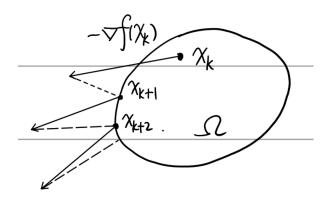
Projected Gradient Descent

1 Projection operator and projected gradient descent

To solve the inequality constrained problems, we introduce the *projected gradient* descent.

Recall the iteration step in the gradient descent method, $x_{k+1} = x_k - \eta \nabla f(x_k)$. Now we need to minimize f(x) over a feasible set Ω . If $x_k - \eta \nabla f(x_k)$ is feasible, then we can run the gradient descent iteration. If $x_k - \eta \nabla f(x_k)$ is infeasible, a simple idea is to project it onto Ω . This method is called the *projected gradient descent*.



Definition (*Projection*)

The projection of a point onto a set is the point in the set with minimum distance to the given point. Namely, the *projection operator* is defined by

$$\mathcal{P}_{\Omega}(oldsymbol{y}) = rg \min_{oldsymbol{x} \in \Omega} \|oldsymbol{x} - oldsymbol{y}\| \, .$$

The the projected gradient descent step can be given by

$$oldsymbol{x}_{k+1} = \mathcal{P}_{\Omega}ig(oldsymbol{x}_k - \eta\,
abla f(oldsymbol{x}_k)ig)\,.$$

Let

$$oldsymbol{g}(oldsymbol{x}) = rac{1}{\eta} \Big(oldsymbol{x} - \mathcal{P}_{\Omega} ig(oldsymbol{x} - \eta \,
abla f(oldsymbol{x}) \Big) \, \Big) \, ,$$

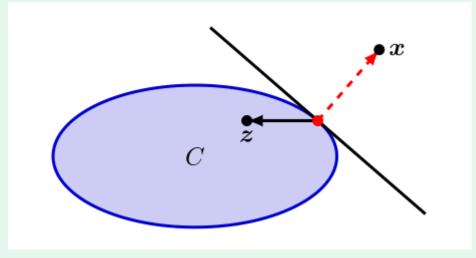
the iteration step can be expressed as

$$oldsymbol{x}_{k+1} = oldsymbol{x}_k - \eta \, oldsymbol{g}(oldsymbol{x}_k) \, .$$

Recall that, in <u>Lecture 4</u>, we show the following lemma.

Lemma

Let C be a nonempty, closed and convex set. Given \boldsymbol{x} and $\boldsymbol{y} = \mathcal{P}_C(\boldsymbol{x})$, for any $\boldsymbol{z} \in C$, it holds that $\langle \boldsymbol{x} - \boldsymbol{y}, \boldsymbol{z} - \boldsymbol{y} \rangle \leq 0$.



Conversely, if there exists $y \in C$ such that $\langle x - y, z - y \rangle \leq 0$, we have $y = \mathcal{P}_C(x)$. Otherwise, let $w = \mathcal{P}_C(x)$. Then we have

$$\langle oldsymbol{x} - oldsymbol{w}, oldsymbol{y} - oldsymbol{w}
angle \leq 0$$
 .

However, we also have $\langle {\boldsymbol x} - {\boldsymbol y}, {\boldsymbol w} - {\boldsymbol y} \rangle \leq 0$, which implies that

$$\langle \boldsymbol{x}-\boldsymbol{w}, \boldsymbol{w}-\boldsymbol{y} \rangle = \langle \boldsymbol{x}-\boldsymbol{y}, \boldsymbol{w}-\boldsymbol{y} \rangle + \langle \boldsymbol{y}-\boldsymbol{w}, \boldsymbol{w}-\boldsymbol{y} \rangle < 0$$

if $\boldsymbol{y} \neq \boldsymbol{w}$. Contradiction.

Thus, $\boldsymbol{y} = \mathcal{P}_C(\boldsymbol{x})$ if and only if $\langle \boldsymbol{x} - \boldsymbol{y}, \boldsymbol{z} - \boldsymbol{y} \rangle$ for any $\boldsymbol{z} \in C$.

Applying this lemma, we can show that g(x) plays a similar role as $\nabla f(x)$ in the gradient descent.

Lemma

For any $\boldsymbol{x} \in \Omega$,

$$\langle
abla f(oldsymbol{x}),\, oldsymbol{g}(oldsymbol{x})
angle \geq 0$$
 .

