

Scalable optimization algorithms for the Structured SVM

Frederic Boileau Elyes Lamouchi William St-Arnaud

28th April 2019

1 Extra Gradient Algorithm

1.1 Preliminaries

In the typical setting for SVM's, we aim to optimize a linear objective over a set of constraints. There are many alternatives to solve that problem, the most popular ones being gradient-based methods. One of the issues with this type of algorithm is that they cannot be applied to matchings and min-cuts. In addition, variants that aim at generating constraints on the fly and incrementally solving new QP's do not typically scale well. The method proposed in this paper, which is in essence a modification of the extra-gradient method by Nesterov, leverages the min-max formulation and takes the dual ; the exponential number of constraints is then transformed an exponential number of variables. However, the paper highlights that we can use only a tractable number of these dual variables to solve the optimization problem. The method proposed in this paper also generalizes the applicability of the usual extra-gradient method to non-euclidean, Bergman projections. We thus end up with a framework that is more flexible and also efficient to implement. Since estimating the maximum likelihood over graphical models is often impractical or infeasible for a wide class of problems, we focus our attention on large margin estimation. For a dataset $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$, where each \mathbf{x}_i is an object with a structure (e.g. sequence of words in french), we attempt to find the optimal parameter \mathbf{w} of a linear classifier:

$$\mathbf{y}_i = \arg \max_{\mathbf{y}'_i \in \nabla} \mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{y}'_i) \quad (1)$$

The function f gives a feature mapping of a structured object with its corresponding label \mathbf{y}_i . The error of a prediction is measured by a loss function l . To make the loss convex, another term is introduced in the form a hinge loss. Since this gives an upper bound to the loss, it is natural to minimize it. We end up with a problem of the form:

$$\min_{\mathbf{w} \in \Delta} \sum_i \max_{\mathbf{y}'_i \in \nabla_i} [\mathbf{w}^T \mathbf{f}_i(\mathbf{y}'_i) + l_i(\mathbf{y}'_i)] - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i) \quad (2)$$

The parameters \mathbf{w} are also regularized with parameter λ . Since we are optimizing over \mathbf{y}'_i , we can drop the term from equation 17 and we end up with a loss-augmented inference problem inside the min function. The three types of structure that are presented in the paper have a general formulation that can better be expressed as:

$$\min_{\mathbf{w} \in \Delta} \max_{\mathbf{z} \in \mathcal{Z}} \sum_i \left(\mathbf{w}^T \mathbf{F}_i \mathbf{z}_i + \mathbf{c}_i^T \mathbf{z}_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i) \right) \quad (3)$$

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:

$$\begin{aligned} \min_{\mathbf{w} \in \Delta, (\lambda, \mu) \geq 0} \quad & \sum_i \left(\mathbf{b}_i^T \lambda_i + \mathbf{1}^T \mu_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i) \right) \\ \text{s.t.} \quad & \mathbf{F}_i^T \mathbf{w} + \mathbf{c}_i \leq \mathbf{A}_i^T \lambda_i + \mu_i \quad i = 1, \dots, m \end{aligned} \quad (4)$$

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 19, the term that is optmized is defined as:

$$\mathcal{J}(\mathbf{w}, \mathbf{z}) \triangleq \sum_i \mathbf{w}^T \mathbf{F}_i \mathbf{z}_i + \mathbf{c}_i^T \mathbf{z}_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i) \quad (5)$$

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.

$$\begin{pmatrix} \nabla_{\mathbf{w}} \mathcal{J}(\mathbf{w}, \mathbf{z}) \\ -\nabla_{\mathbf{z}_1} \mathcal{J}(\mathbf{w}, \mathbf{z}) \\ \vdots \\ -\nabla_{\mathbf{z}_m} \mathcal{J}(\mathbf{w}, \mathbf{z}) \end{pmatrix} = \begin{pmatrix} 0 & \mathbf{F}_1 & \dots & \mathbf{F}_m \\ -\mathbf{F}_1^T & & & \\ \vdots & & \mathbf{0} & \\ -\mathbf{F}_m^T & & & \end{pmatrix} \begin{pmatrix} \mathbf{w} \\ \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_m \end{pmatrix} - \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix} = \mathbf{F}\mathbf{u} - \mathbf{a} \quad (6)$$

