

Gradient descent

Summary

A classical problem of function minimization is considered.

$$x_{k+1} = x_k - \eta_k \nabla f(x_k) \quad (\text{GD})$$

- The bottleneck (for almost all gradient methods) is choosing step-size, which can lead to the dramatic difference in method's behavior.
- One of the theoretical suggestions: choosing stepsize inversely proportional to the gradient Lipschitz constant $\eta_k = \frac{1}{L}$
- In huge-scale applications the cost of iteration is usually defined by the cost of gradient calculation (at least $\mathcal{O}(p)$)
- If function has Lipschitz-continuous gradient, then method could be rewritten as follows:

$$\begin{aligned} x_{k+1} &= x_k - \frac{1}{L} \nabla f(x_k) = \\ &= \arg \min_{x \in \mathbb{R}^n} \left\{ f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|_2^2 \right\} \end{aligned}$$

Intuition

Direction of local steepest descent

Let's consider a linear approximation of the differentiable function f along some direction h , $\|h\|_2 = 1$:

$$f(x + \eta h) = f(x) + \eta \langle f'(x), h \rangle + o(\eta)$$

We want h to be a decreasing direction:

$$f(x + \eta h) < f(x)$$

$$f(x) + \eta \langle f'(x), h \rangle + o(\eta) < f(x)$$

and going to the limit at $\eta \rightarrow 0$:

$$\langle f'(x), h \rangle \leq 0$$

Also from Cauchy–Bunyakovsky–Schwarz inequality:

$$|\langle f'(x), h \rangle| \leq \|f'(x)\|_2 \|h\|_2 \rightarrow \langle f'(x), h \rangle \geq -\|f'(x)\|_2 \|h\|_2 = -\|f'(x)\|_2$$

Thus, the direction of the antigradient

$$h = -\frac{f'(x)}{\|f'(x)\|_2}$$

gives the direction of the **steepest local** decreasing of the function f .

The result of this method is

$$x_{k+1} = x_k - \eta f'(x_k)$$

Gradient flow ODE

Let's consider the following ODE, which is referred as Gradient Flow equation.

$$\frac{dx}{dt} = -f'(x(t))$$

and discretize it on a uniform grid with η step:

$$\frac{x_{k+1} - x_k}{\eta} = -f'(x_k),$$

where $x_k \equiv x(t_k)$ and $\eta = t_{k+1} - t_k$ - is the grid step.

From here we get the expression for x_{k+1}

$$x_{k+1} = x_k - \eta f'(x_k),$$

which is exactly gradient descent.

Necessary local minimum condition

$$\begin{aligned} f'(x) &= 0 \\ -\eta f'(x) &= 0 \\ x - \eta f'(x) &= x \\ x_k - \eta f'(x_k) &= x_{k+1} \end{aligned}$$

This is, surely, not a proof at all, but some kind of intuitive explanation.

Minimizer of Lipschitz parabola

Some general highlights about Lipschitz properties are needed for explanation. If a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is continuously differentiable and its gradient satisfies Lipschitz conditions with constant L , then $\forall x, y \in \mathbb{R}^n$:

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

which geometrically means, that if we'll fix some point $x_0 \in \mathbb{R}^n$ and define two parabolas:

$$\phi_1(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle - \frac{L}{2} \|x - x_0\|^2,$$

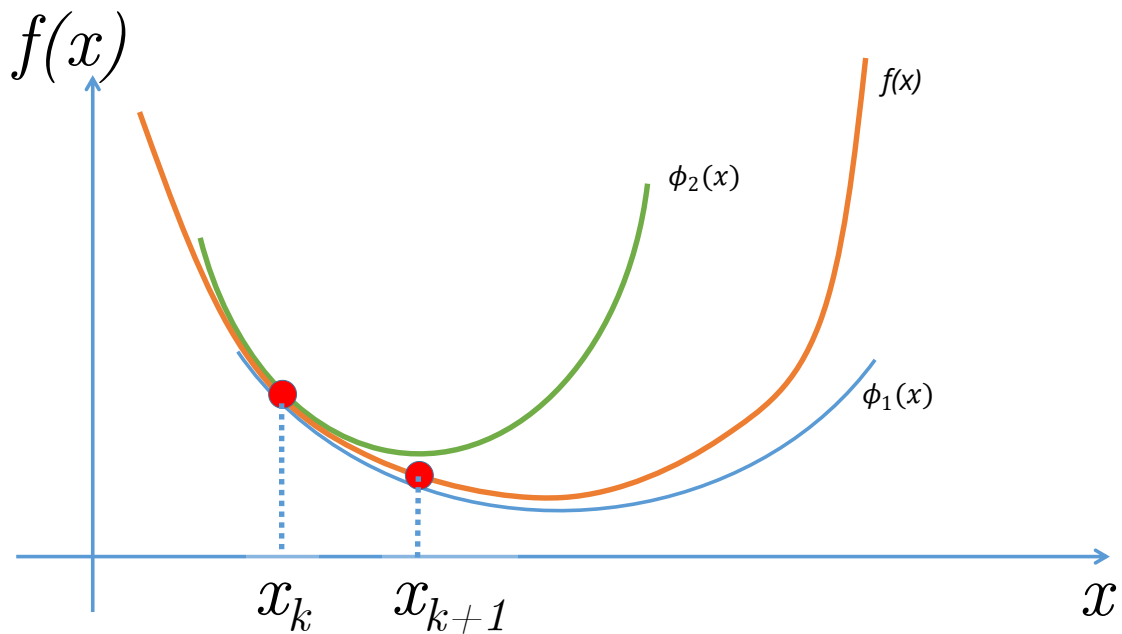
$$\phi_2(x) = f(x_0) + \langle \nabla f(x_0), x - x_0 \rangle + \frac{L}{2} \|x - x_0\|^2.$$

