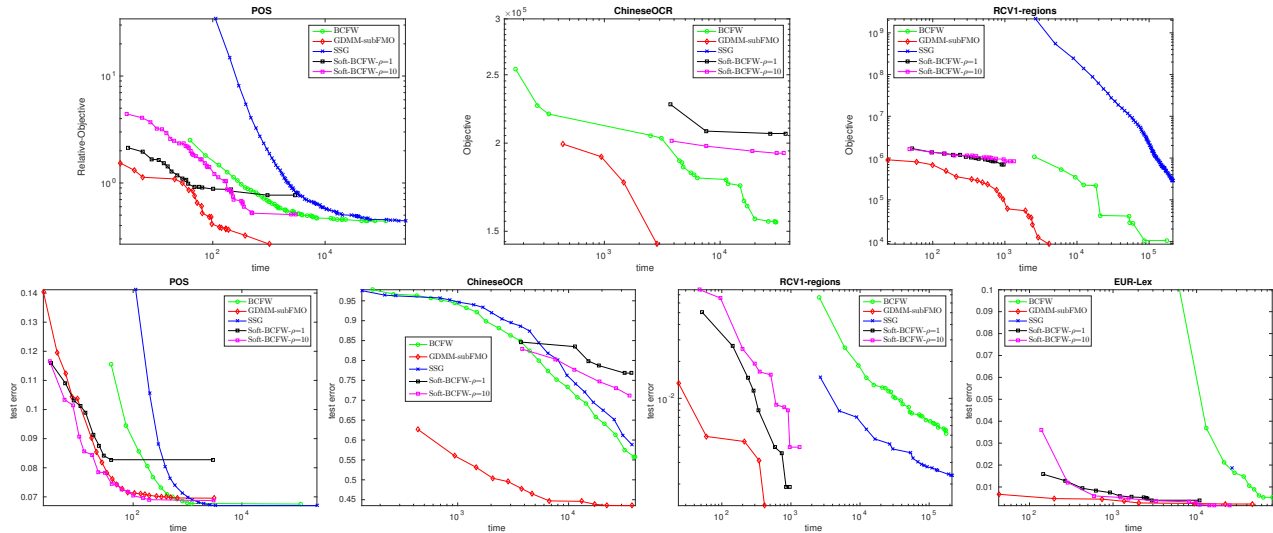


Figure 3: Primal Objective v.s. Time and Test error v.s. Time plots. Note that subfigures of objective have showed that SSG converges to a objective value much higher than all other methods, this is also observed in [9].

GDMM-subFMO is competitive to other solvers. As the size of output domain grows (ChineseOCR, RCV1, EUR-Lex), the complexity of structural maximization oracle grows linearly or even quadratically, while the complexity of GDMM-subFMO only grows sublinearly in the experiments. Therefore, GDMM-subFMO achieves orders-of-magnitude speedup over other methods. In particular, when running on ChineseOCR and EUR-Lex, each iteration of SSG, GDMM, BCFW and Soft-BCFW take over 10^3 seconds, while it only takes a few seconds in GDMM-subFMO.

References

- [1] D. Das, D. Chen, A. F. Martins, N. Schneider, and N. A. Smith. Frame-semantic parsing. *Computational linguistics*, 40(1):9–56, 2014.
- [2] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. Liblinear: A library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874, 2008.
- [3] K. Gimpel and N. A. Smith. Structured ramp loss minimization for machine translation. In *NAACL*, pages 221–231. Association for Computational Linguistics, 2012.
- [4] A. Hoffman. On approximate solutions of systems of linear inequalities. *Journal of Research of the National Bureau of Standards*, 1952.
- [5] M. Hong and Z.-Q. Luo. On the linear convergence of the alternating direction method of multipliers. *arXiv preprint arXiv:1208.3922*, 2012.
- [6] T. Joachims, T. Finley, and C.-N. J. Yu. Cutting-plane training of structural svms. *Machine Learning*, 77(1):27–59, 2009.
- [7] M. P. Kumar, V. Kolmogorov, and P. H. Torr. An analysis of convex relaxations for map estimation. *Advances in Neural Information Processing Systems*, 20:1041–1048, 2007.
- [8] S. Lacoste-Julien and M. Jaggi. On the global linear convergence of frank-wolfe optimization variants. In *Advances in Neural Information Processing Systems*, pages 496–504, 2015.

- [9] S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-coordinate frank-wolfe optimization for structural svms. In *ICML 2013 International Conference on Machine Learning*, pages 53–61, 2013.
- [10] O. Meshi, M. Mahdavi, and D. Sontag. On the tightness of lp relaxations for structured prediction. *arXiv preprint arXiv:1511.01419*, 2015.
- [11] O. Meshi, D. Sontag, T. Jaakkola, and A. Globerson. Learning efficiently with approximate inference via dual losses. 2010.
- [12] O. Meshi, N. Srebro, and T. Hazan. Efficient training of structured svms via soft constraints. In *AISTAT*, 2015.
- [13] A. Osokin, J.-B. Alayrac, I. Lukasewitz, P. K. Dokania, and S. Lacoste-Julien. Minding the gaps for block frank-wolfe optimization of structured svms. *arXiv preprint arXiv:1605.09346*, 2016.
- [14] P. Ravikumar and J. Lafferty. Quadratic programming relaxations for metric labeling and markov random field map estimation. In *ICML*, 2006.
- [15] R. Samdani and D. Roth. Efficient decomposed learning for structured prediction. *ICML*, 2012.
- [16] S. Shalev-Shwartz, Y. Singer, N. Srebro, and A. Cotter. Pegasos: Primal estimated sub-gradient solver for svm. *Mathematical programming*, 2011.
- [17] T. Su, T. Zhang, and D. Guan. Corpus-based hit-mw database for offline recognition of general-purpose chinese handwritten text. *IJDAR*, 10(1):27–38, 2007.
- [18] B. Taskar, V. Chatalbashev, D. Koller, and C. Guestrin. Learning structured prediction models: A large margin approach. In *ICML*, 2005.
- [19] B. Taskar, C. Guestrin, and D. Koller. Max-margin markov networks. In *Advances in neural information processing systems*, volume 16, 2003.
- [20] I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun. Support vector machine learning for interdependent and structured output spaces. In *ICML*, 2004.
- [21] P. Woodland and D. Povey. Large scale discriminative training for speech recognition. In *ASR2000-Automatic Speech Recognition Workshop (ITRW)*, 2000.
- [22] I. E. Yen, X. Lin, E. J. Zhang, E. P. Ravikumar, and I. S. Dhillon. A convex atomic-norm approach to multiple sequence alignment and motif discovery. 2016.
- [23] I. E. Yen, D. Malioutov, and A. Kumar. Scalable exemplar clustering and facility location via augmented block coordinate descent with column generation. In *AISTAT*, 2016.
- [24] I. E.-H. Yen, K. Zhong, C.-J. Hsieh, P. K. Ravikumar, and I. S. Dhillon. Sparse linear programming via primal and dual augmented coordinate descent. In *NIPS*, 2015.

