



# بهینه‌سازی محدب برای مسائل گراف

محمدهادی فروغمنداعرابی  
زمستان ۱۳۹۵

## موضوع: کاهش گرادیان تصویر شده

جلسه سوم

نگارندگان: سهیلا فرخی و معصومه رحیمی

### ۱ مقدمه

در این جلسه الگوریتم کاهش گرادیان را بررسی می‌کنیم. این الگوریتم یک روش بهینه‌سازی مرتبه اول است که به ما کمک می‌کند از همواری تابع کمینه‌شده استفاده کنیم و نرخ همگرایی بهتری نسبت به الگوریتم کاهش زیرگرادیان (تصویر شده) که پیش از این مورد بررسی قرار گرفت بدست آوریم.

### ۲ مسئله شار بیشینه به عنوان یک مسئله بهینه‌سازی

مسئله بهینه‌سازی ترکیبیاتی موردنظر ما، مسئله شار بیشینه است. برای گراف بدون جهت  $G(V, E)$  داده شده با  $|V|$  رأس و  $|E|$  یال، و  $s, t \in V$  مسئله به این صورت فرموله می‌شود:

$$\begin{aligned} \min \quad & \|f\|_{\infty} \\ \text{s.t.} \quad & Bf = \chi_{s,t}, \end{aligned} \quad (1)$$

که  $B \in \mathbb{R}^{n \times m}$  یک ماتریس مجاورت رأس-یال از گراف  $G$  است که به صورت زیر تعریف می‌شود:

$$B_{v,e} := \begin{cases} -1 & \text{سر یال } e \text{ است } v \\ 1 & \text{ته یال } e \text{ است } v \\ 0 & \text{در غیر این صورت} \end{cases} \quad (2)$$

و  $\chi_{s,t} \in \mathbb{R}^V$  نیز به صورت زیر تعریف می‌شود:

$$\chi_{s,t}(v) := \begin{cases} -1 & v = s \\ 1 & v = t \\ 0 & \text{در غیر این صورت} \end{cases} \quad (3)$$

در اینجا می‌خواهیم جزئیات را کنار بگذاریم و مسئله را به صورت یک برنامه‌ریزی محدب کلی در نظر بگیریم، یعنی:

$$\min_{s.t.} f(x) \quad x \in K. \quad (4)$$

که  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  یک تابع محدب و  $K \subset \mathbb{R}^n$  یک مجموعه محدب است.

همانطور که پیش از این اشاره کردیم، می‌توانستیم از روش‌های کلی بهینه‌سازی محدب همچون الگوریتم بیضی برای حل این مسئله استفاده کنیم. اما، این الگوریتم تقریباً کند است. بنابراین، ما تلاش می‌کنیم از استراتژی‌های کاهش گرادیانی برای بدست آوردن کارایی بهتر استفاده کنیم، اگرچه این کار ما را به تقریب بدتری از جواب می‌رساند.

### ۳ الگوریتم‌های کاهش گرادیان و کاهش گرادیان تصویر شده

تا کنون با الگوریتم‌های کاهش زیرگرادیان و کاهش زیرگرادیان تصویر شده آشنا شدیم. حال فرض کنید تابع موردنظر مشتق‌پذیر بوده و مسئله نامقید است (یعنی داریم  $K = \mathbb{R}^n$ )، در این حالت می‌توانیم به جای زیرگرادیان در یک نقطه، از گرادیان در آن نقطه استفاده کنیم و همچنین به دلیل نامقید بودن به تصویرسازی نیازی نداریم. به این ترتیب الگوریتم کاهش گرادیان به فرم الگوریتم ۱ درخواهد آمد.

#### الگوریتم ۱ - الگوریتم کاهش گرادیان

$$\begin{aligned} x_1 &\leftarrow \vec{0} \\ \text{به ازای } s = 1, \dots, T-1 & \\ x_{s+1} &\leftarrow x_s - \eta \nabla f(x_s) \\ x_T &\text{ را برگردان} \end{aligned}$$

به عبارت دیگر، ما تخمین خود را با نقطه شدنی دلخواه  $x_1$ ، همچون بردار تماماً صفر، شروع می‌کنیم و  $T$  گام برای بهبود تخمین برمی‌داریم. در هر یک از این گام‌ها در جهت مخالف گرادیان تخمین کنونی،  $-\nabla f(x_s)$ ، حرکت می‌کنیم و اندازه هر گام به وسیله پارامتر  $\eta$  تنظیم می‌شود. از آنجا که گرادیان همواره در جهت تندترین شیب به سمت بالا است، می‌دانیم که گام برداشتن در جهت مخالف گرادیان بهترین بهبود موضعی در جهت کمینه‌سازی را سبب می‌شود.

