

Extragradient Algorithm for structured predictions

William St-Arnaud

April 16, 2019

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:

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

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:

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

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 1 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:

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

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:

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

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 3, the term that is optimized is defined as:

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

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.

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

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

$$\mathcal{G}(\mathbf{w}, \mathbf{z}) \triangleq \left[\max_{\mathbf{z}' \in \mathcal{Z}} \mathcal{L}(\mathbf{w}, \mathbf{z}') - \mathcal{L}^* \right] + \left[\mathcal{L}^* - \min_{\mathbf{w}' \in \mathcal{W}} \mathcal{L}(\mathbf{w}', \mathbf{z}) \right] \quad (7)$$

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

$$\mathcal{G}_{D_{\mathbf{w}}, D_{\mathbf{z}}}(\mathbf{w}, \mathbf{z}) = \max_{\mathbf{z}' \in \mathcal{Z}} [\mathcal{L}(\mathbf{w}', \mathbf{z}') : d(\mathbf{z}, \mathbf{z}') \leq D_{\mathbf{z}}] - \left[\min_{\mathbf{w}' \in \mathcal{W}} \mathcal{L}(\mathbf{w}', \mathbf{z}) : d(\mathbf{w}, \mathbf{w}') \leq D_{\mathbf{w}} \right] \quad (8)$$

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.

2 Dual Extragradient algorithm

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

```

Initialize: Choose  $\hat{\mathbf{u}} \in \mathcal{U}$ , set  $\mathbf{s}^{-1} = 0$ .
for  $t = 0$  to  $t = \tau$  do
   $\mathbf{v} = \Pi_{\mathcal{U}}(\hat{\mathbf{u}} + \eta \mathbf{s}^{t-1})$ 
   $\mathbf{u}^t = \Pi_{\mathcal{U}}(\mathbf{v} - \eta(\mathbf{F}\mathbf{v} - \mathbf{a}))$ 
   $\mathbf{s}^t = \mathbf{s}^{t-1} - (\mathbf{F}\mathbf{u}^t - \mathbf{a})$ 
end for
return  $\bar{\mathbf{u}}^{\tau} = \frac{1}{1+\tau} \sum_{t=0}^{\tau} \mathbf{u}^t$ 

```

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{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 $\Pi_{\mathcal{U}}$ in the algorithm simply projects a point back to the set \mathcal{U} by finding the nearest point with respect to the distance metric used.

2.1 Proximal step operator

We define the proximal step operator as follows:

$$\mathcal{T}_\eta(\mathbf{u}, \mathbf{s}) = \max_{\mathbf{u} \in \mathcal{U}} \left\{ \langle \mathbf{s}, \mathbf{u}' - \mathbf{u} \rangle - \frac{1}{\eta} d(\mathbf{u}, \mathbf{u}') \leq D \right\} \quad (9)$$

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

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

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^* 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, $h(\mathbf{u}) = \frac{1}{2} \|\mathbf{u}'\|^2$.

2.2 Convergence analysis

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

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

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

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

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:

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

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

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

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 $\max_{\mathbf{u} \in \mathcal{U}} \|\mathbf{F}(\mathbf{u} - \mathbf{u}')\|_2 / \|\mathbf{u} - \mathbf{u}'\|_2 \leq \|\mathbf{F}\|_2$. 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 13. 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} .

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.

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 $|\mathcal{Z}| \gg |\mathcal{W}|$ since $\bar{u}^t = \{u_w^t, u_z^t\}$ and we only care about the part that corresponds to w . \bar{u}_z^t is maintained implicitly by storing a vector of size $|\mathcal{W}|$ (although we now need to store s_w^t). It can be reconstructed using u_w^t . The following figure (**ADD FIGURE 5 FROM PAPER**) illustrates the various dependencies.