5 Appendix A: Block-Coordinate Frank-Wolfe (BCFW)— Practically Faster-Convergent Variant

Algorithm 4 Block-Coordinate Frank Wolfe (improving upon Algorithm 2)

for $s = 1$ to $|\mathcal{F}|$ **do**

1. Draw $f \in \mathcal{F}$ uniformly at random.
2. Find a greedy direction \mathbf{v}_f^+ satisfying (15).
3. $\mathcal{A}_f^{s+1} = \mathcal{A}_f^s \cup \{\mathbf{v}_f^+\}$.
4. Solve (25) with active set \mathcal{A}_f^{s+1} .
5. Maintain $\mathbf{w}_F(\boldsymbol{\alpha})$.

end for

The *Block-Coordinate Frank-Wolfe (BCFW)* (Algorithm 4) differs from *Frank-Wolfe* (Algorithm 2) in that it updates dual variables $\boldsymbol{\alpha}_f$ of each factor sequentially, and the bases \mathbf{v}_f and active sets \mathcal{A}_f are maintained for each factor f separately. For each iteration of BCFW, we find a greedy direction \mathbf{v}_f^+ in the same way (15) as AFW, but for one factor at a time. Then we add \mathbf{v}_f^+ to an active set \mathcal{A}_f maintained for each factor. Since in BCFW we update one factor at a time, we can minimize the following *block active-set subproblem*

$$\mathbf{d}_{\mathcal{A}_f}^* = \underset{\boldsymbol{\alpha}_f + \mathbf{d}_{\mathcal{A}_f} \in \Delta^{|\mathcal{Y}_f|}}{\operatorname{argmin}} \quad \langle \nabla_{\boldsymbol{\alpha}_f} \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda}^t), \mathbf{d}_{\mathcal{A}_f} \rangle + \frac{Q_f}{2} \|\mathbf{d}_{\mathcal{A}_f}\|^2 \quad (25)$$

where Q_f is an upper bound on the Hessian of variables in the active set (discussed in section 5.1). The active-set subproblem (25) can be solved via a simplex projection in time $O(|\mathcal{A}_f|)$ [2]. Furthermore, by maintaining $\mathbf{w}_F(\boldsymbol{\alpha})$ after solving each sub-problem (25), we can compute the gradient

$$\nabla_{\boldsymbol{\alpha}_f(\mathbf{y}_f)} \mathcal{L} = \langle \mathbf{w}_F, \boldsymbol{\phi}_f(\mathbf{x}_f, \mathbf{y}_f) \rangle - \delta_f(\mathbf{y}_f, \bar{\mathbf{y}}_f) + \rho_f \sum_{j \in \mathcal{N}(f), \mathbf{y}_j = [\mathbf{y}_f]_j} m_{jf}(\mathbf{y}_j) \quad (26)$$

for $\mathbf{y}_f \in \mathcal{A}_f$ in time $O(|\mathcal{A}_f| \operatorname{nnz}(\boldsymbol{\phi}_f))$, where $\rho_f = -\rho$, $\delta_f = \delta_j$ for $f \in \mathcal{U}$ and $\rho_f = \rho$, $\delta_f = 0$ for $f \notin \mathcal{U}$. Note the size of active set $|\mathcal{A}_f|$ is bounded by the number of GDMM iterations, and in practice $|\mathcal{A}_f|$ converges to a constant much smaller than $|\mathcal{Y}_f|$ for problems of large output domains. Therefore, the bottleneck of the BCFW algorithm lies in the step (15), which as we show in Section 3.2, can be computed in time sublinear to $|\mathcal{Y}_f|$ via an efficient FMO.

5.1 Constant Q_f in Problem (25)

The constant Q_f is an upper bound on the maximum eigenvalue of the Hessian submatrix for variables in the active set \mathcal{A}_f , that is, $\|[\Phi_f \Phi_f^T]_{\mathcal{A}_f}\| + \rho \sum_{j \in \mathcal{N}(f)} \|[M_{jf}^T M_{jf}]_{\mathcal{A}_f}\|$, where the notation $[\cdot]_{\mathcal{A}}$ denotes the sub-matrix formed by row and column indexes in \mathcal{A} and $\|\cdot\|$ is the spectral norm of a matrix.

For many types of factors used in practice, $Q_f = O(|\mathcal{A}_f|)$ and is easy to compute in the beginning.

In particular, for unigram factor, we have $M_{jf} = I$ and thus $\|[M_{jf}^T M_{jf}]_{\mathcal{A}}\| = 1$, and for higher-order factor we have $\|[M_{jf}^T M_{jf}]_{\mathcal{A}}\| = |\mathcal{A}|$.

As for the term $\|[\Phi_f \Phi_f^T]_{\mathcal{A}}\|$. In most of applications, Φ_f is a $|\mathcal{Y}_f| \times (|\mathcal{Y}_f|d)$ block-diagonal matrix that duplicates $1 \times d$ feature vector $\boldsymbol{\phi}_f(\mathbf{x}_f)^T$ for $|\mathcal{Y}_f|$ times, for which we have $\|[\Phi_f \Phi_f^T]_{\mathcal{A}}\| = \|\boldsymbol{\phi}_f(\mathbf{x}_f)\|^2$. Note in this case, the quadratic upper bound in (25) is tight for unigram factors.