$\underbrace{\hspace{10em}}_{\mathbf{F}} \quad \underbrace{\hspace{2em}}_{\mathbf{u}} \quad \underbrace{\hspace{2em}}_{\mathbf{a}}$

We can measure the “goodness” of the parameters using the gap function \mathbb{M} :

$$\mathbb{M}(\mathbf{w}, \mathbf{z}) \triangleq \left[\max_{\mathbf{z}' \in \mathcal{Z}} \mathcal{J}(\mathbf{w}, \mathbf{z}') - \mathcal{J}^* \right] + \left[\mathcal{J}^* - \min_{\mathbf{w}' \in \Delta} \mathcal{J}(\mathbf{w}', \mathbf{z}) \right] \quad (7)$$

where \mathcal{J}^* gives the result of the min-max of the objective \mathcal{J} . 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 $(\hat{\mathbf{u}}_{\mathbf{w}}, \hat{\mathbf{u}}_{\mathbf{z}}) \in \nabla$:

$$\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\mathbf{w}, \mathbf{z}) = \max_{\mathbf{z}' \in \nabla} [\mathcal{J}(\mathbf{w}', \mathbf{z}') : d(\mathbf{z}, \mathbf{z}') \leq D_{\mathbf{z}}] - \left[\min_{\mathbf{w}' \in \Delta} \mathcal{J}(\mathbf{w}', \mathbf{z}) : d(\mathbf{w}, \mathbf{w}') \leq D_{\mathbf{w}'} \right] \quad (8)$$

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.

1.2 Dual Extragradient algorithm

The dual extragradient algorithm from Nesterov gives a convergence guarantee for the objective \mathcal{J} . The algorithm can be given by:

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

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

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

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

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

end for

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

This algorithm has a lookahead step (i.e. v) that serves to perform the actual gradient update 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 u^t up to iteration n . The function g corresponds to the objective \mathcal{J} 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_{∇} in the algorithm simply projects a point back to the set ∇ by finding the nearest point with respect to the distance metric used.

1.2.1 Proximal step operator

We define the proximal step operator as follows:

$$\blacktriangleleft_{\eta}(\mathbf{u}, \mathbf{s}) = \max_{\mathbf{u}' \in \nabla} \left\{ \langle \mathbf{s}, \mathbf{u}' - \mathbf{u} \rangle - \frac{1}{\eta} d(\mathbf{u}, \mathbf{u}') \leq D \right\} \quad (9)$$

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

$$h(\mathbf{u}') \geq h(\mathbf{u}) + \langle \nabla h(\mathbf{u}), \mathbf{u}' - \mathbf{u} \rangle + \frac{\sigma}{2} \|\mathbf{u}' - \mathbf{u}\|^2 \quad (10)$$

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:

$$d(\mathbf{u}', \mathbf{u}) \triangleq h(\mathbf{u}') - h(\mathbf{u}) - \langle \nabla h(\mathbf{u}), \mathbf{u}' - \mathbf{u} \rangle \geq \frac{\sigma}{2} \|\mathbf{u}' - \mathbf{u}\|^2 \quad (11)$$

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 $h \star$ 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:

$$\mathbf{s} - \nabla_{\mathbf{u}'} d(\mathbf{u}', \mathbf{u}) = \mathbf{s} - \frac{1}{\eta} \nabla_{\mathbf{u}'} d(\mathbf{u}, \mathbf{u}') = \mathbf{s} - \frac{1}{\eta} [\nabla h(\mathbf{u}') - \nabla h(\mathbf{u})] \quad (12)$$

