

Scalable optimization algorithms for the Structured SVM

Frederic Boileau Elyes Lamouchi William St-Arnaud

28th April 2019

1 Introduction

Structured prediction in machine learning is tasked with learning predictors where the labels on the data-points are more than simple tags but have inherent structure which implies the following:

- The number of potential labels for a given feature vector can grow exponentially with the input which makes traditional classification procedures intractable
- A certain intelligibility of the structure; hence a hope to leverage it to improve tractability

It is often quite hard to know in advance whether we can tackle a structured prediction problem with known approaches.

In this paper we discuss two different approaches to solving structured prediction problems, both focusing on a *large-margin approach* which translates into a support-vector machine - like problem formulation.

Let us introduce some notation which we will use throughout the paper. Define the dataset to be

$$S = \{x^{(i)}, y^{(i)}\}_{i=1}^n \in (\mathcal{X} \times \mathcal{Y})^n \quad (1)$$

If one were to frame the machine learning goal in the concisest and simplest way possible we could say that the goal is to *learn* a (parametrized) *predictor* function:

$$\begin{aligned} h_w &: \mathcal{X} \rightarrow \mathcal{Y} \\ &: \hat{x} \mapsto y \quad y \in \mathcal{Y}(\hat{x}) \end{aligned}$$

where \hat{x} is just some arbitrary sampled \hat{x} and we abuse notation and mean the set valued mapping which outputs the feasible label set for a given x by $\mathcal{Y}(x)$

The traditional probabilistic approach is to compute h or its parameter by calculating the most likely parameter given conditioned on the observed data. However this is more often than not untractable in the structured prediction context as summed up in Taskar, Lacoste-Julien, and Jordan (2006)

In a large margin approach we wish to compute the predictor function the following way :

$$h_w(x) = \arg \max_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad (2)$$

where ϕ is just the feature map for the dataset.

This is a constrained optimization problem and the structure of \mathcal{Y} clearly has a big influence on how well we can solve the problem as well as which method should work well.

The dual extragradient approach presented in Taskar, Lacoste-Julien, and Jordan (ibid.) leverages a saddle-point approach to the problem to tackle the problem with a first order method and moreover “this approach allows us to exploit the structure of W and Z separately, allowing for efficient solutions for a wider range of parameterizations and structures.”(ibid.)

Despite its improvements at the time the dual extragradient approach suffers from two main computational draw-backs:

1. It is not seperable/'stochasizable' which is a severe problem as machine learning has seen much of its success in recent years by optimizing through first-order sampled method, i.e. first order method which use a probabilistic approximation of the gradient through sampling instead of computing the actual gradient (which might be unfeasible). TODO : cite mr simon
2. The algorithm requires a projection step which can be expensive depending on the structure of \mathcal{Y} .
TODO : cite explain why proj is expensive

This leads us to consider the Block Seperable Frank Wolfe algorithm presented in Lacoste-Julien et al. (2013)

2 Preliminaries

2.1 Structured SVM N-slack formulation

Let us recall the basic goal; to construct an accurate linear classifier¹

$$h_w(x) = \arg \max_{y \in \mathcal{Y}} \langle w, \phi(x, y) \rangle \quad (3)$$

In general finding the optimal separating hyperplane is an ill-defined problem. In the support vector machine (SVM) setting we want to find the separating hyperplane with the *largest margin*. The support vectors are the datapoints which lie on the the defined margins. It is not obvious that we should strive to completely separate the sets so we can include some slack variables. Indeed, some set of points might not be linearly separable but adding slack variables might allow one to misclassify some point but get a higher proportion of correctly classified points.

We call the following way to pose problem the n -slack formulation ² of the problem :

$$\max_{w, \xi} \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \varepsilon_i \quad (4)$$

$$s.t. \quad \langle w, \psi_i(y) \rangle \geq L(y_i, y) - \varepsilon_i, \quad \forall i, \forall y \in \mathcal{Y}(x) = \mathcal{Y}_i \quad (5)$$

2.2 Loss-Augmented Decoding

Define $\psi(y) := \phi(x^{(i)}, y^{(i)}) - \phi(x^{(i)}, y)$ and $L_i(y) = L(y^{(i)}, y)$

$$H_i = \max_{y \in \mathcal{Y}_i} \{L_i(y) - \langle w, \psi_i(y) \rangle\} \quad (6)$$

$$\mathcal{H}_i = \arg \max_{y \in \mathcal{Y}_i} \{L_i(y) - \langle w, \psi_i(y) \rangle\} \quad (7)$$

Proposition. The max oracle is a convex upper bound to the task loss.

Proof sketch: The maximum of two convex (linear) functions is convex, and

$$L(y_i, h_w(x_i)) \leq L(y_i, h_w(x_i)) + \underbrace{\langle w, \psi_i(y) \rangle}_{\geq 0 \text{ by definition}} \quad (8)$$

$$\leq \max_{y \in \mathcal{Y}_i} L_i(y) - \langle w, \psi_i(y) \rangle \quad (9)$$

It is not hard to see thus that learning w amounts to the unconstrained problem,

$$\max_w \quad \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \tilde{H}_i(w) \quad (10)$$

¹clearly we can easily extend to non-linear cases, with kernel maps for example

²see Moguerza and Muñoz (2006) for example

2.3 Convex Analysis etc

2.3.1 Proximal step operator

We define the proximal step operator as follows:

$$T_\eta(\mathbf{u}, \mathbf{s}) = \max_{\mathbf{u}' \in U} \left\{ \langle \mathbf{s}, \mathbf{u}' - \mathbf{u} \rangle - \frac{1}{\eta} d(\mathbf{u}, \mathbf{u}') \leq D \right\} \quad (11)$$