Then

$$\phi_1(x) \leq f(x) \leq \phi_2(x) \quad \forall x \in \mathbb{R}^n.$$

Now, if we have global upper bound on the function, in a form of parabola, we can try to go directly to its minimum.

$$\begin{aligned} \nabla \phi_2(x) &= 0 \\ \nabla f(x_0) + L(x^* - x_0) &= 0 \\ x^* &= x_0 - \frac{1}{L} \nabla f(x_0) \\ x_{k+1} &= x_k - \frac{1}{L} \nabla f(x_k) \end{aligned}$$



This way leads to the $\frac{1}{L}$ stepsize choosing. However, often the L constant is not known.

But if the function is twice continuously differentiable and its gradient has Lipschitz constant L , we can derive a way to estimate this constant $\forall x \in \mathbb{R}^n$:

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

or

$$-LI_n \preceq \nabla^2 f(x) \preceq LI_n$$

Stepsize choosing strategies

Stepsize choosing strategy η_k significantly affects convergence. General [link.html](#) title='Line search algorithms might help in choosing scalar parameter.'

Constant stepsize

For $f \in C_L^{1,1}$:

$$\eta_k = \eta$$

$$f(x_k) - f(x_{k+1}) \geq \eta \left(1 - \frac{1}{2}L\eta\right) \|\nabla f(x_k)\|^2$$

With choosing $\eta = \frac{1}{L}$, we have:

$$f(x_k) - f(x_{k+1}) \geq \frac{1}{2L} \|\nabla f(x_k)\|^2$$

Fixed sequence

$$\eta_k = \frac{1}{\sqrt{k+1}}$$

The latter 2 strategies are the simplest in terms of implementation and analytical analysis. It is clear that this approach does not often work very well in practice (the function geometry is not known in advance).

Exact line search aka steepest descent

$$\eta_k = \arg \min_{\eta \in \mathbb{R}^+} f(x_{k+1}) = \arg \min_{\eta \in \mathbb{R}^+} f(x_k - \eta \nabla f(x_k))$$

More theoretical than practical approach. It also allows you to analyze the convergence, but often exact line search can be difficult if the function calculation takes too long or costs a lot.

Interesting theoretical property of this method is that each following iteration is orthogonal to the previous one:

$$\eta_k = \arg \min_{\eta \in \mathbb{R}^+} f(x_k - \eta \nabla f(x_k))$$

Optimality conditions:

$$\nabla f(x_{k+1})^\top \nabla f(x_k) = 0$$

Goldstein-Armijo

This strategy of inexact line search works well in practice, as well as it has the following geometric interpretation:

Let's consider the following scalar function while being at a specific point of x_k :

$$\phi(\eta) = f(x_k - \eta \nabla f(x_k)), \eta \geq 0$$

consider first order approximation of $\phi(\eta)$:

$$\phi(\eta) \approx f(x_k) - \eta \nabla f(x_k)^\top \nabla f(x_k)$$

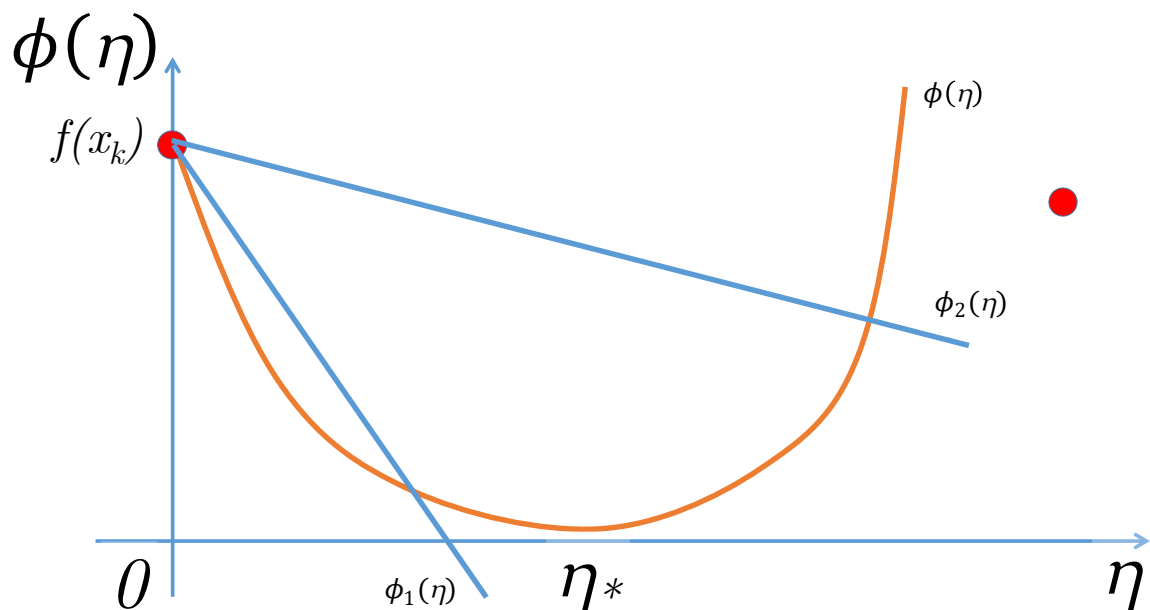
Let's consider also 2 linear scalar functions $\phi_1(\eta), \phi_2(\eta)$:

$$\phi_1(\eta) = f(x_k) - \alpha \eta \|\nabla f(x_k)\|^2$$

and

$$\phi_2(\eta) = f(x_k) - \beta \eta \|\nabla f(x_k)\|^2$$

Note, that Goldstein-Armijo conditions determine the location of the function $\phi(\eta)$ between $\phi_1(\eta)$ and $\phi_2(\eta)$. Typically, we choose $\alpha = \rho$ and $\beta = 1 - \rho$, while $\rho \in (0.5, 1)$



Convergence analysis

Quadratic case

Bounds

Conditions	$\ f(x_k) - f(x^*)\ \leq$	Type of convergence	$\ x_k - x^*\ \leq$
Convex Lipschitz-continuous function(G)	$\mathcal{O}\left(\frac{1}{k}\right) \frac{GR}{k}$	Sublinear	
Convex Lipschitz-continuous gradient (L)	$\mathcal{O}\left(\frac{1}{k}\right) \frac{LR^2}{k}$	Sublinear	
μ -Strongly convex Lipschitz-continuous gradient(L)		Linear	$(1 - \eta\mu)^k R^2$
μ -Strongly convex Lipschitz-continuous hessian(M)		Locally linear $R < \bar{R}$	$\frac{\bar{R}R}{\bar{R} - R} \left(1 - \frac{2\mu}{L + 3\mu}\right)$

- $R = \|x_0 - x^*\|$ - initial distance
- $\bar{R} = \frac{2\mu}{M}$

Materials

- [The zen of gradient descent. Moritz Hardt](#)
- [Great visualization](#)