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

1.2.2 Convergence analysis

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

$$\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\overline{\mathbf{w}^\tau}, \overline{\mathbf{z}^\tau}) \leq \frac{(D_{\mathbf{w}} + D_{\mathbf{z}})L}{\tau + 1} \quad (13)$$

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

$$f_D(\mathbf{x}) = \max_{\mathbf{y} \in \nabla} \{\langle g(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle : d(\mathbf{x}, \mathbf{y})\} \quad (14)$$

where the set ∇ 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 $\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$. This comes out as:

$$\begin{aligned} \mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\mathbf{w}, \mathbf{z}) &= \sum_i \mathbf{w}^T \mathbf{F}_i \mathbf{z}_i^* - (\mathbf{w}^*)^T \mathbf{F}_i \mathbf{z}_i - \sum_i (\mathbf{w}^T - (\mathbf{w}^*)^T) \mathbf{f}_i(\mathbf{y}_i) - \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \\ &= \sum_i (\mathbf{z}_i^*)^T \mathbf{F}_i^T \mathbf{w} - (\mathbf{w}^*)^T \mathbf{F}_i \mathbf{z}_i - \sum_i (\mathbf{f}_i(\mathbf{y}_i))^T (\mathbf{w} - \mathbf{w}^*) - \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \\ &= \sum_i (\mathbf{z}_i^*)^T \mathbf{F}_i^T (\mathbf{w} - \mathbf{w}^*) - (\mathbf{w}^*)^T \mathbf{F}_i (\mathbf{z}_i - \mathbf{z}_i^*) - \sum_i (\mathbf{f}_i(\mathbf{y}_i))^T (\mathbf{w} - \mathbf{w}^*) - \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \\ &= \begin{pmatrix} \sum_i \mathbf{F}_i \mathbf{z}_i^* \\ -\mathbf{F}_1^T \mathbf{w}^* \\ \vdots \\ -\mathbf{F}_m^T \mathbf{w}^* \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} - \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} \end{aligned} \quad (15)$$

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

$$g(\mathbf{w}, \mathbf{z}) = \begin{pmatrix} \sum_i \mathbf{F}_i \mathbf{z}_i^* \\ -\mathbf{F}_i^T \mathbf{w}^* \\ \vdots \\ -\mathbf{F}_m^T \mathbf{w}^* \end{pmatrix} - \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix} = \mathbf{F} \mathbf{u}^* - \mathbf{a} \quad (16)$$

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 $\max_{\mathbf{u} \in \mathbf{V}} \|\mathbf{F}(\mathbf{u} - \mathbf{u}')\|_2 / \|\mathbf{u} - \mathbf{u}'\|_2 \leq \|\mathbf{F}\|_2$. 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 $f_D \geq \mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$. Thus, the function $\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$ is upper bounded by the right-hand side of equation 30. We can also observe that the function $g(\mathbf{w}, \mathbf{z})$ is exactly the gradient of the objective $\mathfrak{J}(\mathbf{w}, \mathbf{z})$ at the point \mathbf{w}, \mathbf{z} .

1.2.3 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.

1.2.4 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 $|\mathbf{B}| \gg |\mathbf{A}|$ 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 $|\mathbf{A}|$ (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.

2 Frank Wolfe and Variants

2.1 Extragradient Algorithm for structured predictions

2.1.1 Preliminaries

In the typical setting for SVM's, we aim to optimize a linear objective over a set of constraints. There are many alternatives to solve that problem, the most popular ones being gradient-based methods. One of the

