

CS794/CO673: Optimization for Data Science

Lec 15: Projection Algorithms

Yaoliang Yu



UNIVERSITY OF
WATERLOO

FACULTY OF MATHEMATICS
DAVID R. CHERITON SCHOOL
OF COMPUTER SCIENCE

November 04, 2022

Problem

Constrained minimization problem:

$$\begin{aligned} \inf_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \bigcap_{i \in I} C_i, \end{aligned}$$

- Each $C_i \subseteq \mathbb{R}^d$ is closed, convex and simple
- Projector $P_i = P_{C_i}$ can be easily computed
- However, projecting to the intersection C is usually much harder
- Function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is convex

Problem

Constrained minimization problem:

$$\begin{aligned} \inf_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \bigcap_{i \in I} C_i, \end{aligned}$$

- Each $C_i \subseteq \mathbb{R}^d$ is closed, convex and simple
- Projector $P_i = P_{C_i}$ can be easily computed
- However, projecting to the intersection C is usually much harder
- Function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is convex

Problem

Constrained minimization problem:

$$\begin{aligned} \inf_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \bigcap_{i \in I} C_i, \end{aligned}$$

- Each $C_i \subseteq \mathbb{R}^d$ is closed, convex and simple
- Projector $P_i = P_{C_i}$ can be easily computed
- However, projecting to the intersection C is usually much harder
- Function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is convex

Problem

Constrained minimization problem:

$$\begin{aligned} \inf_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \bigcap_{i \in I} C_i, \end{aligned}$$

- Each $C_i \subseteq \mathbb{R}^d$ is closed, convex and simple
- Projector $P_i = P_{C_i}$ can be easily computed
- However, projecting to the intersection C is usually much harder
- Function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is convex

Problem

Constrained minimization problem:

$$\begin{aligned} \inf_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) \\ \text{s.t. } \mathbf{w} \in \bigcap_{i \in I} C_i, \end{aligned}$$

- Each $C_i \subseteq \mathbb{R}^d$ is closed, convex and simple
- Projector $P_i = P_{C_i}$ can be easily computed
- However, projecting to the intersection C is usually much harder
- Function $f : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is convex

Perceptron and SVM revisited

Recall the perceptron problem:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) &\equiv 0 \\ \text{s.t. } \mathbf{w} &\in \bigcap_{i=1}^n C_i, \quad \text{where } C_i := \{\mathbf{w} : \langle y_i \mathbf{x}_i, \mathbf{w} \rangle \geq 1\} \end{aligned}$$

Similarly, we may rewrite the hard-margin SVM problem as:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|_2^2 \quad \text{s.t. } \mathbf{w} \in \bigcap_{i=1}^n C_i.$$

We note that the projector P_{C_i} is available in closed-form:

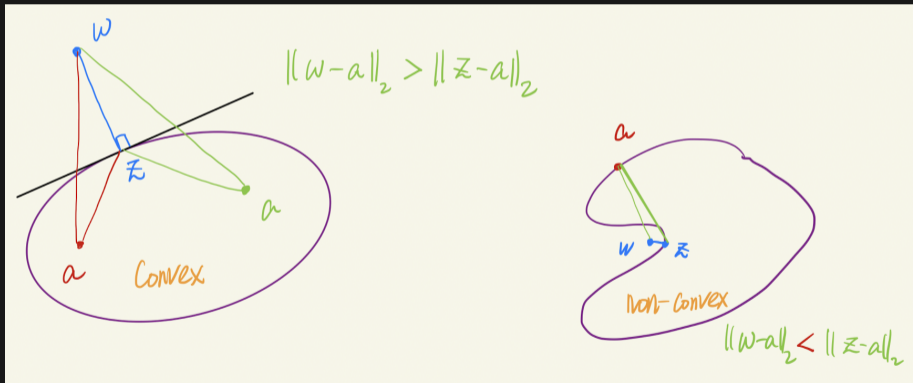
$$P_{C_i}(\mathbf{z}) := \left[\operatorname{argmin}_{\mathbf{w} \in C_i} \|\mathbf{w} - \mathbf{z}\|_2 \right] = \mathbf{z} + \frac{(1 - \langle y_i \mathbf{x}_i, \mathbf{z} \rangle)_+}{\|\mathbf{x}_i\|_2^2} y_i \mathbf{x}_i.$$

A nonconvex example

	2			3		9		7
	1							
4		7				2		8
		5	2				9	
			1	8		7		
	4				3			
				6			7	1
	7							
9		3		2		6		5

Theorem: Fejér's characterization of the closed convex hull

Let $A \subseteq \mathbb{R}^d$. Then, $w \notin \overline{\text{conv}}A$ iff there exists $z \in \mathbb{R}^d$ such that for all $a \in A$ (hence all $a \in \overline{\text{conv}}A$) we have $\|w - a\|_2 > \|z - a\|_2$.