6 Appendix B: Convergence of GDMM

6.1 Proof of Theorem (1)

Recall that the Augmented Lagrangian $\mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda})$ is of the form

$$\mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda}) := G(\boldsymbol{\alpha}) + \langle \boldsymbol{\lambda}, M\boldsymbol{\alpha} \rangle + \frac{\rho}{2} \|M\boldsymbol{\alpha}\|^2.$$

where M is "the number of consistency constraints" by "the number of variables" matrix and $M\boldsymbol{\alpha} = \mathbf{0}$ encodes all constraints of the form

$$M_{jf}\boldsymbol{\alpha}_f - \boldsymbol{\alpha}_j = \begin{bmatrix} M_{jf} & -I_j \end{bmatrix} \begin{bmatrix} \boldsymbol{\alpha}_f \\ \boldsymbol{\alpha}_j \end{bmatrix} = \mathbf{0}.$$

The function

$$G(\boldsymbol{\alpha}) = \frac{1}{2} \sum_{F \in \mathcal{F}} \|\mathbf{w}_F(\boldsymbol{\alpha}_F)\|^2 - \sum_{j \in \mathcal{U}} \boldsymbol{\delta}_j^T \boldsymbol{\alpha}_j$$

can be written in a compact form as

$$\begin{aligned} G(\boldsymbol{\alpha}) &= \frac{1}{2} \|\mathbf{w}(\boldsymbol{\alpha})\|^2 + \boldsymbol{\delta}^T \boldsymbol{\alpha} \\ &= \frac{1}{2} \|\Phi^T \boldsymbol{\alpha}\|^2 + \boldsymbol{\delta}^T \boldsymbol{\alpha} \end{aligned} \tag{27}$$

where Φ is the "number of variables (in $\boldsymbol{\alpha}$)" by "number of parameters (in \mathbf{w})" design matrix.

Now let $\boldsymbol{\alpha}$ be the "primal variables" and denote

$$\boldsymbol{\alpha}(\boldsymbol{\lambda}) := \{\boldsymbol{\alpha} \mid \boldsymbol{\alpha} = \arg \min_{\boldsymbol{\alpha}} \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda})\} \tag{28}$$

with

$$\bar{\boldsymbol{\alpha}}^t := \underset{\bar{\boldsymbol{\alpha}} \in \boldsymbol{\alpha}(\boldsymbol{\lambda}^t)}{\operatorname{argmin}} \|\bar{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^t\|,$$

and let $\mathcal{M} = \{\boldsymbol{\alpha} \mid \boldsymbol{\alpha}_f \in \Delta^{|\mathcal{Y}_f|}, \forall f \in \mathcal{F}\}$. The dual objective of the augmented problem is

$$d(\boldsymbol{\lambda}) = \min_{\boldsymbol{\alpha} \in \mathcal{M}} \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda})$$

and

$$d^* = \max_{\boldsymbol{\lambda}} d(\boldsymbol{\lambda})$$

is the optimal dual objective value.

Then we measure the sub-optimality of iterates $\{(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t)\}_{t=1}^T$ given by GDMM in terms of dual function difference

$$\Delta_d^t = d^* - d(\boldsymbol{\lambda}^t)$$

and the primal function difference for a given dual iterate $\boldsymbol{\lambda}^t$:

$$\Delta_p^t = \mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t) - d(\boldsymbol{\lambda}^t)$$

yielded by $\boldsymbol{\alpha}^{t+1}$ obtained from one pass of FC-BCFW algorithm on $\boldsymbol{\alpha}$. Then we have following lemma.

Lemma 1 (Dual Progress). *Each iteration of GDMM (Algorithm 1) has*

$$\Delta_d^t - \Delta_d^{t-1} \leq -\eta (M\boldsymbol{\alpha}^t)^T (M\bar{\boldsymbol{\alpha}}^t). \tag{29}$$

Proof.

$$\begin{aligned}
\Delta_d^t - \Delta_d^{t-1} &= d^* - d(\boldsymbol{\lambda}^t) - d^* - d(\boldsymbol{\lambda}^{t-1}) \\
&= \mathcal{L}(\bar{\boldsymbol{\alpha}}^{t-1}, \boldsymbol{\lambda}^{t-1}) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t) \\
&\leq \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^{t-1}) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t) \\
&= \langle \boldsymbol{\lambda}^{t-1} - \boldsymbol{\lambda}^t, M\bar{\boldsymbol{\alpha}}^t \rangle \\
&= -\eta \langle M\boldsymbol{\alpha}^t, M\bar{\boldsymbol{\alpha}}^t \rangle
\end{aligned}$$

where the first inequality follows the optimality of $\bar{\boldsymbol{\alpha}}^{t-1}$ for the function $\mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda}^{t-1})$ defined by $\boldsymbol{\lambda}^{t-1}$, and the last equality follows the dual update in GDMM (13). \square

On the other hand, the following lemma gives an expression on the primal progress that is independent of the algorithm used for minimizing Augmented Lagrangian

Lemma 2 (Primal Progress). *Each iteration of GDMM (Algorithm 1) has*

$$\begin{aligned}
\Delta_p^t - \Delta_p^{t-1} &\leq \mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t) - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) \\
&\quad + \eta \|M\boldsymbol{\alpha}^t\|^2 - \eta \langle M\boldsymbol{\alpha}^t, M\bar{\boldsymbol{\alpha}}^t \rangle
\end{aligned}$$

Proof.

$$\begin{aligned}
&\Delta_p^t - \Delta_p^{t-1} \\
&= \mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t) - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^{t-1}) - (d(\boldsymbol{\lambda}^t) - d(\boldsymbol{\lambda}^{t-1})) \\
&\leq \mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t) - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) + \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^{t-1}) + (d(\boldsymbol{\lambda}^{t-1}) - d(\boldsymbol{\lambda}^t)) \\
&\leq \mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t) - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) + \eta \|M\boldsymbol{\alpha}^t\|^2 - \eta \langle M\boldsymbol{\alpha}^t, M\bar{\boldsymbol{\alpha}}^t \rangle
\end{aligned}$$

where the last inequality uses Lemma 1 on $d(\boldsymbol{\lambda}^{t-1}) - d(\boldsymbol{\lambda}^t) = \Delta_d^t - \Delta_d^{t-1}$. \square

By combining results of Lemma 1 and 2, we can obtain a joint progress of the form

$$\begin{aligned}
&\Delta_d^t - \Delta_d^{t-1} + \Delta_p^t - \Delta_p^{t-1} \\
&\leq \mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t) - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) + \eta \|M\boldsymbol{\alpha}^t - M\bar{\boldsymbol{\alpha}}^t\|^2 - \eta \|M\bar{\boldsymbol{\alpha}}^t\|^2
\end{aligned} \tag{30}$$

Note the only positive term in (30) is the second one. To guarantee the descent of joint progress, we upper bound the three terms in (30) with the following lemmas.