issues with this type of algorithm is that they cannot be applied to matchings and min-cuts. In addition, variants that aim at generating constraints on the fly and incrementally solving new QP's do not typically scale well. The method proposed in this paper, which is in essence a modification of the extra-gradient method by Nesterov, leverages the min-max formulation and takes the dual ; the exponential number of constraints is then transformed an exponential number of variables. However, the paper highlights that we can use only a tractable number of these dual variables to solve the optimization problem. The method proposed in this paper also generalizes the applicability of the usual extra-gradient method to non-euclidean, Bergman projections. We thus end up with a framework that is more flexible and also efficient to implement. Since estimating the maximum likelihood over graphical models is often impractical or infeasible for a wide class of problems, we focus our attention on large margin estimation. For a dataset $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^m$, where each \mathbf{x}_i is an object with a structure (e.g. sequence of words in french), we attempt to find the optimal parameter \mathbf{w} of a linear classifier:

$$\mathbf{y}_i = \arg \max_{\mathbf{y}'_i \in \nabla} \mathbf{w}^T \mathbf{f}(\mathbf{x}_i, \mathbf{y}'_i) \quad (17)$$

The function f gives a feature mapping of a structured object with its corresponding label \mathbf{y}_i . The error of a prediction is mesured by a loss function l . To make the loss convex, another term is introduced in the form a hinge loss. Since this gives an upper bound to the loss, it is natural to minimize it. We end up with a problem of the form:

$$\min_{\mathbf{w} \in \Delta} \sum_i \max_{\mathbf{y}'_i \in \nabla_i} [\mathbf{w}^T \mathbf{f}_i(\mathbf{y}'_i) + l_i(\mathbf{y}'_i)] - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i) \quad (18)$$

The parameters \mathbf{w} are also regularized with parameter λ . Since we are optimizing over \mathbf{y}'_i , we can drop the term from equation 17 and we end up with a loss-augmented inference problem inside the min function. The three types of structure that are presented in the paper have a general formulation that can better be expressed as:

$$\min_{\mathbf{w} \in \Delta} \max_{\mathbf{z} \in \mathbf{Z}} \sum_i (\mathbf{w}^T \mathbf{F}_i \mathbf{z}_i + \mathbf{c}_i^T \mathbf{z}_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i)) \quad (19)$$

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:

$$\min_{\mathbf{w} \in \Delta, (\lambda, \mu) \geq 0} \sum_i (\mathbf{b}_i^T \lambda_i + \mathbf{1}^T \mu_i - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i)) \quad (20)$$

$$\text{s.t.} \quad \mathbf{F}_i^T \mathbf{w} + \mathbf{c}_i \leq \mathbf{A}_i^T \lambda_i + \mu_i \quad i = 1, \dots, m \quad (21)$$

The number of variables and constraints is linear in the number of parameters 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 19, the term that is optimized is defined as:

$$\mathcal{J}(\mathbf{w}, \mathbf{z}) \triangleq \sum_i \mathbf{w}^T \mathbf{F}_i \mathbf{z}_i + \mathbf{c}_i^T - \mathbf{w}^T \mathbf{f}_i(\mathbf{y}_i) \quad (22)$$

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.

$$\begin{aligned} \text{Let } \mathbf{F} &= \begin{pmatrix} 0 & \mathbf{F}_1 & \dots & \mathbf{F}_m \\ -\mathbf{F}_1^T & & & \\ \vdots & & \mathbf{0} & \\ -\mathbf{F}_m^T & & & \end{pmatrix} \\ \mathbf{u} &= \begin{pmatrix} \mathbf{w} \\ \mathbf{z}_1 \\ \vdots \\ \mathbf{z}_m \end{pmatrix} \quad \text{and} \quad \mathbf{a} = \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix} \\ \text{such that } & \begin{pmatrix} \nabla_{\mathbf{w}} \mathcal{J}(\mathbf{w}, \mathbf{z}) \\ -\nabla_{\mathbf{z}_1} \mathcal{J}(\mathbf{w}, \mathbf{z}) \\ \vdots \\ -\nabla_{\mathbf{z}_m} \mathcal{J}(\mathbf{w}, \mathbf{z}) \end{pmatrix} = \mathbf{F} \mathbf{u} - \mathbf{a} \end{aligned}$$