The inequality holds if and only if g(x) = 0.

Proof

Since $\boldsymbol{x} \in \Omega$, we have

$$\langle oldsymbol{x} - \mathcal{P}_{\Omega}(oldsymbol{x} - \eta \,
abla f(oldsymbol{x})), oldsymbol{x} - \eta \,
abla f(oldsymbol{x}) - \mathcal{P}_{\Omega}(oldsymbol{x} - \eta \,
abla f(oldsymbol{x}))
angle \leq 0$$
 ,

which gives that

$$\langle \eta \, oldsymbol{g}(oldsymbol{x}), \eta \, oldsymbol{g}(oldsymbol{x}) - \eta \,
abla f(oldsymbol{x})
angle = \eta^2 \, \langle oldsymbol{g}(oldsymbol{x}), \, oldsymbol{g}(oldsymbol{x}) -
abla f(oldsymbol{x})
angle \leq 0 \, .$$

Thus,

$$\langle
abla f(oldsymbol{x}), \, oldsymbol{g}(oldsymbol{x})
angle \geq \langle oldsymbol{g}(oldsymbol{x}), oldsymbol{g}(oldsymbol{x})
angle \, .$$

So we know that -g(x) is a desceding direction. Now we show that if g(x) = 0 then x is a minimum point.

Lemma

 $m{x}^*$ is a minimum point of f over Ω , iff $m{g}(m{x}) = m{0}$, namely, $m{x}^* = \mathcal{P}_{\Omega}(m{x}^* - \eta \, \nabla f(m{x}^*))$.

Proof

Applying the above lemma, we have ${\boldsymbol x}^* = \mathcal{P}_{\Omega}({\boldsymbol x}^* - \nabla f({\boldsymbol x}^*))$ if and only if

$$\langle oldsymbol{x}^* - \eta \,
abla f(oldsymbol{x}^*) - oldsymbol{x}^*, oldsymbol{z} - oldsymbol{x}^*
angle \leq 0$$

for all ${m z}\in\Omega,$ which is further equivalent to

$$\langle
abla f(oldsymbol{x}^*), oldsymbol{z} - oldsymbol{x}^*
angle \geq 0$$
 .

We conclude this lemma by the first-order optimality conditions of convex functions.

Hence, in the projected gradient descent, we can stop when $g(x_k)$ is small, or equivalently when $x_{k+1} - x_k$ is small.

Tip

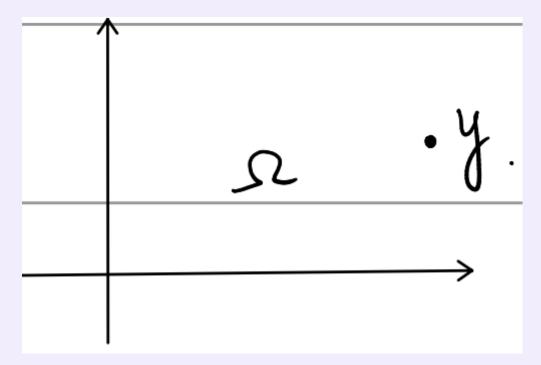
The projected gradient descent $\boldsymbol{x}_{k+1} = \mathcal{P}_{\Omega} (\boldsymbol{x}_k - \eta \nabla f(\boldsymbol{x}_k))$ can be also viewed as $\boldsymbol{x}_{k+1} = \arg\min_{\boldsymbol{z} \in \Omega} \langle \nabla f(\boldsymbol{x}_k), \boldsymbol{z} - \boldsymbol{x}_k \rangle + \frac{1}{2\eta} \|\boldsymbol{z} - \boldsymbol{x}_k\|^2$.

2 Examples of projection operator

Projected gradient descent is useful when the projection operator can be computed efficiently. Here we give some examples.