Lemma 3.

$$\|M\boldsymbol{\alpha}^t - M\bar{\boldsymbol{\alpha}}^t\|^2 \leq \frac{2}{\rho} (\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)) \tag{31}$$

Proof. Let

$$\tilde{\mathcal{L}}(\boldsymbol{\alpha}, \boldsymbol{\lambda}) = h(\boldsymbol{\alpha}) + \frac{\rho}{2} \|M\boldsymbol{\alpha}\|^2,$$

where

$$h(\boldsymbol{\alpha}) = G(\boldsymbol{\alpha}) + \langle \boldsymbol{\lambda}, M\boldsymbol{\alpha} \rangle + \mathbf{I}_{\boldsymbol{\alpha} \in \mathcal{M}}.$$

, $\mathbf{I}_{\boldsymbol{\alpha} \in \mathcal{M}} = 0$ if $\boldsymbol{\alpha} \in \mathcal{M}$ and $\mathbf{I}_{\boldsymbol{\alpha} \in \mathcal{M}} = \infty$ otherwise. Note we have $\tilde{\mathcal{L}}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t) = \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)$ and $\tilde{\mathcal{L}}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) = \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t)$ due to feasible iterates. Due to the optimality of $\bar{\boldsymbol{\alpha}}^t$, we have

$$\mathbf{0} = \boldsymbol{\sigma} + M^T M \bar{\boldsymbol{\alpha}}^t \in \partial_{\boldsymbol{\alpha}} \tilde{\mathcal{L}}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)$$

for some $\boldsymbol{\sigma} \in \partial h(\bar{\boldsymbol{\alpha}}^t)$. And by the convexity of $h(\cdot)$ and the strong convexity of $\frac{\rho}{2} \|\cdot\|^2$, we have

$$h(\boldsymbol{\alpha}^t) - h(\bar{\boldsymbol{\alpha}}^t) \geq \langle \boldsymbol{\sigma}, \boldsymbol{\alpha}^t - \bar{\boldsymbol{\alpha}}^t \rangle$$

and

$$\frac{\rho}{2}\|M\boldsymbol{\alpha}^t\|^2 - \frac{\rho}{2}\|M\bar{\boldsymbol{\alpha}}^t\|^2 \geq \langle M^T M\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\alpha}^t - \bar{\boldsymbol{\alpha}}^t \rangle + \frac{\rho}{2}\|M\boldsymbol{\alpha}^t - M\bar{\boldsymbol{\alpha}}^t\|^2$$

Then the above two together imply

$$\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t) \geq \frac{\rho}{2}\|M(\boldsymbol{\alpha}^t) - M(\bar{\boldsymbol{\alpha}}^t)\|^2$$

which leads to our conclusion. \square

Lemma 4 (Hong and Luo 2012). *There is a constant $\tau > 0$ such that*

$$\Delta_d(\boldsymbol{\lambda}) \leq \tau\|M\bar{\boldsymbol{\alpha}}(\boldsymbol{\lambda})\|^2. \quad (32)$$

for any $\boldsymbol{\lambda}$ and any minimizer $\bar{\boldsymbol{\alpha}}(\boldsymbol{\lambda})$ satisfying (28).

Proof. This is a lemma adapted from [5]. Since our objective (12) satisfies the assumptions $A(a)$ – $A(e)$ and $A(g)$ in [5]. Then Lemma 3.1 of [5] guarantees that, as long as $\|\nabla d(\boldsymbol{\lambda})\|$ is bounded, there is a constant $\tau > 0$ s.t.

$$\Delta_d(\boldsymbol{\lambda}) \leq \tau\|\nabla d(\boldsymbol{\lambda})\|^2 = \|M\bar{\boldsymbol{\alpha}}(\boldsymbol{\lambda})\|^2$$

for all $\boldsymbol{\lambda}$. Note our problem satisfies the condition of bounded gradient magnitude since

$$\|\nabla d(\boldsymbol{\lambda})\| = \|M\bar{\boldsymbol{\alpha}}(\boldsymbol{\lambda})\| \leq \|M\bar{\boldsymbol{\alpha}}(\boldsymbol{\lambda})\|_1 \leq \|M\|_1\|\bar{\boldsymbol{\alpha}}(\boldsymbol{\lambda})\|_1 \leq (\max_f |\mathcal{Y}_f|)|\mathcal{F}|$$

where the last inequality is because $\bar{\boldsymbol{\alpha}}(\boldsymbol{\lambda})$ lies in a simplex domain. \square

The remaining thing is to show that one pass of AFW (Algorithm 2) or BGCD (Algorithm 3) suffices to give a descent amount $\mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t) - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t)$ lower bounded by some constant multiple of the primal sub-optimality $\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)$. If it is true, then by selecting a small enough GDMM step size η , the RHS of (30) would be negative. For AFW (Algorithm 2), this can be achieved by leveraging recent results from [8], which shows linear convergence of AFW, even for non-strongly convex function of the form (34). We thus have the following lemma.

Lemma 5. *The descent amount of Augmented Lagrangian function produced by one pass of AFW (Algorithm 2) (and FMO parameter ν) has*

$$E[\mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t)] - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) \leq -\frac{\mu_{\mathcal{M}}}{4(1+\nu)mQ}(\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)) \quad (33)$$

where $\mu_{\mathcal{M}}$ is the generalized geometric strong convexity constant for function $\mathcal{L}(\boldsymbol{\alpha})$ in domain \mathcal{M} , Q is the Lipschitz-continuous constant of $\nabla_{\boldsymbol{\alpha}}\mathcal{L}(\boldsymbol{\alpha})$ and $m = |\mathcal{F}|$.