توجه کنید که در روش کاهش گرادیان آخرین نقطه، یا  $x_T$  بازگشت داده می‌شود؛ درحالی‌که در روش کاهش زیرگرادیان، میانگین کل نقاط محاسبه‌شده، یا  $\bar{x}_T = \frac{1}{T} \sum_{s=1}^T x_s$  به عنوان خروجی الگوریتم ارائه می‌شود. در واقع، می‌توان در الگوریتم کاهش گرادیان نیز میانگین نقاط را به جای نقطه آخر به عنوان خروجی در نظر گرفت اما در این صورت همگرایی کندتر می‌شود. زیرا می‌دانیم وقتی یک دنباله از نقاط همگرا باشد، دنباله میانگین آنها نیز با نرخ کمتری همگرا خواهد بود.

اکنون، برای اینکه بتوانیم الگوریتم را با حالت مقید (یعنی هنگامی که  $K$  یک زیرمجموعه محض (محدب) از  $\mathbb{R}^n$  است) انطباق دهیم، می‌بایست از تصویرسازی استفاده کنیم.

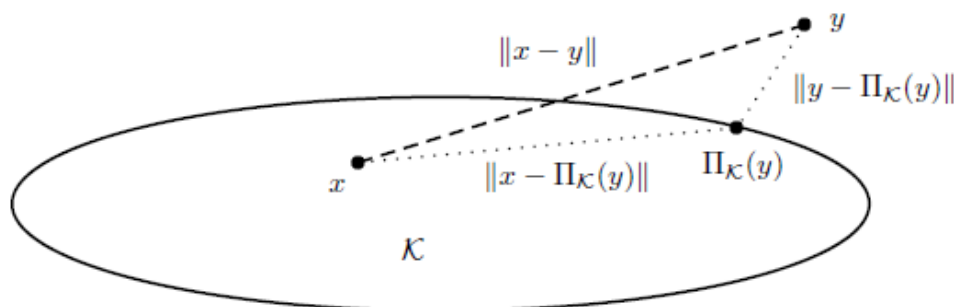
**تعریف ۱.** ( $\ell_2$ -تصویر) برای یک مجموعه محدب  $K \subset \mathbb{R}^n$  و یک نقطه  $y \in K$  تعریف می‌کنیم:

$$\Pi_K(y) := \operatorname{argmin}_{x \in K} \|x - y\|_2.$$

خاصیت مهم تصویرسازی در نکته زیر ارائه شده است.

**نکته ۲.** برای هر  $x \in K$  و  $y \in \mathbb{R}^n$  داریم  $(\Pi_K(y) - x)^T (\Pi_K(y) - y) \leq 0$ .

از نظر هندسی، نکته بالا به ما می‌گوید که زاویه بین خطوط تشکیل شده بین  $x$  و تصویر  $y$ ، و  $y$  و تصویر آن همواره منفرجه است - شکل ۱ را ببینید.



شکل ۱: نمایش نکته ۲.

اکنون با استفاده از تصویرسازی، استراتژی الگوریتم ۱ را اصلاح می‌کنیم تا برای حالت مقید نیز قابل اعمال باشد. الگوریتم کاهش گرادیان تصویر شده در الگوریتم ۲ ارائه شده است.

#### الگوریتم ۲ - الگوریتم کاهش گرادیان تصویر شده

$$\begin{aligned} x_1 &\leftarrow \vec{0} \\ \text{به ازای } s = 1, \dots, T-1: \\ x_{s+1} &\leftarrow \pi_K(x_s - \eta \nabla f(x_s)) \\ x_T &\text{ را برگردان} \end{aligned}$$

## ۴ - همواری $L$

به وضوح، برای اینکه الگوریتم کاهش گرادیان خوش تعریف باشد، می‌بایست تابع هدف  $f$  مشتق‌پذیر باشد. در غیر اینصورت، ممکن است گرادیان تابع  $f$  در برخی نقاط وجود نداشته باشد. در واقع، برای اینکه بتوانیم کران کمی دقیق‌تری برای همگرایی این الگوریتم ارائه کنیم، با نسخه کمی مشتق‌پذیری که در زیر معرفی می‌شود کار می‌کنیم.