The operator is useful to compute projections since when we have a strongly convex function $h(\mathbf{u})$, we can find its convex conjugate $h^*(\mathbf{u}) = \max_{\mathbf{u} \in U} [\langle \mathbf{s}, \mathbf{u} \rangle - h(\mathbf{u})]$. From the definition of a strongly convex function, we have that:

$$h(\mathbf{u}') \geq h(\mathbf{u}) + \langle \nabla h(\mathbf{u}), \mathbf{u}' - \mathbf{u} \rangle + \frac{\sigma}{2} \|\mathbf{u}' - \mathbf{u}\|^2 \quad (12)$$

where σ is the strong convexity parameter. Rearranging, we can define an upper bound on the squared norm of $\mathbf{u}' - \mathbf{u}$. This comes out as:

$$d(\mathbf{u}', \mathbf{u}) \triangleq h(\mathbf{u}') - h(\mathbf{u}) - \langle \nabla h(\mathbf{u}), \mathbf{u}' - \mathbf{u} \rangle \geq \frac{\sigma}{2} \|\mathbf{u}' - \mathbf{u}\|^2 \quad (13)$$

The distance metric d is called the Bregman divergence. The link between the Bregman divergence and the proximal step operator is that if we are given the function h inside the definition of the proximal step update, this induces the Bregman divergence, which in turn induces the update that is performed at each iteration of the extragradient algorithm. For example, if we have $h(\mathbf{u}) = \frac{1}{2} \|\mathbf{u}\|_2^2$, the Bregman divergence becomes $d(\mathbf{u}', \mathbf{u}) = \frac{1}{2} \|\mathbf{u}' - \mathbf{u}\|_2^2$. We might wonder why we care about the Bregman divergence when the definition still includes the usual norm. After all, we still optimize the term $\langle \mathbf{s}, \mathbf{u}' - \mathbf{u} \rangle - \frac{1}{\eta} d(\mathbf{u}', \mathbf{u})$. This is because h^* is differentiable at every point of its domain by the strong convexity of h . Thus, it is easy to compute a projection in the usual fashion: we can compute the derivative of the term inside the projection operator and set it to 0. It is impossible to do for matchings for example as the distance is not even differentiable. We provide the steps to compute a projection:

$$\mathbf{s} - \nabla_{\mathbf{u}'} d(\mathbf{u}', \mathbf{u}) = \mathbf{s} - \frac{1}{\eta} \nabla_{\mathbf{u}'} d(\mathbf{u}, \mathbf{u}') = \mathbf{s} - \frac{1}{\eta} [\nabla h(\mathbf{u}') - \nabla h(\mathbf{u})] \quad (14)$$

By setting this equation to 0, it is possible to recover the optimal \mathbf{y}' when, let's say, $h(\mathbf{u}) = \frac{1}{2} \|\mathbf{u}'\|^2$.

3 Extra Gradient Algorithm

3.1 Saddle-Point

$$\min_{\mathbf{w} \in W} \max_{\mathbf{z} \in Z} \sum_i (\mathbf{w}^T \mathbf{F}_i \mathbf{z}_i + \mathbf{c}_i^T \mathbf{z}_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i)) \quad (15)$$

where the \mathbf{z}_i 's can be identified with the edge and node potentials of a markov network and satisfy the constraints of the structured problem. The terms \mathbf{F}_i correspond to the feature mapping for over all labels \mathbf{y}_i when multiplied by \mathbf{z}_i 's. The \mathbf{c}_i 's correspond to the costs of a \mathbf{z}_i and can be identified with the loss l for a label \mathbf{y}'_i . Taking the dual, we end up with the following:

$$\begin{aligned} \min_{\mathbf{w} \in W, (\lambda, \mu) \geq \mathbf{0}} \quad & \sum_i (\mathbf{b}_i^T \lambda_i + \mathbf{1}^T \mu_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i)) \\ \text{s.t.} \quad & \mathbf{F}_i^T \mathbf{w} + \mathbf{c}_i \leq \mathbf{A}_i^T \lambda_i + \mu_i \quad i = 1, \dots, m \end{aligned} \quad (16)$$

The number of variables and constraints is linear in the number of paramters and training data. We already see that this formulation is much more efficient. We do have a set of constraints that is tractable, as is the number of parameters to update.

In equation 15, the term that is opitmized is defined as:

$$L(\mathbf{w}, \mathbf{z}) \triangleq \sum_i \mathbf{w}^T \mathbf{F}_i \mathbf{z}_i + \mathbf{c}_i^T \mathbf{z}_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i) \quad (17)$$

It is bilinear in w and z . We can then imagine two players represented by \mathbf{w} and \mathbf{z} that play a zero-sum game. They perform updates using gradients of the objective w.r.t. their parameters. They then project the result to the set of feasible points given by the constraints imposed on the structure. We usually consider Euclidean projections, as there are well-known problems where they are efficient to compute. However, as seen later, this is not the case for all problem. This is why Bregman projections will be introduced. Going back to the zero-sum game, we have the following operator that is used to perform the updates for both players at the same time.

$$\begin{pmatrix} \nabla_{\mathbf{w}} L(\mathbf{w}, \mathbf{z}) \\ -\nabla_{\mathbf{z}_1} L(\mathbf{w}, \mathbf{z}) \\ \vdots \\ -\nabla_{\mathbf{z}_m} L(\mathbf{w}, \mathbf{z}) \end{pmatrix} = \underbrace{\begin{pmatrix} 0 & \mathbf{F}_1 & \dots & \mathbf{F}_m \\ -\mathbf{F}_1^T & & & \\ \vdots & & \mathbf{0} & \\ -\mathbf{F}_m^T & & & \end{pmatrix}}_{\mathbf{F}} \underbrace{\begin{pmatrix} \mathbf{w} \\ \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_m \end{pmatrix}}_{\mathbf{u}} - \underbrace{\begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix}}_{\mathbf{a}} = \mathbf{F}\mathbf{u} - \mathbf{a} \quad (18)$$