We can measure the “goodness” of the parameters using the gap function \mathbb{M} :

$$\mathbb{M}(\mathbf{w}, \mathbf{z}) \triangleq \left[\max_{\mathbf{z}' \in \mathbf{\Delta}} \mathcal{J}(\mathbf{w}, \mathbf{z}') - \mathcal{J}^* \right] + \left[\mathcal{J}^* - \min_{\mathbf{w}' \in \Delta} \mathcal{J}(\mathbf{w}', \mathbf{z}) \right] \quad (23)$$

where \mathcal{J}^* gives the result of the min-max of the objective \mathcal{J} . 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 $(\hat{\mathbf{u}}_{\mathbf{w}}, \hat{\mathbf{u}}_{\mathbf{z}}) \in \mathbf{\nabla}$:

$$\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\mathbf{w}, \mathbf{z}) = \max_{\mathbf{z}' \in \mathbf{\Delta}} [\mathcal{J}(\mathbf{w}', \mathbf{z}') : d(\mathbf{z}, \mathbf{z}') \leq D_{\mathbf{z}}] \quad (24)$$

$$- \left[\min_{\mathbf{w}' \in \Delta} \mathcal{J}(\mathbf{w}', \mathbf{z}) : d(\mathbf{w}, \mathbf{w}') \leq D_{\mathbf{w}'} \right] \quad (25)$$

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.

2.1.2 Dual Extragradient algorithm

The dual extragradient algorithm from Nesterov gives a convergence guarantee for the objective \mathfrak{J} . The algorithm can be given by:

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

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

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

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

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

end for

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

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 $f_D(\bar{\mathbf{u}}^n) = \max_y \{\langle g(y), \bar{\mathbf{u}}^n - y \rangle : d(\hat{\mathbf{u}}, y) \leq D\}$, where $\bar{\mathbf{u}}^n$ is the weighted average over all the updates \mathbf{u}^t up to iteration n . The function g corresponds to the objective \mathfrak{J} in our setting. When value of $f_D(\bar{\mathbf{u}}^n)$ 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_{\blacktriangledown} in the algorithm simply projects a point back to the set \blacktriangledown by finding the nearest point with respect to the distance metric used.

2.1.3 Proximal step operator

We define the proximal step operator as follows:

$$\blacktriangleleft_{\eta}(\mathbf{u}, \mathbf{s}) = \max_{\mathbf{u}' \in \blacktriangledown} \left\{ \langle \mathbf{s}, \mathbf{u}' - \mathbf{u} \rangle - \frac{1}{\eta} d(\mathbf{u}, \mathbf{u}') \leq D \right\} \quad (26)$$

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

$$h(\mathbf{u}') \geq h(\mathbf{u}) + \langle \nabla h(\mathbf{u}), \mathbf{u}' - \mathbf{u} \rangle + \frac{\sigma}{2} \|\mathbf{u}' - \mathbf{u}\|^2 \quad (27)$$

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:

$$d(\mathbf{u}', \mathbf{u}) \triangleq h(\mathbf{u}') - h(\mathbf{u}) - \langle \nabla h(\mathbf{u}), \mathbf{u}' - \mathbf{u} \rangle \geq \frac{\sigma}{2} \|\mathbf{u}' - \mathbf{u}\|^2 \quad (28)$$

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 $h \star$ 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:

$$\begin{aligned} \mathbf{s} - \nabla_{\mathbf{u}'} d(\mathbf{u}', \mathbf{u}) &= \mathbf{s} - \frac{1}{\eta} \nabla_{\mathbf{u}'} d(\mathbf{u}, \mathbf{u}') \\ &= \mathbf{s} - \frac{1}{\eta} [\nabla h(\mathbf{u}') - \nabla h(\mathbf{u})] \end{aligned} \tag{29}$$