**تعریف ۳.** ( $L$ -همواری) گوییم تابع  $f$ ، برای یک  $L \geq 0$ ،  $L$ -هموار است، اگر و تنها اگر

$$\forall x, y, \quad \|\nabla f(x) - \nabla f(y)\|_2 \leq L \|x - y\|_2.$$

خاصیت بالا می‌بایست با خاصیتی از  $f$  که برای آنالیز الگوریتم کاهش زیرگرادیان تصویر شده به آن نیاز داشتیم مقایسه شود. در آنجا، ما به یک کران  $G$  نیاز داشتیم که به عنوان یک ثابت لیبشیتز برای تابع  $f$  عمل می‌کرد، در حالی که در اینجا  $L$  یک ثابت لیبشیتز برای  $\nabla f$  است. این خاصیت به این موضوع اشاره دارد که وقتی از نقطه  $x$  کمی دور می‌شویم گرادیان نقطه جدید حداکثر  $L$  برابر آن از گرادیان در نقطه  $x$  دور می‌شود. یعنی، اگر به نقطه بهینه یا  $x^*$  نزدیک باشیم، چون در این نقطه گرادیان صفر است می‌بایست گرادیان کنونی نزدیک به صفر باشد. بنابراین،  $L$ -همواری ما را راهنمایی می‌کند که به نقطه بهینه نزدیک هستیم یا خیر. توجه کنید که توابع  $L$ -هموار مشتق‌پذیر هستند و تابع نرم بینهایت برای هیچ مقدار  $L$ ،  $L$ -هموار نیست. بنابراین، می‌بایست این مشکل را در آینده رفع کنیم.

**لم ۴.** اگر  $f$  محدب و  $L$ -هموار باشد، داریم

$$\forall x, y \in \mathcal{D}, \quad 0 \leq f(y) - f(x) - \nabla f(x)^T(y - x) \leq \frac{L}{2} \|x - y\|_2^2.$$

برای به دست آوردن شهودی در مورد این عبارت، توجه کنید که در عبارت  $f(x) + \nabla f(x)^T(y - x)$ ، عبارت مذکور فاصله تابع از این صفحه در نقطه  $y$  را نشان می‌دهد. بنابراین، نامساوی سمت چپ به این نکته اشاره دارد که تابع در هر نقطه  $y$  بالاتر از صفحه آفین مذکور قرار دارد و نامساوی سمت راست نیز یک کران بالا برای فاصله آنها از یکدیگر ارائه می‌کند. این کران بالا بیانگر این است که فاصله تابع از صفحه آفین ضریبی از مربع فاصله دو نقطه است و رشد بیش از حد ندارد. به عبارت دیگر، اگر از نقطه  $x$  زیاد دور نشویم تابع خیلی بالاتر از صفحه آفین گذرنده از  $(x, f(x))$  قرار نمی‌گیرد که به این معناست که این صفحه تخمین خوبی برای تابع است.

اثبات. نامساوی سمت چپ نتیجه مستقیمی از محدب بودن تابع است. بسط تیلور  $f(y)$  حول نقطه  $x$  را به یاد آورید:

$$f(y) = f(x) + \nabla f(x)^T(y - x) + \dots$$

محدب بودن  $f$  به این معناست که اگر جمع‌وندهای سوم به بعد را از سمت راست این معادله حذف کنیم، عبارت کاهش می‌یابد، یعنی

$$f(y) \geq f(x) + \nabla f(x)^T(y - x).$$

بنابراین، نامساوی سمت چپ را نتیجه می‌گیریم.

برای بدست آوردن نامساوی سمت راست از نامساوی کشی-شوارتز استفاده می‌کنیم و داریم:

$$\begin{aligned} |f(y) - f(x) - \nabla f(x)^T(y - x)| &= \left| \int_0^1 (\nabla f(x + t(y - x))^T(y - x) - \nabla f(x)^T(y - x)) dt \right| \\ &\leq \int_0^1 |\nabla f(x + t(y - x)) - \nabla f(x)|^T(y - x)| dt \\ &\leq \int_0^1 \|\nabla f(x + t(y - x)) - \nabla f(x)\| \cdot \|y - x\| dt \\ &\leq \int_0^1 L \|x + t(y - x) - x\| \cdot \|y - x\| dt \\ &= \frac{L}{2} \|y - x\|^2, \end{aligned}$$

□

که آخرین نامساوی مستقیماً از تعریف  $L$ -همواری نتیجه می‌شود.