We can measure the “goodness” of the parameters using the gap function G :

$$G(\mathbf{w}, \mathbf{z}) \triangleq \left[\max_{\mathbf{z}' \in Z} L(\mathbf{w}, \mathbf{z}') - L^* \right] + \left[L^* - \min_{\mathbf{w}' \in W} L(\mathbf{w}', \mathbf{z}) \right] \quad (19)$$

where L^* gives the result of the min-max of the objective L . When we have a non-optimal point (i.e. not a saddle point), the gap is strictly positive. At an optimal point, the gap is exactly equal to 0. Now

the restricted gap is exactly the same but the min and max are computed over a set of parameters that are within a certain distance of the start point $(\hat{\mathbf{u}}_{\mathbf{w}}, \hat{\mathbf{u}}_{\mathbf{z}}) \in U$:

$$G_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\mathbf{w}, \mathbf{z}) = \max_{\mathbf{z}' \in Z} [L(\mathbf{w}', \mathbf{z}') : d(\mathbf{z}, \mathbf{z}') \leq D_{\mathbf{z}}] - \left[\min_{\mathbf{w}' \in W} L(\mathbf{w}', \mathbf{z}) : d(\mathbf{w}, \mathbf{w}') \leq D_{\mathbf{w}'} \right] \quad (20)$$

The motivation for using this restricted gap function is that if we start “close” to an optimal point, of course we will converge more rapidly to it. This can be seen in the convergence analysis of the method.

We mainly follow the intuition and proofs of Taskar, Lacoste-Julien, and Jordan, n.d.

The dual extragradient algorithm from Nesterov gives a convergence guarantee for the objective L .

We present a simple formulation of the algorithm:

Initialize: Choose $\hat{\mathbf{u}} \in U$, set $\mathbf{s}^{-1} = 0$.

for $t = 0$ to $t = \tau$ **do**

$\mathbf{v} = \Pi_U(\hat{\mathbf{u}} + \eta \mathbf{s}^{t-1})$

$\mathbf{u}^t = \Pi_U(\mathbf{v} - \eta(\mathbf{F}\mathbf{v} - \mathbf{a}))$

$\mathbf{s}^t = \mathbf{s}^{t-1} - (\mathbf{F}\mathbf{u}^t - \mathbf{a})$

end for

return $\overline{\mathbf{u}}^\tau = \frac{1}{1+\tau} \sum_{t=0}^{\tau} \mathbf{u}^t$

This algorithm has a lookahead step (i.e. \mathbf{v}) that serves to perform the actual gradient update \mathbf{u}^t . The intuition behind the lookahead step is that given a function to optimize that is Lipschitz, Nesterov was able to show that we can upper bound $f_D(\bar{u}^n) = \max_y \{\langle g(y), \bar{u}^n - y \rangle : d(\hat{u}, y) \leq D\}$, where \bar{u}^n is the weighted average over all the updates \mathbf{u}^t up to iteration n . The function g corresponds to the objective L in our setting. When value of $f_D(\bar{u}^n)$ gets close to 0, we have that the value $g(y^*)$ for an optimal y^* is close to 0, which signifies that we have reached saddle point (i.e. what we wanted). Nesterov goes on to show that this upper bound indeed goes to 0. We then get convergence to a saddle point. Note that in the definition of f_D , we used a distance metric d . This corresponds to the Euclidean distance (or Bregman distance in non-Euclidean setting). The rojection operator Π_U in the algorithm simply projects a point back to the set U by finding the nearest point with respect to the distance metric used.

3.1.1 Non-Euclidean setting

The main problem with the Euclidean projection operator is that for many problems, it is hard to compute the projection. Indeed for min-cut, we need to compute the partition function first, which is #P-complete. Thus, the authors of the paper introduced the Bregman operator, which computes the projection using the Bregman divergence. Using this operator has the great advantage of being easier to compute. We can see this for $L1$ regularization. Computing a projection using $L1$ distance is hard since it is not differentiable. Using the negative entropy, we get that the Bregman divergence is the KL divergence. This implies that we can differentiate the divergence to get the parameter that minimizes it.

3.1.2 Memory-efficient tweak

In the dual extragradient algorithm, both a vector s^t and a vector \bar{u}^t are maintained. However, we can observe that the s_t 's can be found using the running average \bar{u}^t since $s^t = -(t+1) \sum_{i=0}^t (F\bar{u}^i - a)$. We only have to store the vector \bar{u}^t . We can even do better when $|Z| \gg |W|$ since $\bar{u}^t = \{\bar{u}_w^t, \bar{u}_z^t\}$ and we only care about the part that corresponds to w . \bar{u}_z^t is maintained implicitly by storing a vector of size $|W|$ (although we now need to store s_w^t). It can be reconstructed using \bar{u}_w^t . The following figure (**ADD FIGURE 5 FROM PAPER**) illustrates the various dependencies.

4 Frank Wolfe and Variants

4.1 Frank-Wolfe, the conditional gradient algorithm

4.2 Block Coordinate Frank Wolfe

4.2.1 BCFW variant in the structured SVM setting

Due to the exponential number of dual variables in the structured SVM setting, classical algorithms, like projected gradient are intractable.

Stochastic subgradient methods, on the other hand, achieve a sublinear convergence rate while only requiring a single call to the maximization oracle every step. They are nonetheless very sensitive to the sequence of stepsizes and it is unclear when to terminate the iterations.

Frank-Wolfe methods address these problems by giving an adaptive stepsize $\gamma = \frac{2}{k+2}$ and a computable duality gap while still retaining a sublinear convergence rate. Moreover, despite the exponential number of constraints, the algorithm has sparse iterates alleviating the memory issues which come with the exponential number of dual variables.