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

2.1.4 Convergence analysis

The restricted gap function $\mathbb{M}_{D_{\mathbf{v}}, D_{\mathbf{z}}}$ is upper bounded by:

$$\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\overline{\mathbf{w}^\tau}, \overline{\mathbf{z}^\tau}) \leq \frac{(D_{\mathbf{w}} + D_{\mathbf{z}})L}{\tau + 1} \tag{30}$$

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

$$f_D(\mathbf{x}) = \max_{\mathbf{y} \in \nabla} \{\langle g(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle : d(\mathbf{x}, \mathbf{y})\} \tag{31}$$

where the set ∇ 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 $\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$. This comes out as:

$$\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\mathbf{w}, \mathbf{z}) = \sum_i \mathbf{w}^T \mathbf{F}_i \mathbf{z}_i^* - (\mathbf{w}^*)^T \mathbf{F}_i \mathbf{z}_i - \sum_i (\mathbf{w}^T - (\mathbf{w}^*)^T) \mathbf{f}_i(\mathbf{y}_i) \quad (32)$$

$$- \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \quad (33)$$

$$= \sum_i (\mathbf{z}_i^*)^T \mathbf{F}_i^T \mathbf{w} - (\mathbf{w}^*)^T \mathbf{F}_i \mathbf{z}_i - \sum_i (\mathbf{f}_i(\mathbf{y}_i))^T (\mathbf{w} - \mathbf{w}^*) \quad (34)$$

$$- \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \quad (35)$$

$$= \sum_i (\mathbf{z}_i^*)^T \mathbf{F}_i^T (\mathbf{w} - \mathbf{w}^*) - (\mathbf{w}^*)^T \mathbf{F}_i (\mathbf{z}_i - \mathbf{z}_i^*) \quad (36)$$

$$- \sum_i (\mathbf{f}_i(\mathbf{y}_i))^T (\mathbf{w} - \mathbf{w}^*) - \sum_i \mathbf{c}_i^T (\mathbf{z}_i - \mathbf{z}_i^*) \quad (37)$$

$$= \begin{pmatrix} \sum_i \mathbf{F}_i \mathbf{z}_i^* \\ -\mathbf{F}_1^T \mathbf{w}^* \\ \vdots \\ -\mathbf{F}_m^T \mathbf{w}^* \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} - \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} \quad (38)$$

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

$$g(\mathbf{w}, \mathbf{z}) = \begin{pmatrix} \sum_i \mathbf{F}_i \mathbf{z}_i^* \\ -\mathbf{F}_i^T \mathbf{w}^* \\ \vdots \\ -\mathbf{F}_m^T \mathbf{w}^* \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} - \begin{pmatrix} \sum_i \mathbf{f}_i(\mathbf{y}_i) \\ \mathbf{c}_1 \\ \vdots \\ \mathbf{c}_m \end{pmatrix}^T \begin{pmatrix} \mathbf{w} - \mathbf{w}^* \\ \mathbf{z}_1 - \mathbf{z}_1^* \\ \vdots \\ \mathbf{z}_m - \mathbf{z}_m^* \end{pmatrix} = \mathbf{F} \mathbf{u}^* - \mathbf{a} \quad (39)$$

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 $\max_{\mathbf{u} \in \mathbf{v}} \|\mathbf{F}(\mathbf{u} - \mathbf{u}')\|_2 / \|\mathbf{u} - \mathbf{u}'\|_2 \leq \|\mathbf{F}\|_2$. 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 $f_D \geq \mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$. Thus, the function $\mathbb{M}_{D_{\mathbf{w}}, D_{\mathbf{z}}}$ is upper bounded by the right-hand side of equation 30. We can also observe that the function $g(\mathbf{w}, \mathbf{z})$ is exactly the gradient of the objective $\mathcal{J}(\mathbf{w}, \mathbf{z})$ at the point \mathbf{w}, \mathbf{z} .

2.1.5 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.

