

Lecture 12

Approximate Caratheodory Theorem using Mirror Descent

Caratheodory's Theorem (exact version)

Let $u \in \mathbb{R}^n$ ^{dimension of the space} be a point which lies in the convex hull of a set of points $\{v_1, v_2, \dots, v_m\}$. ^{number of points which define the convex hull}
Thus $u \in \text{conv}(\{v_1, \dots, v_m\})$.

Then u can be written as a convex combination of at most $n+1$ points of $\{v_1, \dots, v_m\}$.

That is $u = \sum_{i=1}^m \lambda_i v_i$, $\lambda_i \geq 0 \leq \sum \lambda_i = 1$
 $\# \text{non zero } \lambda_i \leq n+1$

Note that $n \leq m$

The Theorem is tight. That's why the natural question to ask is if can we less points if we permit an ϵ -error.

Approximate Carathéodory's Theorem

for any ℓ_p -norm with $p \geq 2$,
convex hull $\text{conv}(\{v_1, v_2, \dots, v_m\}) \subseteq \mathbb{R}^n$
with $\|v_i\|_p \leq 1$ ($v_i \in B_p^n \forall i \in [m]$), and
point $u \in \text{conv}(\{v_1, \dots, v_m\})$ there exists
 \tilde{u} s.t. $\|u - \tilde{u}\|_p \leq \varepsilon$ where \tilde{u} is in
the convex hull of at most $\frac{4p}{\varepsilon^2}$
points of $\{v_1, v_2, \dots, v_m\}$. That means

just to denote a permutation of the set $\{v_1, \dots, v_m\}$

$$\tilde{u} \in \text{conv}(\{\tilde{v}_1, \dots, \tilde{v}_{\frac{4p}{\varepsilon^2}}\}) \subseteq \text{conv}(\{v_1, \dots, v_m\})$$

Observations

- ① if $\max \|v_i\|_p \leq d$ then the # of points $\leq \frac{4pd^2}{\varepsilon^2}$
- ② $\frac{4p}{\varepsilon^2}$ is independent of the dimension n of the space.
- ③ the original proof of the theorem is probabilistic and non-constructive.

Optimization formulation

$$V = \begin{bmatrix} | & | & & | \\ v_1 & v_2 & \dots & v_m \\ | & | & & | \end{bmatrix}, \quad x \in \Delta^m \leftarrow \text{simplex}$$

we want to solve $\min \|Vx - u\|_p$

s.t. x is sparse
 $x \in \Delta^m$

Problem: we do not know how to model the sparsity requirement

Idea: solve the unconstrained problem iteratively using gradient descent / mirror descent and hope that the solution is sparse.

↙ This idea will fail since the gradient update may not be sparse

we need to formulate the problem differently.

Reformulation using Sion's minimax theorem

Remember that the dual norm of $\|\cdot\|_p$ is $\|\cdot\|_q$ where $\frac{1}{p} + \frac{1}{q} = 1$. From the dual norm definition we get that

$$\begin{aligned}\|Vx - u\|_p &= \max \{ \langle y, Vx - u \rangle \mid \|y\|_q \leq 1 \} = \\ &= \max \{ \langle y, Vx - u \rangle \mid y \in B_q^u \}\end{aligned}$$

$$\Rightarrow \min_{x \in \Delta^u} \|Vx - u\|_p = \min_{x \in \Delta^u} \max_{y \in B_q^u} \langle y, Vx - u \rangle =$$

Sion's

$\xRightarrow{\text{with max theorem}}$

$$\max_{y \in B_q^u} \min_{x \in \Delta^u} \langle y, Vx - u \rangle =$$

$$= \max_{y \in B_q^u} \left(- \max_{x \in \Delta^u} \langle y, u - Vx \rangle \right) =$$

$$= - \min_{y \in B_q^u} \max_{x \in \Delta^u} \langle y, u - Vx \rangle =$$

$$= - \min_{y \in B_q^u} f(y)$$

Now the unconstrained problem is equivalent to solving:

$$\min_{y \in \mathbb{B}^q} f(y), \quad f(y) = \max_{x \in \Delta^m} \langle y, u - Vx \rangle$$