## ۵ آنالیز الگوریتم کاهش گرادیان

اکنون می‌توانیم کارایی الگوریتم کاهش گرادیان را تحلیل کنیم. در اینجا توجه خود را به نسخه نامقید برنامه محدود می‌کنیم (یعنی  $K = \mathbb{R}^n$ ).

**قضیه ۵.** فرض کنید  $f$  یک تابع  $L$ -هموار است، اگر قرار دهیم  $\eta = \frac{1}{L}$ ، آنگاه خروجی  $x_T$  الگوریتم ۱ در نامساوی زیر صدق می‌کند:

$$f(x_T) - f(x^*) \leq O\left(\frac{L \cdot R^2}{T}\right),$$

که  $R = \|x_1 - x^*\|$ .

توجه کنید که تعریف  $R$  در اینجا از شعاع فضای شدنی  $K$  به فاصله تخمین اولیه از نقطه بهینه تبدیل شده است، زیرا در اینجا  $K = \mathbb{R}^n$  و فضای  $K$  فشرده نبوده و شعاع متناهی ندارد. پیش از پرداختن به جزئیات اثبات، به نکته زیر توجه کنید:

**مشاهده ۶.** برای هر  $s$ ,

$$f(x_s) - f(x_{s+1}) \geq \frac{1}{2L} \|\nabla f(x_s)\|_2^2.$$

ابتدا توجه کنید که این عبارت بیانگر این مسئله است که اگر در نقطه‌ای باشید که گرادیان تابع در آن نقطه مقدار کمی ندارد، فاصله از نقطه بعدی نیز خیلی کم نخواهد بود. یعنی در نقاط دور از نقطه بهینه، قدم‌های بزرگتری برداشته می‌شود چون گرادیان مقدار بزرگتری دارد و در نقاط نزدیک بهینه، به دلیل نزدیک شدن گرادیان به صفر، قدم‌های کوچکتری برداشته می‌شود و فاصله دو نقطه متوالی کاهش می‌یابد.

**اثبات.** از لم ۴ استفاده کرده و قرار می‌دهیم  $x = x_s$  و  $y = x_{s+1}$ . توجه کنید که  $\frac{1}{L} \nabla f(x_s) = \eta \nabla f(x_s) = x_s - x_{s+1}$  بنابراین داریم:

$$\begin{aligned} f(x_{s+1}) - f(x_s) + \nabla f(x_s)^T\left(\frac{1}{L} \nabla f(x_s)\right) &\leq \frac{L}{2} \left\| \frac{1}{L} \nabla f(x_s) \right\|_2^2 \\ \Rightarrow f(x_{s+1}) - f(x_s) + \frac{1}{L} \|\nabla f(x_s)\|_2^2 &\leq \frac{1}{2L} \|\nabla f(x_s)\|_2^2 \end{aligned}$$

□

که به سادگی نتیجه مطلوب از آن بدست می‌آید.

## ۱.۵ اثبات قضیه ۵

برای  $s = 1, \dots, T-1$ ، تعریف می‌کنیم:

$$\delta_s := f(x_s) - f(x^*).$$

داریم:

$$\delta_s = f(x_s) - f(x^*) \leq \nabla f(x_s)^T (x_s - x^*) \leq \|\nabla f(x_s)\|_2 \|x_s - x^*\|_2, \quad (۵)$$

که نامساوی سمت چپ با قرار دادن  $x = x_s$  و  $y = x^*$  در نامساوی سمت چپ لم ۴ نتیجه می‌شود و نامساوی سمت راست از نامساوی کشی-شوارتز بدست می‌آید. بنابراین:

$$\delta_s - \delta_{s+1} \geq \frac{1}{2L} \|\nabla f(x_s)\|_2^2 \geq \frac{1}{2L} \cdot \frac{\delta_s^2}{\|x_s - x^*\|_2^2}, \quad (۶)$$

که نامساوی سمت چپ مشاهده ۶ است و نامساوی سمت راست از معادله (۵) بدست می‌آید. اکنون مشکل ما عبارت  $\|x_s - x^*\|_2^2$  در مخرج کسر است. علاقمندیم که این عبارت را با  $R$  کراندار کنیم، اما تعریف  $R = \|x_1 - x^*\|_2$  به طور کلی نتیجه نمی‌دهد که برای هر  $s$  داریم  $\|x_s - x^*\|_2^2 \leq R^2$ . اگرچه در اینجا می‌توانیم این کران را با اتکا به  $L$ -هموار بودن تابع  $f$  و انتخاب خوب اندازه گام  $\eta = \frac{1}{L}$  بدست آوریم.