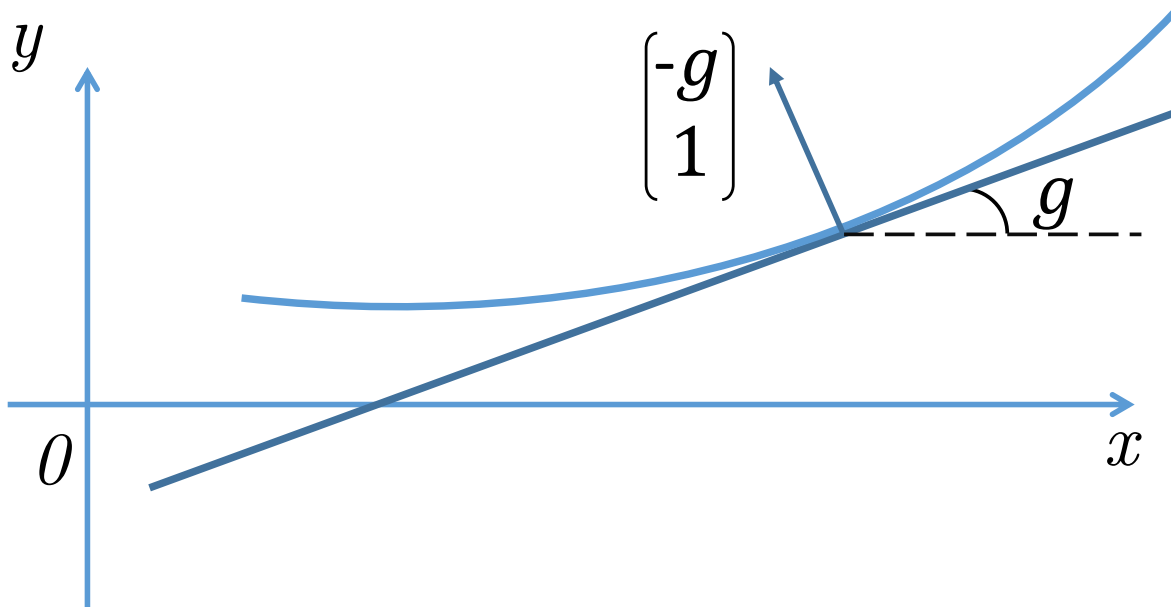
Subgradient and subdifferential

Motivation

Важным свойством непрерывной выпуклой функции $f(x)$ является то, что в выбранной точке x_0 для всех $x \in \text{dom } f$ выполнено неравенство:

$$f(x) \geq f(x_0) + \langle g, x - x_0 \rangle$$

для некоторого вектора g , то есть касательная к графику функции является *глобальной* оценкой снизу для функции.



- Если $f(x)$ - дифференцируема, то $g = \nabla f(x_0)$
- Не все непрерывные выпуклые функции дифференцируемы :)

Не хочется лишаться такого вкусного свойства.

Subgradient

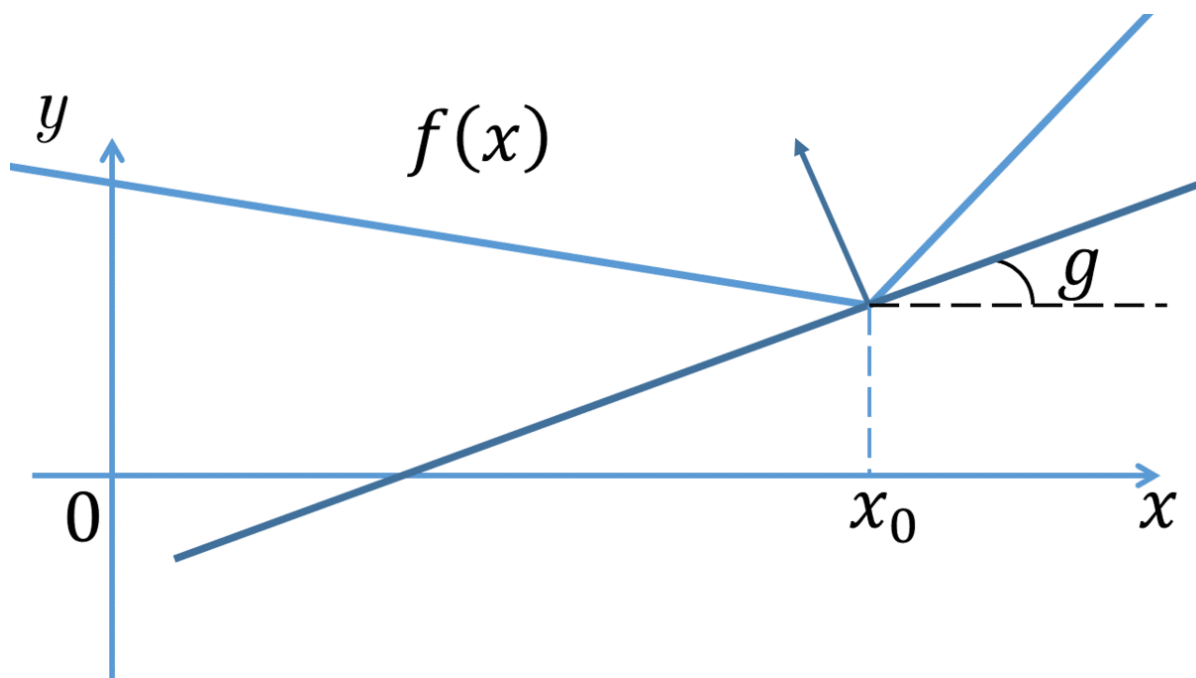
Вектор g называется **субградиентом** функции $f(x) : S \rightarrow \mathbb{R}$ в точке x_0 , если $\forall x \in S$:

$$f(x) \geq f(x_0) + \langle g, x - x_0 \rangle$$

Subdifferential

Множество всех субградиентов функции $f(x)$ в точке x_0 называется **субдифференциалом** f в x_0 и обозначается $\partial f(x_0)$.