Proof. Note the Augmented Lagrangian is of the form

$$H(\boldsymbol{\alpha}) := \mathcal{L}(\boldsymbol{\alpha}, \boldsymbol{\lambda}^t) = g(B\boldsymbol{\alpha}) + \langle \mathbf{b}, \boldsymbol{\alpha} \rangle \quad (34)$$

where

$$B := \begin{bmatrix} \Phi^T \\ M \end{bmatrix}, \quad \mathbf{b} := \boldsymbol{\delta} + M^T \boldsymbol{\lambda}^t$$

and function $g\left(\begin{bmatrix} \mathbf{w} \\ \mathbf{v} \end{bmatrix}\right) = \frac{1}{2}\|\mathbf{w}\|^2 + \frac{\rho}{2}\|\mathbf{v}\|^2 + \text{const.}$ is strongly convex with parameter $\bar{\rho} = \min(1, \rho)$. Without loss of generality, assume $\rho \leq 1$ and thus $\bar{\rho} = \rho$. Since we are minimizing the function subject

to a convex, polyhedral domain \mathcal{M} , by Theorem 10 of [8], we have the *generalized geometrical strong convexity* constant $\mu_{\mathcal{M}}$ of the form

$$\mu_{\mathcal{M}} := \mu(PWidth(\mathcal{M}))^2 \quad (35)$$

where $PWidth(\mathcal{M}) > 0$ is the pyramidal width of the simplex domain \mathcal{M} and μ is the *generalized strong convexity* constant of function (34) (defined by Lemma 9 of [8]). By definition of the geometric strong convexity constant, we have

$$H(\boldsymbol{\alpha}^t) - H^* \leq \frac{g_t^2}{2\mu_{\mathcal{M}}} \quad (36)$$

from (23) in [8], where $g_t := \langle -\nabla H(\boldsymbol{\alpha}^t), \mathbf{v}^F - \mathbf{v}^A \rangle$, and \mathbf{v}^F is the Frank-Wolfe direction

$$\mathbf{v}^F := \arg \min_{\mathbf{v} \in \mathcal{M}} \langle \nabla H(\boldsymbol{\alpha}^t), \mathbf{v} \rangle,$$

\mathbf{v}^A is the away direction

$$\mathbf{v}^A := \arg \max_{\mathbf{v} \in \mathcal{A}^t} \langle \nabla H(\boldsymbol{\alpha}^t), \mathbf{v} \rangle$$

Then let $m = |\mathcal{F}|$ be the number of factors. The FMO returns $\mathbf{v}_f^+ = \mathbf{v}_f^F$ with probability at least $\frac{1}{\nu}$, and suppose we set \mathbf{v}_f^+ to $\boldsymbol{\alpha}_f^t$ whenever $\langle \nabla_{\boldsymbol{\alpha}_f} H, \mathbf{v}_f^+ - \boldsymbol{\alpha}_f^t \rangle \not\leq 0$. We have $\langle \nabla H, \mathbf{d}_F \rangle \leq \frac{1}{\nu} \langle \nabla H, \mathbf{v}^F - \boldsymbol{\alpha}^t \rangle$ and thus

$$(1 + \frac{1}{\nu}) \langle \nabla H, \mathbf{d}^t \rangle \leq \frac{1}{\nu} \langle \nabla H, \mathbf{v}^F - \boldsymbol{\alpha}^t \rangle + \frac{1}{\nu} \langle \nabla H, \boldsymbol{\alpha}^t - \mathbf{v}^A \rangle$$

and $\langle \nabla H, \mathbf{d}^t \rangle \leq -\frac{1}{1+\nu} g_t$. Therefore, for any $\forall \gamma \in [0, 1]$,

$$E[H(\boldsymbol{\alpha}^{t+1})] - H(\boldsymbol{\alpha}^t) \leq -\gamma \frac{g_t}{1+\nu} + \frac{Q}{2} \|\gamma(\boldsymbol{\alpha}^{t+1} - \boldsymbol{\alpha}^t)\|^2 \leq -\gamma \frac{g_t}{1+\nu} + \frac{2mQ}{2} \gamma^2 \quad (37)$$

where Q is an upper bound on the spectral norm of Hessian $\|\nabla^2 H(\boldsymbol{\alpha})\|$ and $2m$ is the square of the radius of domain \mathcal{M} . Now we need to consider two cases. When the greedy direction \mathbf{d}_F in (17) is chosen, we have $\gamma^* = \min(\frac{g_s}{mQ_{\max}}, 1)$, which gives us

$$E[H(\boldsymbol{\alpha}^{t+1})] - H(\boldsymbol{\alpha}^t) \leq -\frac{g_t^2}{4(1+\nu)mQ}. \quad (38)$$

While in case \mathbf{d}_A in (17) is chosen, we have $\gamma^* = \min(\frac{g_s}{mQ_{\max}}, c_{v^-})$. When $\gamma^* = c_{v^-}$, a basis \mathbf{v}^- is removed from the active set and this is called a *drop step* [8] and it is hard to show sufficient descent in this case. Nevertheless, we can ignore those drop steps since the number of them is at most half of the iterates. For a non-drop step t , with the error bound (36), we have

$$E[H(\boldsymbol{\alpha}^{t+1})] - H(\boldsymbol{\alpha}^t) \leq -\frac{\mu_{\mathcal{M}}(H(\boldsymbol{\alpha}^t) - H^*)}{4(1+\nu)mQ}. \quad (39)$$

□

The above Lemma shows a significant progress made by the AFW algorithm. In the following, we provide a similar Lemma for minimizing AL subproblem with the Block-Greedy Coordinate Descent (BGCD) (Algorithm 3). Note that for problem of the form (34), the optimal solution is profiled by a polyhedral set $\mathcal{S} := \{\boldsymbol{\alpha} \mid B\boldsymbol{\alpha} = \mathbf{t}^*, \mathbf{b}^T \boldsymbol{\alpha} = s^*, \boldsymbol{\alpha} \in \mathcal{M}\}$. Therefore, let $\bar{\boldsymbol{\alpha}} := \Pi_{\mathcal{S}}(\boldsymbol{\alpha})$. We can bound the distance of any feasible point $\boldsymbol{\alpha} \in \mathcal{M}$ to its projection $\Pi_{\mathcal{S}}(\boldsymbol{\alpha})$ on \mathcal{S} using the Hoffman's inequality [4]

$$\|\bar{\boldsymbol{\alpha}} - \boldsymbol{\alpha}\|_{2,1}^2 = \sum_{i=1}^n \left(\sum_{f \in \mathcal{F}_i} \|\bar{\boldsymbol{\alpha}}_f - \boldsymbol{\alpha}_f\|_2 \right)^2 \leq \theta_1 (\|B\boldsymbol{\alpha} - \mathbf{t}^*\|^2 + \|\mathbf{b}^T \boldsymbol{\alpha} - s^*\|^2) \quad (40)$$

where θ_1 is a constant depending on the set \mathcal{S} . Then we can establish the following Lemma using the error bound (40).