لم ۷. برای هر  $s$ ،  $\|x_s - x^*\|_2 \leq R$

اثبات لم بالا را پس از تکمیل اثبات قضیه ارائه می‌کنیم.

با استفاده از لم ۷ داریم:

$$\delta_s - \delta_{s+1} \geq \frac{\delta_s^2}{2LR^2},$$

و بنابراین:

$$\frac{1}{\delta_{s+1}} - \frac{1}{\delta_s} = \frac{\delta_s - \delta_{s+1}}{\delta_s^2} \geq \frac{1}{2LR^2}.$$

در عبارت بالا از این نکته استفاده کردیم که  $\delta_s \geq \delta_{s+1}$  بنابراین  $\delta_s - \delta_{s+1} \geq 0$ . با جمع کردن روی همه  $s = 1, \dots, T-1$ ، سمت چپ به صورت تلسکوپی ساده می‌شود و داریم:

$$\frac{1}{\delta_T} - \frac{1}{\delta_1} \geq \frac{T-1}{2LR^2}. \quad (۷)$$

اکنون  $\delta_1$  را کراندار می‌کنیم:

$$\delta_1 = f(x_1) - f(x^*) \leq \nabla f(x^*)^T (x_1 - x^*) + \frac{L}{2} \|x_1 - x^*\|_2^2 \leq \frac{LR^2}{2},$$

که نامساوی اول از لم ۴ و نامساوی دوم از این نکته که  $\nabla f(x^*) = 0$  و  $R = \|x_1 - x^*\|_2$  بدست می‌آید. با قرار دادن این کران در معادله (۷) داریم:

$$\delta_T \leq O\left(\frac{LR^2}{T}\right).$$

## ۲.۵ اثبات لم ۷

ما عبارت زیر را ثابت می‌کنیم که فوراً لم ۷ را نتیجه می‌دهد:

$$\forall s, \|x_{s+1} - x^*\|_2 \leq \|x_s - x^*\|_2. \quad (۸)$$

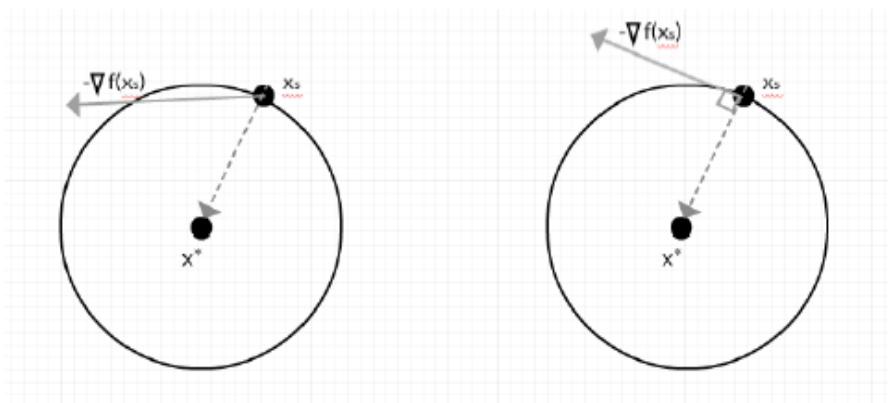
توجه کنید که این لم بیانگر این مسئله است که نیازی نیست فضای شدنی کوچک باشد، تنها کفایت نقطه اولیه از نقطه بهینه فاصله زیادی نداشته باشد، چون همواره با شروع از هر نقطه در هر گام به نقطه بهینه نزدیک می‌شویم.