Note. The main idea here, is that the linear subproblem in Frank-Wolfe and the loss augmented decoding of the structured SVM are equivalent.

Proof of the equivalence. The objective function being differentiable and convex, if we are at a point α such that $f(\alpha)$ is minimized along each coordinate axis, then α is a global minimizer. Therefore,

$$\min_{s \in \mathcal{M}} \langle s, \nabla f(\alpha) \rangle = \sum_i \min_{s_i \in \Delta_{|\mathcal{Y}_i|}} \langle s_i, \nabla_i f(\alpha) \rangle$$

Moreover, with

$$w = A\alpha, A = \left[\frac{1}{n\lambda} \psi_1(y) \dots \frac{1}{n\lambda} \psi_{\sum_i |\mathcal{Y}_i|}(y) \right]$$

$$\text{and } b = \left(\frac{1}{n} L_i(y) \right)_{i \in [n], y \in \mathcal{Y}_i}$$

The gradient of the dual would be,

$$\begin{aligned} \nabla f(\alpha) &= \nabla \left[\frac{\lambda}{2} \|A\alpha\|^2 - b^T \alpha \right] = \lambda A^T A \alpha - b \\ &= \lambda A^T w - b = \frac{1}{n} H_i(y, w) \\ \max_{y_i \in \mathcal{Y}_i} \tilde{H}_i &= - \min_{y_i \in \mathcal{Y}_i} \tilde{H}_i = \min_{y_i \in \mathcal{Y}_i} L_i - \langle w, \psi_i \rangle \\ &= \min_{s_i \in \Delta_{|\mathcal{Y}_i|}} \langle s_i, \nabla_i f(\alpha) \rangle \end{aligned}$$

Thus we can see that, if n = size of the training data, one Frank-Wolfe step is equivalent to n calls to the maximization oracle.

Let $\alpha \in \mathcal{M}$

Let $w^0 = 0, l^0 = 0$

for $k = 0, \dots, K$ **do**

for $i = 1, \dots, n$ **do**

 Solve $y_i^* = \max_{y_i \in \mathcal{Y}_i} H_i(y, w^k) //$

 Let $w_s = \sum_{i=1}^n \frac{1}{n\lambda} \psi_i(y_i^*)$, and $l_s = \frac{1}{n} \sum_{i=1}^n L_i(y_i^*)$

 Let $\gamma = \frac{\lambda(w^k - w_s)^T w^k - l^k + l_s}{\lambda \|w^k - w_s\|^2}$, and clip to $[0, 1]$

 Update $w^{k+1} = (1 - \gamma)w^k + \gamma w_s$, and $l^{k+1} = (1 - \gamma)l^k + \gamma l_s$

end for

end for

Unlike stochastic subgradient and stochastic methods in general, classical Frank-Wolfe requires one call for each training example at each iteration. For large datasets, this can get unpractical.

Hence the stochastic variant of Frank Wolfe, **Block Coordinate Frank Wolfe (BCFW)**.

5 Empirical

5.1 Experiments

In this section, we describe the implementation we performed for this project. The goal was to see how the Dual Extragradient algorithm compared to the Block-coordinate Frank-Wolfe algorithm. The task on which we performed the evaluation was word alignment in machine translation. We extracted the dataset from the Europarl dataset **TODO CITE**. The data consisted of approximately 2 million sentence pairs in both english and french. Each the sentences in each pair were translations of one another. We extracted all sentences and performed a clean by splitting longer sentences into shorter ones. The goal of this step was to reduce the eventual number of matchings in training, which could take long to solve using a LP solver. Of course, each sentence was tokenized beforehand.

We then proceeded to implement the Dual Extragradient and BCFW algorithms using a SVM. We had to define a feature mapping for an input sentence pair. The features were extracted using the fastText library **TODO CITE**. This library included a model that was previously trained to learn embeddings of words in both english and french. We later combined the embeddings in the two languages by applying a transformation found in **TODO CITE**. This transformation consisted in applying a matrix to each vector in each language (matrix W for english and Q for french). These matrices are in fact orthogonal (i.e. $Q^T W = I$). The idea behind such a transformation is that we sort of “put” or “align” both languages in the same vector space, a sort of “middle ground”. This way we can better compare the words “cat” and “chat” by getting their cosine similarity measure. We combined the cosine measure of each pair of words in the alignment by summing. The following blog post **CITE THIS** provides a good intuition using maps that are aligned. As an example, consider the following two sentences:

- This assignment was hard
- Ce travail etait ardu

We would compute the cosine distance for each word tuple (e.g. this/ce, this/travail, ..., hard/ardu). We were then able to get the highest match of each english word for french translation. This was how we extracted the “labels”. As the dataset was not annotated with alignments, we had to compute those according to the procedure mentioned.

For the features, we used the concatenated embeddings of each word pairs in the alignment. This gave us our edge score. Then, we simply performed a weighted combination of these vectors using the edge labels as weights. To clarify what we mean by edge labels, suppose that the edge linking “assignment” and “travail” has a value of 1. Then, we weight the vector extracted from these two words with 1. As another example, if the edge between “ce” and “ardu” has a value of 0 (i.e. no link between the words), we do not include the vector computed from the statistics of this word pair. This is exactly what was done in Taskar **CITE THIS** modulo some other features.

In the implementation of BCFW, we used the solver from scipy **CITE THIS SHIT** with the simplex method. Since the constraint matrix given by the optimization of H_i in the algorithm is unimodular, when we relax the LP, we still get a solution to the ILP without relaxation. Thus, we take advantage of this fact and indeed use the LP solver. The loss that we used was the L_1 distance between the two labels, the proposal and the ground truth.