Lemma 6. *The descent amount of Augmented Lagrangian function given by one pass of BGCD (Algorithm 3) with FMO multiplicative-approximation parameter ν has*

$$E[\mathcal{L}(\boldsymbol{\alpha}^{t+1}, \boldsymbol{\lambda}^t)] - \mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) \leq \frac{-1}{1 + \nu Q_{\max}/\mu_1} (\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)) \quad (41)$$

where

$$\mu_1 := \frac{1}{\max\{16\theta_1\Delta\mathcal{L}^0, 2\theta_1(1 + 4L_g^2)\}}.$$

is the generalized strong convexity constant for function $\mathcal{L}(\boldsymbol{\alpha})$ with feasible domain \mathcal{M} , $\Delta\mathcal{L}^0$ is a bound on $\mathcal{L}(\boldsymbol{\alpha}^0) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^0)$, L_g is the local Lipschitz-continuous constant of $g(\cdot)$ and $Q_{\max} = \max_{f \in \mathcal{F}} Q_f$.

Proof. For each iteration s of Algorithm 3, let i be the chosen sample and suppose that out of ν partitions the one containing greedy factor satisfying (18) is chosen. We have

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha}^{s+1}) - \mathcal{L}(\boldsymbol{\alpha}^s) &\leq \min_{\boldsymbol{\alpha}_{f^*}^s + \mathbf{d}_{f^*} \in \Delta^{|\mathcal{Y}_{f^*}|}} \langle \nabla_{\boldsymbol{\alpha}_{f^*}} \mathcal{L}, \mathbf{d}_{f^*} \rangle + \frac{Q_{\max}}{2} \|\mathbf{d}_{f^*}\|^2 \\ &= \min_{\boldsymbol{\alpha}_f^s + \mathbf{d}_f \in \Delta^{|\mathcal{Y}_f|}} \langle \nabla_{\boldsymbol{\alpha}_f} \mathcal{L}, \mathbf{d}_f \rangle + \frac{Q_{\max}}{2} \left(\sum_{f \in \mathcal{F}_i} \|\mathbf{d}_f\| \right)^2 \end{aligned} \quad (42)$$

where the second equality follows from the optimality of f^* w.r.t. (18). Then consider i being uniformly sampled from $[n]$, and consider the probability that the partition containing greedy factor f^* is chosen, the expected descent amount is

$$\begin{aligned} &E[\mathcal{L}(\boldsymbol{\alpha}^{s+1})] - \mathcal{L}(\boldsymbol{\alpha}^s) \\ &\leq \frac{1}{n\nu} \left(\min_{\boldsymbol{\alpha}_f^s + \mathbf{d}_f \in \Delta^{|\mathcal{Y}_f|}} \sum_{f \in \mathcal{F}} \langle \nabla_{\boldsymbol{\alpha}_f} \mathcal{L}, \mathbf{d}_f \rangle + \frac{Q_{\max}}{2} \sum_{i=1}^n \left(\sum_{f \in \mathcal{F}_i} \|\mathbf{d}_f\| \right)^2 \right) \\ &\leq \frac{1}{n\nu} \left(\min_{\boldsymbol{\alpha}_f^s + \mathbf{d}_f \in \Delta^{|\mathcal{Y}_f|}} \mathcal{L}(\boldsymbol{\alpha}^s + \mathbf{d}) - \mathcal{L}(\boldsymbol{\alpha}^s) + \frac{Q_{\max}}{2} \sum_{i=1}^n \left(\sum_{f \in \mathcal{F}_i} \|\mathbf{d}_f\| \right)^2 \right) \\ &\leq \frac{1}{n\nu} \left(\min_{\beta \in [0,1]} \mathcal{L}(\boldsymbol{\alpha}^s + \beta(\bar{\boldsymbol{\alpha}}^s - \boldsymbol{\alpha}^s)) - \mathcal{L}(\boldsymbol{\alpha}^s) + \frac{Q_{\max}\beta^2}{2} \sum_{i=1}^n \left(\sum_{f \in \mathcal{F}_i} \|\bar{\boldsymbol{\alpha}}_f^s - \boldsymbol{\alpha}_f^s\| \right)^2 \right) \\ &\leq \frac{1}{n\nu} \left(\min_{\beta \in [0,1]} \beta(\mathcal{L}(\bar{\boldsymbol{\alpha}}^s) - \mathcal{L}(\boldsymbol{\alpha}^s)) + \frac{Q_{\max}\beta^2}{2} \|\bar{\boldsymbol{\alpha}}^s - \boldsymbol{\alpha}^s\|_{2,1}^2 \right) \end{aligned} \quad (43)$$

where $\bar{\boldsymbol{\alpha}}^s = \Pi_{\mathcal{S}}(\boldsymbol{\alpha}^s)$ is the projection of $\boldsymbol{\alpha}^s$ to the optimal solution set \mathcal{S} . The second and last inequality is due to convexity, and the third inequality is due to a confinement of optimization domain. Then we discuss two cases in the following.

Case 1: $4L_g^2\|B\boldsymbol{\alpha}^s - \mathbf{t}^*\|^2 < (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)^2$.

In this case, by the Hoffman inequality (40), we have

$$\begin{aligned} \|\boldsymbol{\alpha}^s - \bar{\boldsymbol{\alpha}}^s\|_{2,1}^2 &\leq \theta_1 (\|B\bar{\boldsymbol{\alpha}}^s - \mathbf{t}^*\|^2 + (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)^2) \\ &\leq \theta_1 \left(\frac{1}{4L_g^2} + 1 \right) (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)^2 \\ &\leq 2\theta_1 (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)^2, \end{aligned} \quad (44)$$

since $\frac{1}{4L_g^2} \leq 1$. Then

$$|\mathbf{b}^T \boldsymbol{\alpha}^s - s^*| \geq 2L_g \|B\boldsymbol{\alpha}^s - \mathbf{t}^*\| \geq 2|g(B\boldsymbol{\alpha}^s) - g(\mathbf{t}^*)|$$

by the definition of Lipschitz constant L_g . Note that $\mathbf{b}^T \boldsymbol{\alpha}^s - s^*$ is non-negative since otherwise we have contradiction $\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^* = g(B\boldsymbol{\alpha}^s) - g(\mathbf{t}^*) + (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*) \leq |g(B\boldsymbol{\alpha}^s) - g(\mathbf{t}^*)| - |\mathbf{b}^T \boldsymbol{\alpha}^s - s^*| \leq -\frac{1}{2}|\mathbf{b}^T \boldsymbol{\alpha}^s - s^*| < 0$. Therefore, we have

$$\begin{aligned} \mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^* &= g(B\boldsymbol{\alpha}^s) - g(\mathbf{t}^*) + (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*) \\ &\geq -|g(B\boldsymbol{\alpha}^s) - g(\mathbf{t}^*)| + (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*) \\ &\geq \frac{1}{2}(\mathbf{b}^T \boldsymbol{\alpha}^s - s^*). \end{aligned} \quad (45)$$