Example 1 (Box constraints)

$$\Omega = \{x \mid a_i \leq x_i \leq b_i, \quad i = 1, \cdots, n\}$$

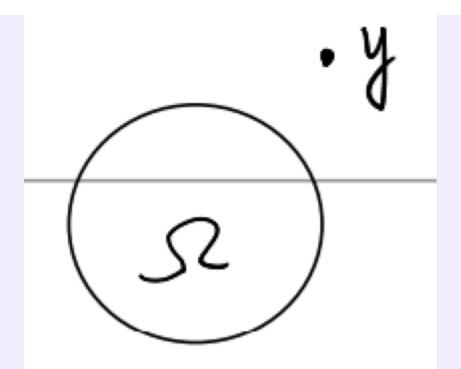


It is easy to see that

$$[\mathcal{P}_\Omega(y)]_i = \min\left\{b_i, \max\{a_i, y_i\}
ight\} = egin{cases} a_i & y_i < a_i \ y_i & a_i \leq y_i \leq b_i \ b_i & y_i > b_i \end{cases}$$

Example 2 (L^2 constraints, ridge regression)

$$\Omega = \{x \mid \|x\|_2 \leq t\}$$



The projection operator $\mathcal{P}_{\Omega}(y)$ is to compute

$$egin{array}{ll} \min & \left\|x-y
ight\|^2 \ \mathrm{subject\ to} & \left\|x
ight\|_2^2 \leq t^2 \end{array}$$

By KKT condition, there exists $\mu \geq 0$ such that

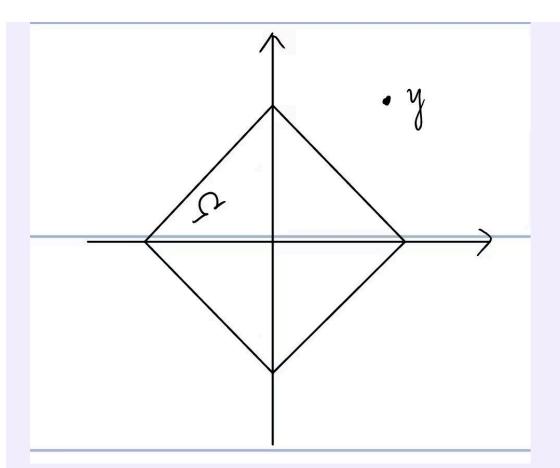
$$2(x-y) + 2\mu x = 0$$
 and $\mu(\|x\|^2 - t) = 0$

Then we have $y = (1 + \mu)x$.

Hence, $\mathcal{P}_{\Omega}(y) = \min\left\{1, rac{t}{\|y\|_2}
ight\}y$.

Example 3 (L^1 constraints, LASSO)

$$\Omega = \{x: \|x\|_1 \leq t\}$$



Unfortunately, there is no closed form for the projection operator $\mathcal{P}_{\Omega}(y)$. But we can compute it efficiently.

By symmetry, we only need to consider the case where $y_i \geq 0$ for all i. Now $\mathcal{P}_{\Omega}(y)$ is equivalent to the following optimization problem:

$$egin{array}{ll} \min & \|x-y\|^2 \ \mathrm{subject\ to} & \sum_i x_i \leq t \ & x_i \geq 0, orall\ i \ . \end{array}$$

By KKT condition, assume there exist KKT multipliers μ_0, \dots, μ_n such that

$$egin{cases} 2(x_i-y_i)+\mu_0-\mu_i=0, orall i\ \mu_0(\sum x_i-t)=0\ \mu_ix_i=0\ \sum x_i\leq t, x_i\geq 0 \end{cases}$$

- Case 1. $||y||_1 \le t$, then $\mu_0 = \mu_i = 0$. Hence x = y.
- Case 2. $\|y\|_1 > t$, then $\sum 2(x_i y_i) + \mu_0 \mu_1 = 2(\sum x_i \sum y_i) + n\mu_0 \sum \mu_i = 0$, hence $\mu_0 > 0$. Since $\mu_0(\sum x_i t) = 0$, we have $\sum x_i = t$.

• If
$$\mu_i = 0$$
, by $2(x_i - y_i) + \mu_0 - \mu_i = 0$, we have $x_i = y_i - \frac{1}{2}\mu_0$.

• If $\mu_i > 0$, by $\mu_i x_i = 0$, we have $x_i = 0$.

Now we have

$$x_i = egin{cases} y_i - rac{1}{2} \mu_0 & ext{if } y_i \geq rac{1}{2} \mu_0 \ 0 & ext{otherwise} \end{cases}$$

and
$$\sum x_i = t$$
.

We may use the binary search to find μ_0 , where the lower bound is 0 and the upper bound is max y_i .