5.1.1 Results

Since running the experiments was computationally intensive and we did not dispose of a lot of computing power, we had to restrict the training set to a number a relatively small number of sentence pairs (100). This sanity check was simply a hint to the general applicability of the method as we scale up the number of training examples. We used the default λ value of 0.01 as our regularizer since, we did not want to train for too many iterations before getting decent results as per Theorem 3 **CITE THIS**. We were able to get the following results for the BCFW algorithm:

For the dual extragradient, we obtained instead the following:

6 Discussion

7 Conclusion and Further Work

8 Citations and References

References

- [1] Simon Lacoste-Julien et al. “Block-Coordinate Frank-Wolfe Optimization for Structural SVMs”. In: (2013), p. 31.
- [2] Javier M. Moguerza and Alberto Muñoz. “Support Vector Machines with Applications”. In: *Statist. Sci.* 21.3 (Aug. 2006), pp. 322–336. ISSN: 0883-4237. DOI: 10.1214/088342306000000493. arXiv: math/0612817. URL: <http://arxiv.org/abs/math/0612817> (visited on 04/29/2019).
- [3] Ben Taskar, Simon Lacoste-Julien, and Michael I Jordan. “Structured Prediction via the Extragradient Method”. In: (), p. 12.
- [4] Ben Taskar, Simon Lacoste-Julien, and Michael I Jordan. “Structured Prediction, Dual Extragradient and Bregman Projections”. In: (2006), p. 27.

9 Appendix

9.0.1 Convergence analysis

The restricted gap function $\mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$ is upper bounded by:

$$\mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\overline{\mathbf{w}^T}, \overline{\mathbf{z}^T}) \leq \frac{(D_{\mathbf{w}} + D_{\mathbf{z}}) L}{\tau + 1} \quad (21)$$

In his proof on the convergence of the extragradient algorithm, Nesterov uses a function f_D instead of $\mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$, where f_D is defined as:

$$f_D(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Q}} \{ \langle g(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle : d(\mathbf{x}, \mathbf{y}) \} \quad (22)$$

where the set \mathcal{Q} is the set of parameters and g is a monotone operator. We can already see a link between the function f_D and the gap $\mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$. This comes out as:

$$\begin{aligned} \mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\mathbf{w}, \mathbf{z}) &= \sum_i \mathbf{w}^T \mathbf{F}_i \mathbf{z}_i^* - (\mathbf{w}^*)^T \mathbf{F}_i \mathbf{z}_i - \sum_i (\mathbf{w}^T - (\mathbf{w}^*)^T) \mathbf{f}_i(\mathbf{y}_i) - \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \\ &= \sum_i (\mathbf{z}_i^*)^T \mathbf{F}_i^T \mathbf{w} - (\mathbf{w}^*)^T \mathbf{F}_i \mathbf{z}_i - \sum_i (\mathbf{f}_i(\mathbf{y}_i))^T (\mathbf{w} - \mathbf{w}^*) - \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \\ &= \sum_i (\mathbf{z}_i^*)^T \mathbf{F}_i^T (\mathbf{w} - \mathbf{w}^*) - (\mathbf{w}^*)^T \mathbf{F}_i (\mathbf{z}_i - \mathbf{z}_i^*) - \sum_i (\mathbf{f}_i(\mathbf{y}_i))^T (\mathbf{w} - \mathbf{w}^*) - \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \\ &= \begin{pmatrix} \sum_i \mathbf{F}_i \mathbf{z}_i^* \\ -\mathbf{F}_1^T \mathbf{w}^* \\ \vdots \\ -\mathbf{F}_m^T \mathbf{w}^* \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} - \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} \end{aligned} \quad (23)$$

From this, we deduce that the function g from the definition of f_D corresponds to:

$$g(\mathbf{w}, \mathbf{z}) = \begin{pmatrix} \sum_i \mathbf{F}_i \mathbf{z}_i^* \\ -\mathbf{F}_i^T \mathbf{w}^* \\ \vdots \\ -\mathbf{F}_m^T \mathbf{w}^* \end{pmatrix} - \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix} = \mathbf{F} \mathbf{u}^* - \mathbf{a} \quad (24)$$

It is constant and thus monotone as required by Nesterov's proof of the convergence of the algorithm. Its Lipschitz constant L is equal to $\max_{\mathbf{u} \in \mathcal{U}} \|\mathbf{F}(\mathbf{u} - \mathbf{u}')\|_2 / \|\mathbf{u} - \mathbf{u}'\|_2 \leq \|\mathbf{F}\|_2$. Of course, a point \mathbf{w}, \mathbf{z} that satisfies $\|\mathbf{w}\|_2 \leq D_{\mathbf{w}}$ and $\|\mathbf{z}\|_2 \leq D_{\mathbf{z}}$ also satisfies $\|(\mathbf{w}, \mathbf{z})\|_2 \leq D$ when $D = \sqrt{D_{\mathbf{w}}^2 + D_{\mathbf{z}}^2}$ since $(\mathbf{w}, 0) \perp (0, \mathbf{z})$. It is then easy to see that $f_D \geq \mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$. Thus, the function $\mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$ is upper bounded by the right-hand side of equation 21. We can also observe that the function $g(\mathbf{w}, \mathbf{z})$ is exactly the gradient of the objective $\mathcal{L}(\mathbf{w}, \mathbf{z})$ at the point \mathbf{w}, \mathbf{z} .

9.1 From classical Frank-Wolfe to more sophisticated variants

9.1.1 Classical Frank-Wolfe

Consider a linear minimization oracle,

$$LMO_{\mathcal{A}}(\nabla f(x_t)) \in \operatorname{argmin}_{s' \in \mathcal{A}} \langle s', \nabla f(x_t) \rangle$$

Starting with an active set consisting of only an initial feasible point $S^0 = \{x_0\}$, the Frank-Wolfe algorithm adds an "atom" $s_t = LMO_{\mathcal{A}}(\nabla f(x_t))$, to the active set in a convex combination with its elements while maintaining this combination sparse.