L. Fejér. "Über die Lage der Nullstellen von Polynomen, die aus Minimumforderungen gewisser Art entspringen". *Mathematische Annalen*, vol. 85, no. 1 (1922), pp. 41–48.

Algorithmic Significance of Fejér's Result

Can be used to solve the convex feasibility problem:

$$\text{find } \mathbf{w} \in C,$$

where the closed (and convex) set $C \subseteq \mathbb{R}^d$ represents the solutions set of any problem. Indeed, starting from an arbitrary point \mathbf{w}_0 , if it is in C then we are done; if not then according to Fejér's Theorem there exists some \mathbf{w}_1 such that $\|\mathbf{w}_1 - \mathbf{w}\| < \|\mathbf{w}_0 - \mathbf{w}\|$ for all $\mathbf{w} \in C$.

- We need to be able to certify if $\mathbf{w}_0 \in C$, which may be trivial when the set C is defined by *explicit* inequalities, such as $C = \{\mathbf{w} : g(\mathbf{w}) \leq 0\}$.
- If $\mathbf{w}_0 \notin C$, we need to be able to *explicitly and efficiently* find \mathbf{w}_1 .
- We also need sufficient decrease so that $\text{dist}(\mathbf{w}_t, C) \rightarrow 0$.
- We may also want to prove the convergence (rate) of the whole sequence \mathbf{w}_t .

Algorithmic Significance of Fejér's Result

Can be used to solve the convex feasibility problem:

$$\text{find } \mathbf{w} \in C,$$

where the closed (and convex) set $C \subseteq \mathbb{R}^d$ represents the solutions set of any problem. Indeed, starting from an arbitrary point \mathbf{w}_0 , if it is in C then we are done; if not then according to Fejér's Theorem there exists some \mathbf{w}_1 such that $\|\mathbf{w}_1 - \mathbf{w}\| < \|\mathbf{w}_0 - \mathbf{w}\|$ for all $\mathbf{w} \in C$.

- We need to be able to certify if $\mathbf{w}_0 \in C$, which may be trivial when the set C is defined by *explicit* inequalities, such as $C = \{\mathbf{w} : g(\mathbf{w}) \leq 0\}$.
- If $\mathbf{w}_0 \notin C$, we need to be able to *explicitly and efficiently* find \mathbf{w}_1 .
- We also need sufficient decrease so that $\text{dist}(\mathbf{w}_t, C) \rightarrow 0$.
- We may also want to prove the convergence (rate) of the whole sequence \mathbf{w}_t .

Algorithmic Significance of Fejér's Result

Can be used to solve the convex feasibility problem:

$$\text{find } \mathbf{w} \in C,$$

where the closed (and convex) set $C \subseteq \mathbb{R}^d$ represents the solutions set of any problem. Indeed, starting from an arbitrary point \mathbf{w}_0 , if it is in C then we are done; if not then according to Fejér's Theorem there exists some \mathbf{w}_1 such that $\|\mathbf{w}_1 - \mathbf{w}\| < \|\mathbf{w}_0 - \mathbf{w}\|$ for all $\mathbf{w} \in C$.

- We need to be able to certify if $\mathbf{w}_0 \in C$, which may be trivial when the set C is defined by *explicit* inequalities, such as $C = \{\mathbf{w} : g(\mathbf{w}) \leq 0\}$.
- If $\mathbf{w}_0 \notin C$, we need to be able to *explicitly and efficiently* find \mathbf{w}_1 .
- We also need sufficient decrease so that $\text{dist}(\mathbf{w}_t, C) \rightarrow 0$.
- We may also want to prove the convergence (rate) of the whole sequence \mathbf{w}_t .

Algorithmic Significance of Fejér's Result

Can be used to solve the convex feasibility problem:

$$\text{find } \mathbf{w} \in C,$$

where the closed (and convex) set $C \subseteq \mathbb{R}^d$ represents the solutions set of any problem. Indeed, starting from an arbitrary point \mathbf{w}_0 , if it is in C then we are done; if not then according to Fejér's Theorem there exists some \mathbf{w}_1 such that $\|\mathbf{w}_1 - \mathbf{w}\| < \|\mathbf{w}_0 - \mathbf{w}\|$ for all $\mathbf{w} \in C$.