2.1.6 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 $|\mathbf{P}| \gg |\Delta|$ 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 $|\Delta|$ (although we now need to store s_w^t). It can be reconstructed using \bar{u}_w^t . Figure 1 illustrates the various dependencies.

2.2 From classical Frank-Wolfe to more sophisticated variants

2.2.1 Classical Frank-Wolfe

Consider the problem of minimizing a convex objective function f over the convex hull of its domain $\mathcal{F} = \text{conv}(\mathcal{C})$.

By minimizing the first-order approximation of the objective function, the Frank-Wolfe algorithm takes a convex combination of the immediate iterate with the previous one. Consider a linear minimization oracle,

$$LMO_{\mathcal{C}}(\nabla f(x_t)) \in \text{argmin}_{s \in \mathcal{C}} \langle s, \nabla f(x_t) \rangle$$

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

2.2.2 Convergence results

We define the duality gap

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

By first order convexity of the objective, we have

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

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{F} ,

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

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} iteraton 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:

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

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

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

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

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

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

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

$$C_f \leq L_q \max_{x, s \in C} \|x - s\|_p^2 \quad \square.$$

$$\stackrel{\triangle}{=} \operatorname{diam}_p^2(C)$$

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

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

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

2.2.3 Optimality in terms of sparsity of the iterates

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

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

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

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  $\alpha \in \mathbb{F}$ 
for  $k = 0$  to  $K$  do
  Compute  $s = \text{argmin}_{s \in \mathbb{F}} \langle s, \nabla f(\alpha^k) \rangle$ 
  Let  $\gamma = \frac{2}{k+2}$ , or optimize  $\gamma$  by line search
  Update  $\alpha^{k+1} = (1 - \gamma)\alpha^k + \gamma s$ 
end for

```

2.2.4 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  $x_0 \in \mathcal{O}$  and  $S_0 := \{x_0\}$ 
for  $t = 0 \dots T$  do
  Let  $s_t := LMO_{\mathcal{O}}(\nabla f(x_t))$  and  $d_t^{FW} := s_t - x_t$ 
  Let  $v_t \in \text{argmax}_{v \in S_t} \langle \nabla f(x_t), v \rangle$  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:  $\gamma_t \in \text{argmin}_{\gamma} f(x_t + \gamma d_t)$ 
  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 $LMO_{\mathbf{A}}$ in the active set. While this slows down each iteration it should be noted that the added step is easier than $LMO_{\mathcal{O}}$ given that we maximize over a subset of \mathcal{O} . 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.

2.2.5 Block Coordinate Frank Wolfe

2.2.6 Structured SVM context

Given a training set $\mathcal{U} = \{(x_i, y_i)\}_{i=1}^n$ where $y \in \nabla$ is a multi-label output, and a feature map $\phi : \mathcal{X} \times \nabla \rightarrow \mathbb{R}$, which encodes a similarity measure between \mathcal{X} and ∇ , such that if y_i is the ground truth (target) for an input x_i , then

$$\forall y \in \nabla \setminus \{y_i\} \quad \text{we have} \quad \psi_i(y) = \phi(x_i, y_i) - \phi(x_i, y) > 0$$

The aim is to construct an accurate linear classifier, $h_w(x) = \underset{y \in \nabla(x)}{\operatorname{argmax}} \langle w, \phi(x, y) \rangle$.

To learn w , consider the task loss $L : \nabla \times \nabla \rightarrow \mathbb{R}_+$, where $L(y, y') = 0 \iff y = y'$.

The n -slack formulation of the problem would be,

$$\begin{aligned} \max_{w, \xi} \quad & \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \varepsilon_i \\ \text{s.t.} \quad & \langle w, \psi_i(y) \rangle \geq L(y_i, y) - \varepsilon_i, \quad \forall i, \forall y \in \nabla(x) = \nabla_i \end{aligned}$$

Problems: (1) The zero-one loss is not differentiable and (2) we have an exponential number of constraints.

Solutions: (1) Minimizing an upper bound to the task loss gives us a worst case guarantee.

Consider the **max oracle**, $\tilde{H} = \max_{y \in \nabla_i} L_i(y) - \langle w, \psi_i(y) \rangle$.
 $\underbrace{\hspace{1.5cm}}_{=H_i(y,w)} \quad \underbrace{\hspace{1.5cm}}_{\text{the hinge loss}}$

(2) The exponential number of constraints are replaced by n piecewise linear ones.

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