9.1.2 Convergence results

We define the duality gap

$$g(\alpha^k) = \max_{s \in \mathcal{M}} \langle \alpha^k - s, \nabla f(\alpha^k) \rangle$$

By first order convexity of the objective, we have

$$\begin{aligned} f(s) &\geq f(\alpha^k) + \langle \alpha^k - s, \nabla f(\alpha^k) \rangle \\ \implies g(\alpha^k) &= -\min_{s \in \mathcal{M}} \langle \alpha^k - s, \nabla f(\alpha^k) \rangle \geq f(\alpha^k) - f^* \end{aligned}$$

We can thus see that the duality gap gives us a computable optimality guarantee.

Definition. The curvature constant C_f is given by the maximum relative deviation of the objective function f from its linear approximations, over the domain \mathcal{M} ,

$$C_f = \sup_{\substack{x, s \in \mathcal{M} \\ \gamma \in [0, 1], y = x + \gamma(s - x)}} \frac{2}{\gamma^2} \left(f(y) - f(x) - \langle y - x, \nabla f(x) \rangle \right)$$

Intuitively, the curvature constant can be seen as a measure of how flat the objective function is. For example, if the objective is linear, say $f(x) = ax + b$ and $x \in [e, f]$ then $\nabla f(x) = a$ and the curvature constant is zero:

$$C_f = \frac{2}{\gamma^2} \left(ay + b - ax - b + (-ay + ax) \right) = 0$$

Moreover $s = \operatorname{argmin}_{s \in [e, f]} \langle s, a \rangle = \frac{e}{a}$. Hence we reach the minimum in one F-W step.

Thus, we can observe that for flatter functions, that is with smaller curvature constants, Frank-Wolfe should converge faster.

Theorem. The duality gap obtained in the t^{th} iteration of the Frank-Wolfe algorithm satisfies

$$g(x_t) \leq 2\beta \frac{C_f}{t+2} (1 + \delta)$$

Where $\beta = \frac{27}{8}$ and δ is the approximation error tolerated in the *LMO*.

Definition. A function f has Lipschitz continuous gradient if:

$$\forall x, y \in \text{dom}(f), \exists L > 0 \quad \text{such that}$$

$$\|\nabla f(x) - \nabla f(y)\|^2 \leq L\|x - y\|^2$$

Theorem. If a convex function f on C has Lipschitz gradient, i.e $\|\nabla f(x) - \nabla f(y)\|_p \leq L_q\|x - y\|_p, \quad \forall x, y \in C$, then

$$C_f \leq L_q \cdot \text{diam}_p^2(C)$$

Proof. f has Lipschitz gradient therefore by the fundamental descent lemma we have,

$$f(y) - f(x) - \langle y - x, \nabla f(x) \rangle \leq \frac{L_q}{2} \|y - x\|_p^2$$

$$C_f \leq \max_{\substack{y=(1-\gamma)x+\gamma s \\ x, s, y \in C}} \frac{2}{\gamma^2} \frac{L_q}{2} \underbrace{\|y - x\|_p^2}_{=\gamma^2\|x-s\|_p^2}$$

$$C_f \leq L_q \underbrace{\max_{x, s \in C} \|x - s\|_p^2}_{\triangleq \text{diam}_p^2(C)} \quad \square.$$

Therefore, assuming $\delta = 0$, we get the following optimality certificate

$$f(\alpha^k) - f^* \leq g(\alpha^k) \leq 2\beta \frac{C_f}{t+2} \leq 2\beta \frac{L_q \cdot \text{diam}_p^2(C)}{t+2}$$

Thus, we see that the Frank-Wolfe algorithm has a sublinear convergence rate.

9.1.3 Optimality in terms of sparsity of the iterates

Lemma. For $f(x) = \|x\|_2^2$ and $1 \leq k \leq n$, it holds that

$$\min_{\substack{x \in \Delta_n \\ \text{card}(x) \leq k}} f(x) = \frac{1}{k}, \quad \text{and}$$

$$g(x) \geq \frac{2}{k} \quad \forall x \in \Delta_n \quad \text{s.t.} \quad \text{card}(x) \leq k.$$

By the first equality we have, for any vector x s.t. $\text{card}(x) = k$, we get $g(x) \leq \frac{1}{k} - \frac{1}{n}$. Thus, combining the upper and lower bound, we have that the sparsity (number of used atoms) by the Frank-Wolfe algorithm is worst case optimal.

Let $\alpha \in \mathcal{M}$

for $k = 0$ **to** K **do**

 Compute $s = \text{argmin}_{s \in \mathcal{M}} \langle s, \nabla f(\alpha^k) \rangle$

 Let $\gamma = \frac{2}{k+2}$, or optimize γ by line search

Update $\alpha^{k+1} = (1 - \gamma)\alpha^k + \gamma s$

end for

Theorem. Given a convex, differentiable objective $f : \mathcal{M}^1 \times \dots \times \mathcal{M}^n \rightarrow \mathbb{R}$, where $\forall i \in \{1..n\}$, each factor $\mathcal{M}^i \subseteq \mathbb{R}^n$ is convex and compact, if we are at a point x such that $f(x)$ is minimized along each coordinate axis, then x is a global minimum.

As in coordinate descent, we minimize the objective function one coordinate (block) at a time. At each iteration, BCFW picks the i^{th} block (from n) uniformly at random and updates the i^{th} coordinate of the corresponding weight, by calling the maximization oracle on the chosen block.