- Если $x_0 \in \text{ri}S$, то $\partial f(x_0)$ выпуклое компактное множество.
- Выпуклая функция $f(x)$ дифференцируема в точке $x_0 \iff \partial f(x_0) = \nabla f(x_0)$
- Если $\partial f(x_0) \neq \emptyset \quad \forall x_0 \in S$, то $f(x)$ - выпукла на S .



Moreau - Rockafellar theorem (subdifferential of a linear combination)

Пусть $f_i(x)$ - выпуклые функции на выпуклых множествах S_i , $i = \overline{1, n}$.

Тогда, если $\bigcap_{i=1}^n \text{ri}S_i \neq \emptyset$ то функция $f(x) = \sum_{i=1}^n a_i f_i(x)$, $a_i > 0$ имеет субдифференциал

$\partial_S f(x)$ на множестве $S = \bigcap_{i=1}^n S_i$ и

$$\partial_S f(x) = \sum_{i=1}^n a_i \partial_{S_i} f_i(x)$$

Dubovitsky - Milutin theorem (subdifferential of a point-wise maximum)

Пусть $f_i(x)$ - выпуклые функции на открытом выпуклом множестве $S \subseteq \mathbb{R}^n$, $x_0 \in S$, а поточечный максимум определяется как $f(x) = \max_i f_i(x)$. Тогда:

$$\partial_S f(x_0) = \text{conv} \left\{ \bigcup_{i \in I(x_0)} \partial_S f_i(x_0) \right\},$$

где $I(x) = \{i \in [1 : m] : f_i(x) = f(x)\}$

Chain rule for subdifferentials

Пусть g_1, \dots, g_m - выпуклые функции на открытом выпуклом множестве $S \subseteq \mathbb{R}^n$, $g = (g_1, \dots, g_m)$ - образованная из них вектор - функция, φ - монотонно неубывающая выпуклая функция на открытом выпуклом множестве $U \subseteq \mathbb{R}^m$, причем $g(S) \subseteq U$. Тогда субдифференциал функции $f(x) = \varphi(g(x))$ имеет вид:

$$\partial f(x) = \bigcup_{p \in \partial \varphi(u)} \left(\sum_{i=1}^m p_i \partial g_i(x) \right),$$

где $u = g(x)$

В частности, если функция φ дифференцируема в точке $u = g(x)$, то формула запишется так:

$$\partial f(x) = \sum_{i=1}^m \frac{\partial \varphi}{\partial u_i}(u) \partial g_i(x)$$

Subdifferential calculus

- $\partial(\alpha f)(x) = \alpha \partial f(x)$, for $\alpha \geq 0$
- $\partial(\sum f_i)(x) = \sum \partial f_i(x)$, f_i - выпуклые функции
- $\partial(f(Ax + b))(x) = A^T \partial f(Ax + b)$, f - выпуклая функция
- $z \in \partial f(x)$ if and only if $x \in \partial f^*(z)$.

Examples

Концептуально, различают три способа решения задач на поиск субградиента:

- Теоремы Моро - Рокафеллара, композиции, максимума
- Геометрически
- По определению

1

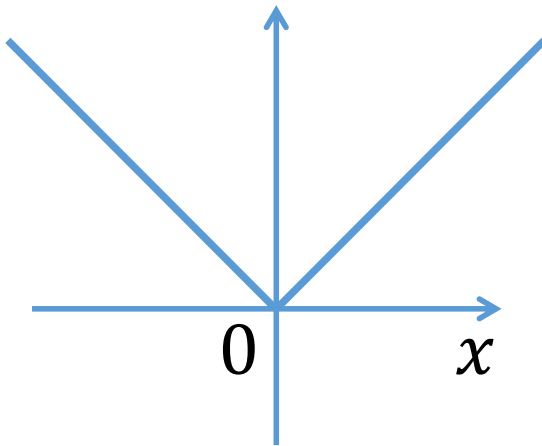
Найти $\partial f(x)$, если $f(x) = |x|$

Решение:

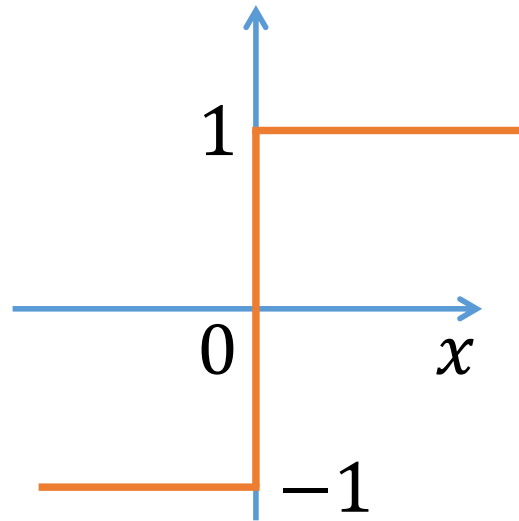
Решить задачу можно либо геометрически (в каждой точке числовой прямой указать угловые коэффициенты прямых, глобально подпирающих функцию снизу), либо по теореме Моро - Рокафеллара, рассмотрев $f(x)$ как композицию выпуклых функций:

$$f(x) = \max\{-x, x\}$$

$$f(x) = |x|$$



$$\partial f(x)$$



2

Найти $\partial f(x)$, если $f(x) = |x - 1| + |x + 1|$

Решение:

Совершенно аналогично применяем теорему Моро - Рокафеллара, учитывая следующее:

$$\partial f_1(x) = \begin{cases} -1, & x < 1 \\ [-1; 1], & x = 1 \\ 1, & x > 1 \end{cases} \quad \partial f_2(x) = \begin{cases} -1, & x < -1 \\ [-1; 1], & x = -1 \\ 1, & x > -1 \end{cases}$$

Таким образом:

$$\partial f(x) = \begin{cases} -2, & x < -1 \\ [-2; 0], & x = -1 \\ 0, & -1 < x < 1 \\ [0; 2], & x = 1 \\ 2, & x > 1 \end{cases}$$

3

Найти $\partial f(x)$, если $f(x) = [\max(0, f_0(x))]^q$. Здесь $f_0(x)$ - выпуклая функция на открытом выпуклом множестве S , $q \geq 1$.

Решение:

Согласно теореме о композиции (функция $\varphi(x) = x^q$ - дифференцируема), а $g(x) = \max(0, f_0(x))$ имеем: $\partial f(x) = q(g(x))^{q-1} \partial g(x)$

Тогда по теореме Дубовицкого - Милютина, в каждой точке $\partial f = \text{conv} \left(\bigcup_{i \in I(x)} \partial g_i(x) \right)$

$$\text{Заметим, что } \partial g(x) = \partial (\max\{s^\top x, -s^\top x\}) = \begin{cases} -s, & s^\top x < 0 \\ \text{conv}(-s; s), & s^\top x = 0. \\ s, & s^\top x > 0 \end{cases}$$

Причем, правило выбора "активной" функции поточечного максимума в каждой точке следующее:

- Если j -ая координата точки отрицательна, $s_i^j = -1$
- Если j -ая координата точки положительна, $s_i^j = 1$
- Если j -ая координата точки равна нулю, то подходят оба варианта коэффициентов и соответствующих им функций, а значит, необходимо включать субградиенты этих функций в объединение в теореме Дубовицкого - Милютина.

В итоге получаем ответ:

$$\partial f(x) = \{g : \|g\|_\infty \leq 1, \quad g^\top x = \|x\|_1\}$$

References

- [Lecture Notes for ORIE 6300: Mathematical Programming I by Darnek Davis](#)

Projection

Projection

Distance between point and set

The distance d from point $\mathbf{y} \in \mathbb{R}^n$ to closed set $S \subset \mathbb{R}^n$:

$$d(\mathbf{y}, S, \|\cdot\|) = \inf\{\|x - \mathbf{y}\| \mid x \in S\}$$

Projection of a point on set

Projection of a point $\mathbf{y} \in \mathbb{R}^n$ on set $S \subseteq \mathbb{R}^n$ is a point $\pi_S(\mathbf{y}) \in S$:

$$\|\pi_S(\mathbf{y}) - \mathbf{y}\| \leq \|\mathbf{x} - \mathbf{y}\|, \forall \mathbf{x} \in S$$

- if set is open, and a point is beyond this set, then its projection on this set does not exist.
- if a point is in set, then its projection is the point itself

- $$\pi_S(\mathbf{y}) = \underset{\mathbf{x}}{\text{argmin}} \|\mathbf{x} - \mathbf{y}\|$$

- Let $S \subseteq \mathbb{R}^n$ - convex closed set. Let the point $\mathbf{y} \in \mathbb{R}^n$ и $\pi \in S$. Then if for all $\mathbf{x} \in S$ the inequality holds:

$$\langle \pi - \mathbf{y}, \mathbf{x} - \pi \rangle \geq 0,$$

then π is the projection of the point \mathbf{y} на S , so $\pi_S(\mathbf{y}) = \pi$

- Let $S \subseteq \mathbb{R}^n$ - affine set. Let we have points $\mathbf{y} \in \mathbb{R}^n$ и $\pi \in S$. Then π is a projection of point \mathbf{y} на S , so $\pi_S(\mathbf{y}) = \pi$ if and only if for all $\mathbf{x} \in S$ the inequality holds:

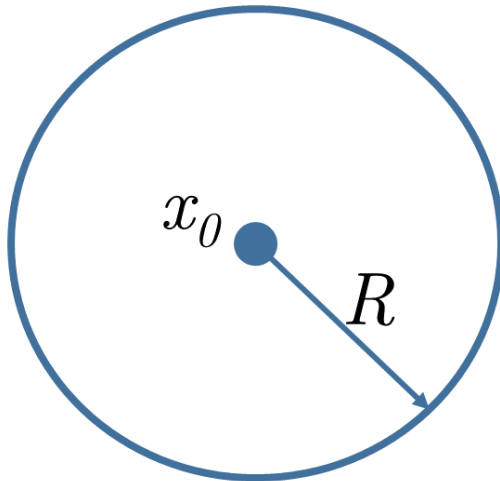
$$\langle \pi - \mathbf{y}, \mathbf{x} - \pi \rangle = 0$$

- **Sufficient conditions of existence of a projection.** If $S \subseteq \mathbb{R}^n$ - closed set, then for all points exist projection on set S .
- **Sufficient conditions of uniqueness of a projection.** Если $S \subseteq \mathbb{R}^n$ - convex set, then projection for all point on set S is unique (if exists).

Example 1

Find $\pi_S(y) = \pi$, if $S = \{x \in \mathbb{R}^n \mid \|x - x_0\| \leq R\}, y \notin S$

Solution:



- Build a hypothesis from the figure: $\pi = x_0 + R \cdot \frac{y - x_0}{\|y - x_0\|}$
- Check the inequality for a convex closed set: $(\pi - y)^T(x - \pi) \geq 0$

$$\begin{aligned} & \left(x_0 - y + R \frac{y - x_0}{\|y - x_0\|} \right)^T \left(x - x_0 - R \frac{y - x_0}{\|y - x_0\|} \right) = \\ & \left(\frac{(y - x_0)(R - \|y - x_0\|)}{\|y - x_0\|} \right)^T \left(\frac{(x - x_0)\|y - x_0\| - R(y - x_0)}{\|y - x_0\|} \right) = \\ & \frac{R - \|y - x_0\|}{\|y - x_0\|^2} (y - x_0)^T ((x - x_0)\|y - x_0\| - R(y - x_0)) = \\ & \frac{R - \|y - x_0\|}{\|y - x_0\|} \left((y - x_0)^T (x - x_0) - R\|y - x_0\| \right) = \\ & (R - \|y - x_0\|) \left(\frac{(y - x_0)^T (x - x_0)}{\|y - x_0\|} - R \right) \end{aligned}$$

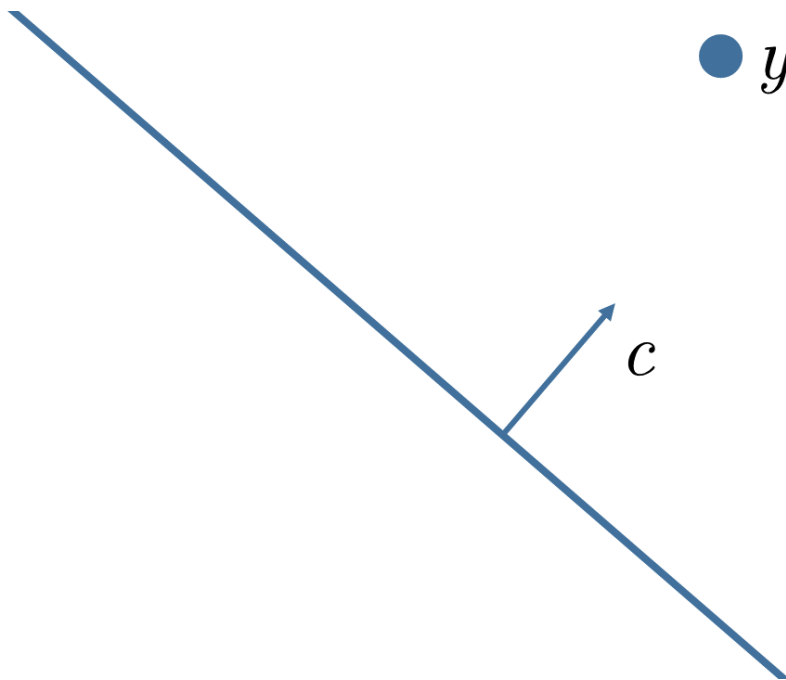
- The first factor is negative for point selection y . The second factor is also negative if we apply the Cauchy - Bunyakovsky theorem to its notation:

$$\begin{aligned} & (y - x_0)^T (x - x_0) \leq \|y - x_0\| \|x - x_0\| \\ & \frac{(y - x_0)^T (x - x_0)}{\|y - x_0\|} - R \leq \frac{\|y - x_0\| \|x - x_0\|}{\|y - x_0\|} - R = \|x - x_0\| - R \leq 0 \end{aligned}$$