- We need to be able to certify if $\mathbf{w}_0 \in C$, which may be trivial when the set C is defined by *explicit* inequalities, such as $C = \{\mathbf{w} : g(\mathbf{w}) \leq 0\}$.
- If $\mathbf{w}_0 \notin C$, we need to be able to *explicitly and efficiently* find \mathbf{w}_1 .
- We also need sufficient decrease so that $\text{dist}(\mathbf{w}_t, C) \rightarrow 0$.
- We may also want to prove the convergence (rate) of the whole sequence \mathbf{w}_t .

Algorithmic Significance of Fejér's Result

Can be used to solve the convex feasibility problem:

$$\text{find } \mathbf{w} \in C,$$

where the closed (and convex) set $C \subseteq \mathbb{R}^d$ represents the solutions set of any problem. Indeed, starting from an arbitrary point \mathbf{w}_0 , if it is in C then we are done; if not then according to Fejér's Theorem there exists some \mathbf{w}_1 such that $\|\mathbf{w}_1 - \mathbf{w}\| < \|\mathbf{w}_0 - \mathbf{w}\|$ for all $\mathbf{w} \in C$.

- We need to be able to certify if $\mathbf{w}_0 \in C$, which may be trivial when the set C is defined by *explicit* inequalities, such as $C = \{\mathbf{w} : g(\mathbf{w}) \leq 0\}$.
- If $\mathbf{w}_0 \notin C$, we need to be able to *explicitly and efficiently* find \mathbf{w}_1 .
- We also need sufficient decrease so that $\text{dist}(\mathbf{w}_t, C) \rightarrow 0$.
- We may also want to prove the convergence (rate) of the whole sequence \mathbf{w}_t .

Let $C = \bigcap_{i \in I} C_i \neq \emptyset$. Suppose $\mathbf{w}_0 \notin C$ (otherwise we are done). Then there exists some $C_i \not\ni \mathbf{w}_0$. Apply the constructive part of Fejér's Theorem by letting

$$\mathbf{w}_1 = P_{C_i}(\mathbf{w}_0),$$

we immediately have

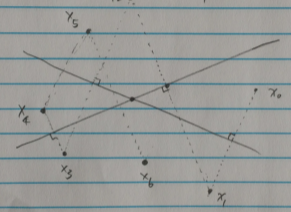
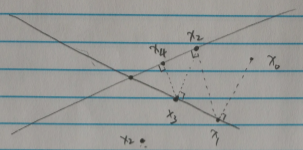
$$\forall \mathbf{w} \in C_i \supseteq C, \|\mathbf{w} - \mathbf{w}_1\|_2 < \|\mathbf{w} - \mathbf{w}_0\|_2.$$

Iterating the above idea leads to the method of alternating projections:

Algorithm 1: Method of alternating projections

Input: \mathbf{w}_0

```
1 for  $t = 0, 1, \dots$  do
2   choose set  $C_{i_t}$  // cyclic, random or greedy
3    $\mathbf{w}_{t+1} \leftarrow (1 - \eta_t)\mathbf{w}_t + \eta_t P_{C_{i_t}}(\mathbf{w}_t)$  //  $\eta_t \in [0, 2]$ 
```



Half Justification

Clearly, we have for any $\mathbf{w} \in C$:

$$\begin{aligned}\|\mathbf{w}_{t+1} - \mathbf{w}\|_2^2 &= \|\mathbf{w}_t - \mathbf{w} - \eta_t(\mathbf{w}_t - \mathbf{P}_{C_{i_t}}(\mathbf{w}_t))\|_2^2 \\ &= \|\mathbf{w}_t - \mathbf{w}\|_2^2 + (\eta_t^2 - 2\eta_t)\|\mathbf{w}_t - \mathbf{P}_{C_{i_t}}(\mathbf{w}_t)\|_2^2 + \\ &\quad 2\eta_t \langle \mathbf{w} - \mathbf{P}_{C_{i_t}}(\mathbf{w}_t), \mathbf{w}_t - \mathbf{P}_{C_{i_t}}(\mathbf{w}_t) \rangle \\ (\text{optimality of projection}) &\leq \|\mathbf{w}_t - \mathbf{w}\|_2^2 + (\eta_t^2 - 2\eta_t)\|\mathbf{w}_t - \mathbf{P}_{C_{i_t}}(\mathbf{w}_t)\|_2^2 \\ (\eta_t \in [0, 2]) &\leq \|\mathbf{w}_t - \mathbf{w}\|_2^2.\end{aligned}$$