معادله (۸) به طور کلی درست نیست (بدون تکیه به  $L$ -همواری و اندازه گام). گوی به مرکز  $x^*$  و شعاع  $\|x_s - x^*\|_2$  را در نظر بگیرید. ما در نقطه  $x_s$  ایستاده‌ایم و می‌خواهیم به نقطه  $x_{s+1}$  برویم. معادله (۸) در صورتی برقرار است که این گام ما را داخل گوی مذکور نگاه دارد. جهت

که در واقع به سمت آن حرکت می‌کنیم  $-\nabla f(x_s)$  است، در حالی که جهت درست به سمت  $x^*$ ، یعنی مرکز گوی، است. محدب بودن  $f$  تضمین می‌کند که زاویه بین جهت درست و جهت واقعی نمی‌تواند منفرجه باشد. به طور مشخص، از شرط زیرگرادیان

$$f(x) - f(y) \leq \nabla f(x)^T(x - y)$$

و جایگذاری  $x = x_s$  و  $y = x^*$  و توجه به اینکه همواره  $f(x_s) - f(x^*) \geq 0$  است، حاده بودن زاویه نتیجه می‌شود. اگرچه این زاویه همچنان می‌تواند (نزدیک به) عمود باشد. در این حالت، هر گام به اندازه کافی بزرگ در جهت  $-\nabla f(x_s)$  فاصله از نقطه بهینه را زیاد و معادله (۸) را نقض می‌کند. شکل ۲ را ببینید. بنابراین تنها گام‌های به اندازه کافی کوچک معادله (۸) را ارضاء می‌کنند.



شکل ۲

بنابراین می‌بایست نشان دهیم که  $\eta = \frac{1}{L}$  به گونه‌ای انتخاب شده است که اندازه گام برای ارضاء معادله (۸) مناسب باشد. این مسئله در لم زیر نشان داده شده است.

لم ۸. برای هر  $x, y \in \mathbb{R}^n$

$$(\nabla f(x) - \nabla f(y))^T(x - y) \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|_2^2.$$

برای اثبات معادله (۸)، از لم بالا و جایگذاری  $x = x_s$  و  $y = x^*$  داریم:

$$\nabla f(x_s)^T(x_s - x^*) \geq \frac{1}{L} \|\nabla f(x_s)\|_2^2, \quad (9)$$

که در آن از این حقیقت استفاده کردیم که  $\nabla f(x^*) = 0$ . اکنون با استفاده از این نامساوی داریم:

$$\begin{aligned} \|x_{s+1} - x^*\|_2^2 &= \|x_s - \frac{1}{L} \nabla f(x_s) - x^*\|_2^2 \\ &= \|x_s - x^*\|_2^2 + \left\| \frac{1}{L} \nabla f(x_s) \right\|_2^2 - 2 \cdot \frac{1}{L} \nabla f(x_s)^T(x_s - x^*) \\ &\leq \|x_s - x^*\|_2^2 + \left\| \frac{1}{L} \nabla f(x_s) \right\|_2^2 - 2 \cdot \frac{1}{L^2} \|\nabla f(x_s)\|_2^2 \\ &\leq \|x_s - x^*\|_2^2, \end{aligned}$$

که لم ۷ را نتیجه می‌دهد.

## ۶ نتیجه‌گیری

در این جلسه، نشان دادیم که الگوریتم کاهش گرادیان نرخ همگرایی بهتری نسبت به الگوریتم کاهش زیرگرادیان (با کران همگرایی  $\frac{RG}{\sqrt{T}}$ ) دارد. زیرا بستگی معکوس به  $T$  از مجذور به خطی بهبود یافته است. اما این الگوریتم محدودیت‌هایی دارد: اول اینکه، الگوریتم برای مسائل نامقید که  $K = \mathbb{R}^n$  ارائه شد، که مسئله مورد نظر ما، شار بیشینه، را پوشش نمی‌دهد. برای اینکه به حالت مقید برگردیم، می‌بایست الگوریتم کاهش گرادیان



تصویرشده را آنالیز کنیم. نکته دوم این است که قضیه ۵ به مشتق‌پذیری تابع  $f$  وابسته است؛ حال آنکه مسئله شار بیشینه از نرم بینهایت استفاده می‌کند که در همه نقاط مشتق‌پذیر نیست. روش ما برای رفع این مشکل این است که تابع نرم بینهایت را با یک تابع هدف هموار تقریب بزنییم که باعث از دست رفتن دقت در بدست آوردن جواب بهینه می‌شود. چالش بوجود آمده ایجاد تعادل بین میزان از دست رفتن دقت در تقریب و همواری تابع هدف تقریبی است.