Combining (43), (44) and (45), we have

$$\begin{aligned} &\mathbb{E}[\mathcal{L}(\boldsymbol{\alpha}^{s+1})] - \mathcal{L}(\boldsymbol{\alpha}^s) \\ &\leq \frac{1}{n\nu} \left(\min_{\beta \in [0,1]} -\frac{\beta}{2}(\mathbf{b}^T \boldsymbol{\alpha}^s - s^*) + \frac{2\theta_1 Q_{\max} \beta^2}{2}(\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)^2 \right) \\ &= \begin{cases} -1/(16\theta_1 Q_{\max} n\nu) & , 1/(4\theta_1 Q_{\max}(\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)) \leq 1 \\ -\frac{1}{4n\nu}(\mathbf{b}^T \boldsymbol{\alpha}^s - s^*) & , o.w. \end{cases} \end{aligned}$$

Furthermore, we have

$$-\frac{1}{16Q_{\max}\theta_1 n\nu} \leq -\frac{1}{16Q_{\max}\theta_1 n\nu(\mathcal{L}^0 - \mathcal{L}^*)} (\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^*) \quad (46)$$

where $\mathcal{L}^0 = \mathcal{L}(\boldsymbol{\alpha}^0)$, and

$$-\frac{1}{4n\nu}(\mathbf{b}^T \boldsymbol{\alpha}^s - s^*) \leq -\frac{1}{6n\nu}(\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^*) \quad (47)$$

since $\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^* \leq |g(B\boldsymbol{\alpha}^s) - g(\mathbf{t}^*)| + \mathbf{b}^T \boldsymbol{\alpha}^s - s^* \leq \frac{3}{2}(\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)$. Since the bound (46) is much smaller than (47). For Case 1, we obtain

$$\mathbb{E}[\mathcal{L}(\boldsymbol{\alpha}^{s+1})] - \mathcal{L}^s \leq -\frac{\mu_1}{n\nu Q_{\max}} (\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^*) \quad (48)$$

where

$$\mu_1 = \frac{1}{16\theta(\mathcal{L}^0 - \mathcal{L}^*)}. \quad (49)$$

Case 2: $4L_g^2 \|B\boldsymbol{\alpha}^s - \mathbf{t}^*\|^2 \geq (\mathbf{b}^T \boldsymbol{\alpha}^s - s^*)^2$.

In this case, we have

$$\|\boldsymbol{\alpha}^s - \bar{\boldsymbol{\alpha}}^s\|_{2,1}^2 \leq \theta_1 (1 + 4L_g^2) \|B\boldsymbol{\alpha}^s - \mathbf{t}^*\|^2, \quad (50)$$

and by strong convexity of $g(\cdot)$,

$$\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^* \geq \mathbf{b}^T (\boldsymbol{\alpha}^s - \boldsymbol{\alpha}^*) + \nabla g(\mathbf{t}^*)^T B(\boldsymbol{\alpha}^s - \bar{\boldsymbol{\alpha}}^s) + \frac{\rho}{2} \|B\boldsymbol{\alpha}^s - \mathbf{t}^*\|^2.$$

Now let $h(\boldsymbol{\alpha})$ be a function that takes value 0 when $\boldsymbol{\alpha}$ is feasible and takes value ∞ otherwise. Adding inequality $0 = h(\boldsymbol{\alpha}^s) - h(\bar{\boldsymbol{\alpha}}^s) \geq \langle \boldsymbol{\sigma}^*, \boldsymbol{\alpha}^s - \bar{\boldsymbol{\alpha}}^s \rangle$ to the above gives

$$\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^* \geq \frac{\rho}{2} \|B\boldsymbol{\alpha}^s - \mathbf{t}^*\|^2 \quad (51)$$

because $\boldsymbol{\sigma}^* + \mathbf{b} + \nabla g(\mathbf{t}^*)^T B = \boldsymbol{\sigma}^* + \nabla \mathcal{L}(\bar{\boldsymbol{\alpha}}^s) = 0$ for some $\boldsymbol{\sigma}^* \in \partial h(\bar{\boldsymbol{\alpha}}^s)$. Combining (43), (50), and (51), we obtain

$$\begin{aligned} & \mathbb{E}[\mathcal{L}(\boldsymbol{\alpha}^{s+1})] - \mathcal{L}(\boldsymbol{\alpha}^s) \\ & \leq \frac{1}{n\nu} \left(\min_{\beta \in [0,1]} -\beta(\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^*) + \frac{\theta_1(1+4L_g^2)Q_{\max}\beta^2}{\rho} (\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^*) \right) \\ & \leq -\frac{\rho}{n\nu\theta_1(1+4L_g^2)Q_{\max}} (\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^*) \end{aligned} \quad (52)$$

Combining results of Case 1 (48) and Case 2 (52), and taking expectation on both sides w.r.t. the history, we have

$$E[\mathcal{L}(\boldsymbol{\alpha}^{s+1})] - \mathcal{L}(\boldsymbol{\alpha}^s) \leq -\frac{\mu_1}{Q_{\max}n\nu} (\mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^*). \quad (53)$$

where

$$\mu_1 := \min\left\{\frac{1}{16\theta(\Delta\mathcal{L}^0)}, \frac{\rho}{\theta_1(1+4L_g^2)}\right\}.$$

Taking summation of (53) over iterates $s = 1 \dots n$, we have

$$\begin{aligned} E[\mathcal{L}(\boldsymbol{\alpha}^{t+1})] - \mathcal{L}(\boldsymbol{\alpha}^t) & \leq -\frac{\mu_1}{Q_{\max}n\nu} \left(\sum_{s=1}^n \mathcal{L}(\boldsymbol{\alpha}^s) - \mathcal{L}^* \right) \\ & \leq -\frac{\mu_1}{Q_{\max}\nu} (\mathcal{L}(\boldsymbol{\alpha}^{t+1}) - \mathcal{L}^*). \end{aligned} \quad (54)$$

Rearranging terms gives the conclusion. \square

Now we provide proof of Theorem 1 as follows.