Theorem: Convergence of alternating projections

Let $C = \bigcap_{i \in I} C_i \neq \emptyset$ where each C_i is closed and convex and $|I| < \infty$. If $0 < \alpha \leq \eta_t \leq 2 - \beta < 2$ for some $\alpha, \beta > 0$, then with the cyclic update order we have

$$\mathbf{w}_t \rightarrow \mathbf{w}_* \in C.$$

L. M. Bregman. "The method of successive projection for finding a common point of convex sets". *Soviet Mathematics Doklady*, vol. 162, no. 3 (1965), pp. 688–692, L. G. Gubin et al. "The Method of Projections for Finding the Common Point of Convex Sets". *USSR Computational Mathematics and Mathematical Physics*, vol. 7, no. 6 (1967), pp. 1–24. [English translation of paper in *Zh. Vychisl. Mat. mat. Fiz.* vol. 7, no. 6, pp. 1211–1228, 1967].

Alternating Bregman Projection

Instead of the Euclidean projection, can also consider the Bregman projection

$$\mathbb{P}_C(\mathbf{z}) = \mathbb{P}_{C,h}(\mathbf{z}) = \operatorname{argmin}_{\mathbf{w} \in C} D_h(\mathbf{w}, \mathbf{z}),$$

where $h : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ is a Legendre function.

Algorithm 2: Alternating Bregman projection

Input: \mathbf{w}_0 , $\operatorname{dom} h \supseteq C$

```
1 for  $t = 0, 1, \dots$  do
2   choose set  $C_{i_t}$  // cyclic, random or greedy
3    $\mathbf{w}_{t+1} \leftarrow (1 - \eta_t)\mathbf{w}_t + \eta_t \mathbb{P}_{C_{i_t}}(\mathbf{w}_t)$  //  $\eta_t \in [0, 2]$ 
```

L. M. Bregman. "A relaxation method of finding a common point of convex sets and its application to problems of optimization". *Soviet Mathematics Doklady*, vol. 171, no. 5 (1966), pp. 1578–1581.

Dykstra's algorithm

We now present a beautiful algorithm for solving:

$$\min_{\mathbf{w}} f(\mathbf{w}) \quad \text{s.t.} \quad \mathbf{w} \in C := \bigcap_{i \in I} C_i,$$

where f is Legendre and each C_i is closed and convex.

Algorithm 3: Dykstra's algorithm

Input: $\mathbf{w}_0 = \operatorname{argmin} f$, $\mathbf{a}_i = \mathbf{0}$, $b_i = 0$ for all $i \in I$

```
1 for  $t = 0, 1, \dots$  do
2   choose set  $C_{i_t}$  // cyclic, random or greedy
3    $\mathbf{w}_{t+1} \leftarrow \operatorname{argmin}_{\mathbf{w} \in C_{i_t}} f(\mathbf{w}) - \langle \mathbf{w}, \nabla f(\mathbf{w}_t) + \mathbf{a}_{i_t} \rangle$  // Bregman projection
4    $\mathbf{a}_{i_t} \leftarrow \mathbf{a}_{i_t} + \nabla f(\mathbf{w}_t) - \nabla f(\mathbf{w}_{t+1})$ 
5    $b_{i_t} \leftarrow \langle \mathbf{a}_{i_t, t+1}, \mathbf{w}_{t+1} \rangle$  // needed only for proof
```

Dykstra = AltMin in the Dual

Apply **Fenchel-Rockafellar duality** we obtain the dual problem:

$$\inf_{\{\mathbf{w}_i^*\}} f^* \left(- \sum_i \mathbf{w}_i^* \right) + \sum_i \sigma_i(\mathbf{w}_i^*),$$

where the (unique) primal solution \mathbf{w} and dual solution $\{\mathbf{w}_i^*\}$ are connected by:

$$\sum_i \mathbf{w}_i^* + \nabla f(\mathbf{w}) = \mathbf{0}.$$

- f is Legendre $\implies f^*$ is smooth and convex so AltMin applies

$$\mathbf{w}_{i,t+1}^* = \operatorname{argmin}_{\mathbf{w}_i^*} f^* \left(- \mathbf{w}_i^* - \sum_{j \neq i} \mathbf{w}_{j,t}^* \right) + \sigma_i(\mathbf{w}_i^*)$$

$$\text{or } \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in C_i} f(\mathbf{w}) + \left\langle \mathbf{w}; \sum_{j \neq i} \mathbf{w}_{j,t}^* \right\rangle$$