Example 2

Find $\pi_S(y) = \pi$, if $S = \{x \in \mathbb{R}^n \mid c^T x = b\}$, $y \notin S$

Solution:



- Build a hypothesis from the figure: $\pi = y + \alpha c$. Coefficient α is chosen so that $\pi \in S$: $c^T \pi = b$, so:

$$c^T(y + \alpha c) = b$$

$$c^T y + \alpha c^T c = b$$

$$c^T y = b - \alpha c^T c$$

- Check the inequality for a convex closed set: $(\pi - y)^T(x - \pi) \geq 0$

$$(y + \alpha c - y)^T(x - y - \alpha c) =$$

$$\alpha c^T(x - y - \alpha c) =$$

$$\alpha(c^T x) - \alpha(c^T y) - \alpha^2(c^T c) =$$

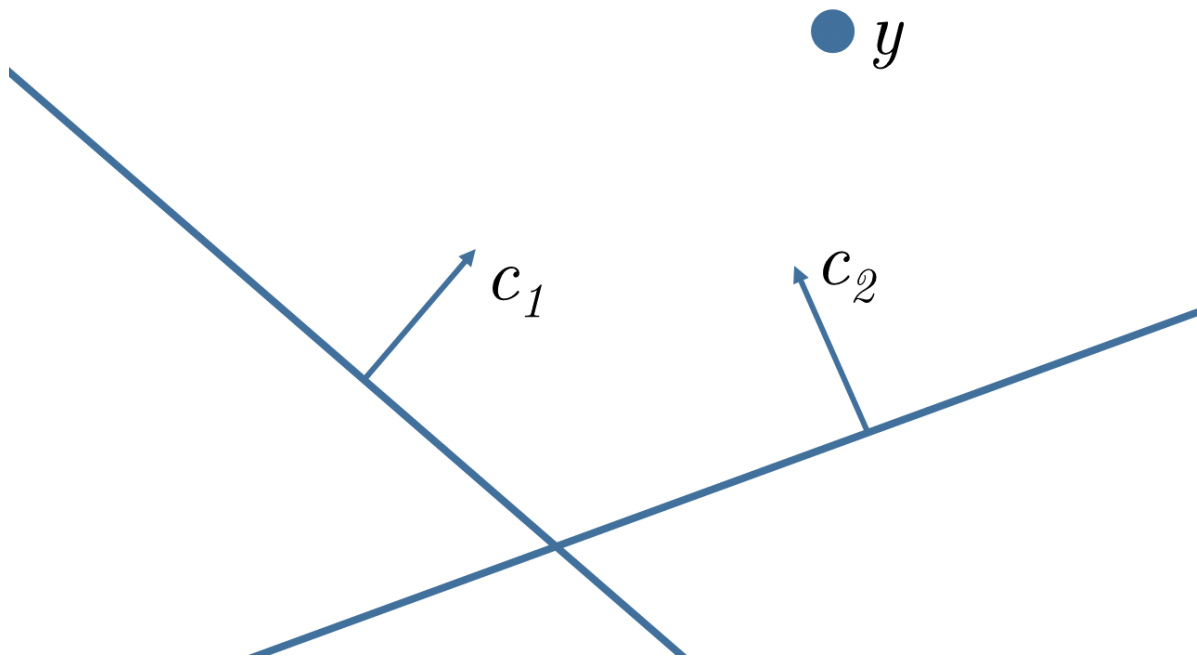
$$\alpha b - \alpha(b - \alpha c^T c) - \alpha^2 c^T c =$$

$$\alpha b - \alpha b + \alpha^2 c^T c - \alpha^2 c^T c = 0 \geq 0$$

Example 3

Find $\pi_S(y) = \pi$, if $S = \{x \in \mathbb{R}^n \mid Ax = b, A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m\}$, $y \notin S$

Solution:



- Build a hypothesis from the figure: $\pi = y + \sum_{i=1}^m \alpha_i A_i = y + A^T \alpha$. Coefficient α is chosen so that $\pi \in S: A\pi = b$, so:

$$A(y + A^T \alpha) = b$$

$$Ay = b - AA^T \alpha$$

- Check the inequality for a convex closed set: $(\pi - y)^T (x - \pi) \geq 0$

$$(y + A^T \alpha - y)^T (x - y - A^T \alpha) =$$

$$\alpha^T A(x - y - A^T \alpha) =$$

$$\alpha^T (Ax) - \alpha^T (Ay) - \alpha^T (AA^T \alpha) =$$

$$\alpha^T b - \alpha^T (b - AA^T \alpha) - \alpha^T AA^T \alpha =$$

$$\alpha^T b - \alpha^T b + \alpha^T AA^T \alpha - \alpha^T AA^T \alpha = 0 \geq 0$$