Proof. Let $\kappa = 4(1+\nu)mQ/\mu_{\mathcal{M}}$. By lemma 3, 5, 4 and (30), we have

$$\begin{aligned} & \Delta_d^t - \Delta_d^{t-1} + E[\Delta_p^t] - \Delta_p^{t-1} \\ & \leq \frac{-1}{1+\kappa} (\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)) + \frac{2\eta}{\rho} (\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)) - \frac{\eta}{\tau} \Delta_d^t. \end{aligned} \quad (55)$$

Then by choosing $\eta < \frac{\rho}{2(1+\kappa)}$, we have guaranteed descent on $\Delta_p + \Delta_d$ for each GDM iteration. By choosing $\eta \leq \frac{\rho}{4(1+\kappa)}$, we have

$$\begin{aligned} & (\Delta_d^t + E[\Delta_p^t]) - (\Delta_d^{t-1} + \Delta_p^{t-1}) \\ & \leq \frac{-1}{2(1+\kappa)} (\mathcal{L}(\boldsymbol{\alpha}^t, \boldsymbol{\lambda}^t) - \mathcal{L}(\bar{\boldsymbol{\alpha}}^t, \boldsymbol{\lambda}^t)) - \frac{\eta}{\tau} \Delta_d^t \\ & \leq -\min\left(\frac{1}{2(1+\kappa)}, \frac{\eta}{\tau}\right) (\Delta_p^t + \Delta_d^t) \end{aligned}$$

which then leads to the conclusion. \square

The proof for Theorem 2 follows the same line of above reasoning with step (55) replaced by application of Lemma 6 instead of Lemma 5.

7 Appendix C: Implementation details of FMO

7.1 C-1: Indicator Factor

Here we assume $\delta(y_j, \bar{y}_j)$ is constant for $\forall y_j \neq \bar{y}_j$ as in the case of Hamming error. Then we find maximizers of the 4 cases as following

- (i) Visit \mathbf{y}_f in descending order of $v(\cdot)$ to find the first $\mathbf{y}_f: m_{if}(y_i) = 0, m_{jf}(y_j) = 0$.
- (ii) $\forall y_j: m_{jf}(y_j) \neq 0$, visit y_i in descending order of $v(\cdot)$ to find the first $y_i: m_{if}(y_i) = 0$.
- (iii) $\forall y_i: m_{if}(y_i) \neq 0$, visit y_j in descending order of $v(\cdot)$ to find the first $y_j: m_{jf}(y_j) = 0$.
- (iv) Evaluate (26) for $\forall (y_i, y_j): m_{if}(y_i) \neq 0, m_{jf}(y_j) \neq 0$.

Then \mathbf{y}_f^* is returned as label ($\neq \bar{\mathbf{y}}_f$) of maximum gradient (26) among the 4 cases. One can verify the above procedure considers all labels that have potential to be \mathbf{y}_f^* . The complexities for (ii)-(iv) are bounded by $O(\text{nnz}(\mathbf{m}_{if})\text{nnz}(\mathbf{m}_{jf}))$, where $\text{nnz}(\mathbf{m}_{jf}) \leq |\hat{\mathcal{A}}_j^t|$. When BCFW adopts sampling without replacement, we have $|\hat{\mathcal{A}}_f^t| \leq t$. In practice, as t keeps increasing, $|\hat{\mathcal{A}}_f^t|$ converges to a constant that depends on the optimal $\text{nnz}(\boldsymbol{\alpha}_f^*)$. Note $\text{nnz}(\boldsymbol{\alpha}_f^*)$ is equivalent to the number of labels \mathbf{y}_f that attains the maximum of hinge loss (8), which is small in general as long as there are few labels with larger responses than the others.

Define $\mathcal{Y}_{NZ} = \{\mathbf{y}_f | m_{if}(y_i) \neq 0 \vee m_{jf}(y_j) \neq 0\}$ as the set of labels with messages from one of the variables involved, and $\mathcal{Y}_{Inc} = \{\mathbf{y}_f | \mathbf{y}_f \in \mathcal{Y}_{NZ} \wedge v(\mathbf{y}_f, \mathbf{x}_f) > v(\mathbf{y}'_f, \mathbf{x}_f), \forall \mathbf{y}'_f \notin \mathcal{Y}_{NZ}\}$ as the subset being inconsistently ranked at the top in the multimap. The complexity of step (i) is $O(|\mathcal{Y}_{Inc}|)$, where

$$|\mathcal{Y}_{Inc}| \leq \max(|\mathcal{Y}_i| |\hat{\mathcal{A}}_j^t|, |\mathcal{Y}_j| |\hat{\mathcal{A}}_i^t|), \quad (56)$$

which is sublinear to the size of factor domain $|\mathcal{Y}_f| = |\mathcal{Y}_i| |\mathcal{Y}_j|$. Although the bound (56) is already sublinear to $|\mathcal{Y}_f|$, it is a *very loose bound*. In our experiments, we observed the average number of elements being visited at stage (i) is no more than 5 for problems of $|\mathcal{Y}_f|$ up to 10^7 , presumably because the inconsistency between factors is small in real applications.

7.2 C-2: Binary-Variable Interaction Factor

Similar to Appendix C-1, we're trying to find active factors with largest gradient. Here is the procedure.

- (i) Visit \mathbf{y}_f in descending order of $v(\cdot)$ to find the first $\mathbf{y}_f: i \notin \mathcal{A}, j \notin \mathcal{A}$.
- (ii) $\forall j: j \notin \mathcal{A}$, visit i in descending order of $v(\cdot)$ to find the first $i: i \notin \mathcal{A}$.
- (iii) $\forall i: i \notin \mathcal{A}$, visit j in descending order of $v(\cdot)$ to find the first $j: j \notin \mathcal{A}$.
- (iv) Compute gradient for $\forall (i, j): i \in \mathcal{A}, j \in \mathcal{A}$.

A similar reasoning as C-1 applies here for complexity analysis.