Let $w^0 = w_i^0 = \bar{w}^0 = 0$, $l^0 = l_i^0 = 0$

for $k = 0 \dots K$ **do**

Pick i at random in $\{1, \dots, n\}$

Solve $y_i^* = \max_{y_i \in \mathcal{Y}_i} H_i(y, w^k)$

Let $w_s = \frac{1}{n\lambda} \psi_i(y_i^*)$, and $l_s = \frac{1}{n} L_i(y_i^*)$

Let $\gamma = \frac{\lambda(w_i^k - w_s)^T w_i^k - l_i^k + l_s}{\lambda \|w_i^k - w_s\|^2}$, and clip to $[0, 1]$

Update $w_i^{k+1} = (1 - \gamma)w_i^k + \gamma w_s$, and $l_i^{k+1} = (1 - \gamma)l_i^k + \gamma l_s$

Update $w^{k+1} = w^k + w_i^{k+1} - w_i^k$, and $l^{k+1} = (1 - \gamma)l^k + \gamma l_s$

end for

9.1.4 Convergence Results

Definition. Over each coordinate block \mathcal{M}^i , let the curvature be given by,

$$C_f^{(i)} = \sup_{\substack{x \in \mathcal{M}, s_i \in \mathcal{M}^i \\ y = x + \gamma(s_{[i]} - x_{[i]}) \\ \gamma \in [0, 1]}} \frac{2}{\gamma^2} \left(f(y) - f(x) - \langle y_i - x_i, \nabla_i f(x) \rangle \right)$$

Where $x_{[i]}$ refers to the zero-padding of i^{th} coordinate of x . And let the global *product curvature constant* be,

$$C_f^\otimes = \sum_{i=1}^n C_f^{(i)}$$

Theorem. For the dual structural SVM objective function over the domain $\mathcal{M} = \Delta_{|\mathcal{Y}_1|} \times \dots \times \Delta_{|\mathcal{Y}_n|}$, the total curvature constant C_f^\otimes , on the product domain \mathcal{M} , is upper bounded by,

$$C_f^\otimes \leq \frac{4R^2}{\lambda n} \quad \text{where} \quad R = \max_{i \in [n], y \in \mathcal{Y}_i} \|\psi_i(y)\|_2$$

Proof. By the second order convexity condition on f at y , we have

$$\begin{aligned}
f(y) &\leq f(x) + \langle y_i - x_i, \nabla_i f(x) \rangle \\
&\quad + (y - x)^T \nabla^2 f(x) (y - x) \\
f(y) - f(x) - \langle y_i - x_i, \nabla_i f(x) \rangle &\leq (y - x)^T \nabla^2 f(x) (y - x) \\
C_f^{(i)} &\leq \sup_{\substack{x \in \mathcal{M}, s_i \in \mathcal{M}^i \\ y = x + \gamma(s_{[i]} - x_{[i]}) \\ \gamma \in [0, 1]}} \left(f(y) - f(x) - \langle y_i - x_i, \nabla_i f(x) \rangle \right) \\
&\leq \sup_{\substack{x, y \in \mathcal{M}, (y-x) \in \mathcal{M}^{[i]} \\ z \in [x, y] \subseteq \mathcal{M}}} (y - x)^T \nabla^2 f(z) (y - x) \\
\text{Moreover} \quad &\sup_{\substack{x, y \in \mathcal{M}, (y-x) \in \mathcal{M}^{[i]} \\ z \in [x, y] \subseteq \mathcal{M}}} (y - x)^T \nabla^2 f(z) (y - x) \\
&= \lambda \sup_{x, y \in \mathcal{M}, (y-x) \in \mathcal{M}^{[i]}} (A(y - x))^T \nabla^2 f(z) (A(y - x)) \\
C_f^{(i)} &\leq \lambda \sup_{v, w \in A\mathcal{M}^{(i)}} \|v - w\|_2^2 \leq \lambda \sup_{v \in A\mathcal{M}^{(i)}} \|2v\|_2^2
\end{aligned}$$

Where $\forall v \in A\mathcal{M}^{(i)}$, v is a convex combination of the feature vectors corresponding to the possible labelings for the i^{th} example of the training data, such that $\|v\|_2 \leq \text{the longest column of } A = \frac{1}{n\lambda}R$. Therefore,

$$C_f^\otimes = \sum_{i=1}^n C_f^{(i)} \leq 4\lambda \sum_{i=1}^n \left(\frac{1}{n\lambda}R \right)^2 = \frac{4}{n\lambda}R^2 \quad \square$$

First, we observe that the curvature constant for BCFW is n times smaller than that of batch Frank Wolfe which is $\leq \frac{4}{\lambda}R^2$. Hence the n times faster convergence rate of BCFW.

9.1.5 Tightening the bound

Definition. Let $\|\cdot\|$ be a norm on \mathbb{R}^n . The associated dual norm, denoted $\|\cdot\|_*$ is defined as,

$$\|z\|_* = \sup_z \{z^T x \mid \|x\| \leq 1\}$$

We denote the dual norm of l_p by l_q . For $p = 2$ we have $q = 2$ and for $p = 1$, $q = \infty$. *Problem.* For $p = q = 2$ we get $\text{diam}_2^2(C) = 2n$, and the Lipschitz constant L_q is the largest eigenvalue of the hessian.

$$\lambda A^T A = \frac{1}{n^2 \lambda} \left(\langle \psi_i(y) - \psi_j(y') \rangle \right)_{(i,y), (j,y')}$$

And say $\langle \psi_i(y) - \psi_j(y') \rangle \approx 1$ for a lot of outputs, we get:

$$\mathbb{1}^T \mathbb{1} \approx \underbrace{\mathbb{1} \text{diam}_2(C)}_{=\sqrt{2n}}$$

Hence the largest eigenvalue the hessian, and therefore the Lipschitz constant, can scale with the dimension of $A^T A$, i.e exponentially with the size of the training data, rendering the bound above very loose, and thus

of little practical use.

Solution. Taking $p = 1$ and therefore $q = \infty$, we get $Ldiam^2(C) \approx \frac{4}{\lambda} R^2$.