Dykstra = AltMin in the Dual

Apply **Fenchel-Rockafellar duality** we obtain the dual problem:

$$\inf_{\{\mathbf{w}_i^*\}} f^* \left(- \sum_i \mathbf{w}_i^* \right) + \sum_i \sigma_i(\mathbf{w}_i^*),$$

where the (unique) primal solution \mathbf{w} and dual solution $\{\mathbf{w}_i^*\}$ are connected by:

$$\sum_i \mathbf{w}_i^* + \nabla f(\mathbf{w}) = \mathbf{0}.$$

- f is Legendre $\implies f^*$ is smooth and convex so AltMin applies

$$\mathbf{w}_{i,t+1}^* = \operatorname{argmin}_{\mathbf{w}_i^*} f^* \left(- \mathbf{w}_i^* - \sum_{j \neq i} \mathbf{w}_{j,t}^* \right) + \sigma_i(\mathbf{w}_i^*)$$

$$\text{or } \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in C_i} f(\mathbf{w}) + \left\langle \mathbf{w}; \sum_{j \neq i} \mathbf{w}_{j,t}^* \right\rangle$$

The primal solution \mathbf{w}_{t+1} and dual solution $\mathbf{w}_{i,t+1}^*$ are now both unique due to the strict convexity in Legendre functions and they are connected by:

$$\nabla f(\mathbf{w}_{t+1}) + \mathbf{w}_{i,t+1}^* + \sum_{j \neq i} \mathbf{w}_{j,t}^* = \mathbf{0} = \nabla f(\mathbf{w}_{t+1}) + \sum_j \mathbf{w}_{j,t+1}^*, \quad (1)$$

since at time t we update $\mathbf{w}_{i,t+1}^*$ and keep $\mathbf{w}_{j,t+1}^* = \mathbf{w}_{j,t}^*$ for all $j \neq i$.

Let us define (and maintain)

$$\forall l = 1, \dots, |I|, \quad \mathbf{a}_{l,t} + \nabla f(\mathbf{w}_t) + \sum_{j \neq l} \mathbf{w}_{j,t}^* = \mathbf{0} \stackrel{(1)}{=} \mathbf{a}_{l,t} - \mathbf{w}_{l,t}^*,$$

where the last inequality follows from (1). Then,

$$\begin{aligned} \mathbf{a}_{i,t+1} &= \mathbf{w}_{i,t+1}^* \stackrel{(1)}{=} -\nabla f(\mathbf{w}_{t+1}) - \sum_{j \neq i} \mathbf{w}_{j,t}^* \stackrel{(1)}{=} -\nabla f(\mathbf{w}_{t+1}) + \mathbf{w}_{j,t}^* + \nabla f(\mathbf{w}_t) \\ &= \mathbf{a}_{i,t} + \nabla f(\mathbf{w}_t) - \nabla f(\mathbf{w}_{t+1}) \end{aligned}$$

while for all $l \neq i$, $\mathbf{a}_{l,t+1} = \mathbf{w}_{l,t}^* = \mathbf{a}_{l,t}$ since $\mathbf{w}_{l,t}^*$ was held fixed.

Entropy-regularized optimal transport

Let $\mathbf{p} \in \Delta_m$ and $\mathbf{q} \in \Delta_n$ be two probability vectors, and we seek a joint distribution $\Pi \in \mathbb{R}_+^{m \times n}$ with \mathbf{p} and \mathbf{q} as marginals such that the transportation cost is minimized:

$$\min_{\Pi \in \mathbb{R}_+^{m \times n}} \langle C, \Pi \rangle \quad \text{s.t.} \quad \Pi \mathbf{1} = \mathbf{p}, \quad \Pi^\top \mathbf{1} = \mathbf{q}.$$

Add a small entropy regularization:

$$\min_{\Pi \in \mathbb{R}_+^{m \times n}} \langle C, \Pi \rangle + \lambda \sum_{ij} \pi_{ij} \log \pi_{ij} \quad \text{s.t.} \quad \Pi \mathbf{1} = \mathbf{p}, \quad \Pi^\top \mathbf{1} = \mathbf{q}.$$

W.l.o.g. let $\Pi_0 \propto \exp(-C/\lambda) \geq \mathbf{0}$ and $\mathbf{1}^\top \Pi_0 \mathbf{1} = 1$ to obtain the equivalent problem:

$$\begin{aligned} \min_{\Pi \in \mathbb{R}_+^{m \times n}} \quad & \text{KL}(\Pi \| \Pi_0) \\ \text{s.t.} \quad & \Pi \mathbf{1} = \mathbf{p}, \quad \Pi^\top \mathbf{1} = \mathbf{q}. \end{aligned}$$