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

$$\begin{aligned} L(y_i, h_w(x_i)) &\leq L(y_i, h_w(x_i)) + \underbrace{\langle w, \psi_i(y) \rangle}_{\geq 0 \text{ by definition}} \\ &\leq \max_{y \in \nabla_i} L_i(y) - \langle w, \psi_i(y) \rangle \end{aligned}$$

Thus learning w amounts to the unconstrained problem,

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

2.2.7 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,

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

Moreover, with

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

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

The gradient of the dual would be,

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

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 $\alpha \in \Gamma$

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

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

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

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

 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)**.

Theorem. Given a convex, differentiable objective $f : \mathbb{R}^1 \times \dots \times \mathbb{R}^n \rightarrow \mathbb{R}$, where $\forall i \in \{1..n\}$, each factor $\mathbb{R}^i \subseteq \mathbb{R}^n$ 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 $y_i^* = \max_{y_i \in \nabla_i} H_i(y, w^k)$

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

 Let $\gamma = \frac{\lambda(w_i^k - w_s)^T w^k - l_i^k + l_s}{\lambda \|w_i^k - w_s\|^2}$, 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 $l_i^{k+1} = (1 - \gamma)l_i^k + \gamma l_s$

end for

2.2.8 Convergence Results

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

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

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 $\mathbb{R} = \Delta_{|\nabla_1|} \times \dots \times \Delta_{|\nabla_n|}$, the total curvature constant C_f^\otimes , on the product domain \mathbb{R} , is upper bounded by,

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

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

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

Where $\forall v \in A\Gamma^{(i)}$, 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.

2.2.9 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:

$$\mathbb{1}^T \mathbb{1} \approx \mathbb{1} \text{diam}_2(C)$$



$= \sqrt{2n}$

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.

2.3 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 \mathfrak{O} 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.

2.3.1 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 $\mathfrak{F} = \text{conv}(\mathfrak{O})$, with \mathfrak{O} 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 $\rho_f = \frac{\mu_f^A}{4C_f^A}$, $\eta \frac{p}{|\mathfrak{O}|}$ 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

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

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

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

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

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

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

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

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

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 $S_{t+1} = \{v \in \mathcal{O} \mid s.t. \alpha_{v_t}^{(t)} > 0\}$

end for

3 Empirical

3.1 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 **CITE THIS SHIT**. 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. The data had been aigned using the GIZA++ algorithm **CITE THIS SHIT**. 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.

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. As in Taskar **CITE THIS SHIT**, the features were composed of the following statistics of the corpora:

- Dice coefficient between pairs of words:

$$\frac{2 * C_{ef}}{C_e + C_f}$$

where C_{ef} is the number of co-occurrences of the pair of words e, f in sentences and C_e, C_f are the number of occurrences of the words in english and french respectively.

- Absolute value of the difference of the relative positions of each words in each sentence

- The square root of that value
- The square of the same value
- The hamming distance between the two words in a pair (padded with zeros when the length is not the same)

We combined these statistics by summing them for all word tuples in each sentence pairs. As an example, consider the following two sentences:

- This assignment was hard
- Ce travail etait ardu

We would compute the statistics for each word tuple (e.g. this/ce, this/travail, ..., hard/ardu). Thus, we end up with a vector of statistics for each tuple. 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 SHIT** 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.

3.1.1 Results