Combined with the convergence results above, we get a sublinear convergence rate for BCFW. And although subgradient methods converge at the same rate, BCFW presents an adaptive stepsize and an indication as to when to terminate, making it a more practical alternative.

9.2 Away steps Frank-Wolfe

When the minimizer of the objective function lies at the boundary of the domain, after a number of iterations, the duality gap starts to stagnate.

As a result of the strong dependency of the immediate iterate on previously accumulated atoms in the active set, as it approaches to boundary, the F-W algorithm starts to zig-zag around the descent direction as can be seen in figure 2.

```

Let  $x_0 \in \mathcal{A}$  and  $S_0 := \{x_0\}$ 
for  $t = 0 \dots T$  do
  Let  $s_t := LMO_{\mathcal{A}}(\nabla f(x_t))$  and  $d_t^{FW} := s_t - x_t$ 
  Let  $v_t \in \underset{v \in S_t}{argmax} \langle \nabla f(x_t), v \rangle$  and  $d_t^A := x_t - v_t$ 
  if  $g_t^{FW} = \langle -\nabla f(x_t), d_t^{FW} \rangle \leq \epsilon$  then return  $x_t$ 
  if  $\langle -\nabla f(x_t), d_t^{FW} \rangle \geq \langle -\nabla f(x_t), d_t^A \rangle$  then
     $d_t = d_t^{FW}$  and  $\gamma^{max} = 1$ 
  else
     $d_t = d_t^A$  and  $\gamma^{max} = \frac{\alpha_{v_t}^{(t)}}{1 - \alpha_{v_t}^{(t)}}$ 
  Line-search:  $\gamma_t \in argmin_{\gamma} f(x_t + \gamma d_t)$ 
  Update  $x_t = x_t + \gamma_t d_t$ 
end for

```

To address this issue an improved variant of F-W named **Away-steps Frank-Wolfe** adds the possibility of moving away (by removing a fraction of) a maximizer of the $LMO_{S_{\square}}$ in the active set. While this slows down each iteration it should be noted that the added step is easier than $LMO_{\mathcal{A}}$ given that we maximize over a subset of \mathcal{A} . Furthermore, given that this variant converges linearly, the algorithm progresses in a fewer number of iterations in the descent direction, making it much faster than the original F-W.

9.3 Randomized Away-step Frank-Wolfe

A crucial assumption in constructing the BCFW is whether the domain is block-separable. While this is true in the context of the structured SVM, this leaves out important cases such as l_1 constrained optimization (e.g. lasso type problems).

Moreover, while being an improvement on the classical variant by being $n = \text{size of the data}$ times cheaper per iteration, BCFW still converges at a sublinear rate unlike the Away-step FW.

The Randomized Away-steps Frank-Wolfe (RAWF) finds a compromise between the two variants. By subsampling a $\eta \in (0, 1]$ portion of the domain \mathcal{A} in the *LMO* and adding an away step at each iteration, we get a linear convergence rate with cheaper oracle calls than that of the original F-W.

9.3.1 Convergence results

Definition. Let the *away curvature* C_f^A and the *geometric strong convexity* constants be, respectively And $s_f, v_f(x)$ are the FW atom and away atom respectively, starting from x .

Theorem. Consider the set $\mathcal{M} = \text{conv}(\mathcal{A})$, with \mathcal{A} a finite set of extreme atoms, after T iterations of RAWF, we have the following convergence rate

$$E[f(x_{T+1})] - f^* \leq (f(x_0) - f^*) \cdot (1 - \eta^2 \rho_f)^{\max\{0, \lfloor \frac{T-s}{2} \rfloor\}}$$

With $\rho_f = \frac{\mu_f^A}{4C_f^A}$, $\eta = \frac{p}{|\mathcal{A}|}$ and $s = |S_0|$.

Proof sketch. First we upper-bound $h_t = f(x_t) - f^*$ by the pairwise dual gap $\tilde{g}_t = \langle \tilde{s}_t - v_t \rangle$, then we lower bound the progress $h_t - h_{t+1}$ by using the away curvature constant in similar way to the proof in (Lacoste-Julien & Jaggi, 2015, Theorem 8). \square

With the above theorem, we get

$$\lim_{t \rightarrow \infty} \frac{Ef(x_{t+1}) - f^*}{Ef(x_t) - f^*} \in (0, 1)$$

Thus proving a linear convergence rate for the Randomized Away-steps Frank-Wolfe.

Let $x_0 = \sum_{v \in \mathcal{A}} \alpha_v^{(0)}$ with $s = |S_0|$, a subsampling parameter $1 \leq p \leq |\mathcal{A}|$.

for $t = 0 \dots T$ **do**

 Get \mathcal{A}_t by sampling $\min\{p, |\mathcal{A} \setminus S_t|\}$ elements uniformly from $|\mathcal{A} \setminus S_t|$

 Compute $s_t = \text{LMO}(\nabla f(x), S_t \cup \mathcal{A}_t)$

 Let $d_t^{FW} = s_t - x_t$ **RFW step**

 Compute $v_t = \text{LMO}(-\nabla f(x), S_t)$

 Let $d_t^A = x_t - v_t$ **Away step**

```

if  $\langle -\nabla f(x_t), d_t^{FW} \rangle \geq \langle -\nabla f(x_t), d_t^A \rangle$  then
     $d_t = d_t^{FW}$  and  $\gamma^{max} = 1$ 
else
     $d_t = d_t^A$  and  $\gamma^{max} = \frac{\alpha_{v_t}^{(t)}}{1 - \alpha_{v_t}^{(t)}}$ 
    Let  $x_{t+1} = x_t + \gamma_t d_t$ 
    Let  $S_{t+1} = \{v \in \mathcal{A} \text{ s.t. } \alpha_{v_t}^{(t)} > 0\}$ 
end for

```