$$= \langle y, u - Vx^y \rangle$$

$$x^y = \arg \max_{x \in \Delta^m} \langle y, u - Vx \rangle$$

Let's calculate some (sub)gradients:

$$x^y = \arg \max_{x \in \Delta^m} -\langle V^T y, x \rangle \Rightarrow x^y \text{ is the vector of all zeros except for a 1 at the coordinate } i = \arg \max_i \{-u_i^T y\}$$

we abuse notation even if it is a subgradient

$$\nabla f(y) = u - Vx^y = u - u_i$$

proof

$$f(z) \geq f(y) + \langle u - Vx^y, z - y \rangle \Rightarrow$$

$$\Rightarrow f(z) \geq f(y) - \langle u - Vx^y, y \rangle + \langle u - Vx^y, z \rangle \Rightarrow$$

$$\Rightarrow \max_{x \in \Delta^m} \langle z, u - Vx \rangle \geq \langle z, u - Vx^y \rangle \quad \square$$

this guy (which is a dual certificate of the new optimization formulation) is the primal

U.P. Lemma

$$D\varphi(x, y) \geq \alpha \|x - y\|_q^2$$

f is α -strongly convex w.r.t $\|\cdot\|_p$ then with an appropriate stepsize we get

$$\frac{1}{T} \sum \langle \nabla f(y_t), y_t - y^* \rangle \leq 2 \frac{D\varphi(x^*, x_0)}{\sqrt{\alpha T}}, \forall y^* \in B_q$$

in the original statement we have

$$\frac{1}{T} \leq (f(y_t) - f(y^*))$$

because of the first order condition and because y^* is a minimizer

$$\frac{1}{T} \sum \langle \nabla f(y_t), y_t - y^* \rangle = \frac{1}{T} \sum \langle u - v_{i(t)}, y_t - y^* \rangle =$$

$$= \frac{1}{T} \sum (\underbrace{\langle u - v_{i(t)}, y_t \rangle}_{f(y_t)} - \langle u - v_{i(t)}, y^* \rangle) =$$

$$= \frac{1}{T} \sum f(y_t) - \frac{1}{T} \sum \langle u - v_{i(t)}, y^* \rangle,$$

$$\geq -\frac{1}{T} \sum \langle u - v_{i(t)}, y^* \rangle, \text{ since } f(y_t) \geq 0$$

$(f(y)) = \max_{x \in \Delta^m} \langle y, u - vx \rangle \geq 0$
because $\exists x^* \text{ s.t. } u = vx^*$

So, using mirror descent we can get a bound of the form

$$-\frac{1}{T} \sum \langle u - v_{i(t)}, y \rangle \leq \dots \quad \forall y \in B_q$$

$$\Rightarrow \langle y, \frac{1}{T} \sum v_{i(t)} - u \rangle \leq \dots \quad \forall y \in B_q$$

by definition of the dual norm

$$\left\| \frac{1}{T} \sum v_{i(t)} - u \right\|_p \leq \dots$$

$$\left(\frac{1}{p} + \frac{1}{q} = 1 \right)$$

(just a reminder)

Now it remains to calculate

the Lipschitz parameter L and find

$$\text{a function } \psi \text{ s.t. } D_\psi(x, y) \geq \alpha \|x - y\|_q^2$$

$$\textcircled{1} \text{ we have } \|\nabla f(y_t)\|_p = \|u - v_{i(t)}\|_p \leq$$

$$\leq \|u\|_p + \|v_{i(t)}\|_p \leq 2 \quad (\text{since } v_i \in B_q \text{ and } u \text{ belongs to the convex hull of } \{v_1, \dots, v_m\})$$

$$\text{so } L = 2$$

② For $1 \leq q \leq 2$, $\psi(y) = \frac{1}{2} \|y\|_q^2$, $y \in \mathcal{B}^q$

(why should we try this function?)

well. In the case of the classic l_2 -norm

we have $\psi(y) = \frac{1}{2} \|y\|_2^2 \Rightarrow D_\psi(x, y) = \frac{1}{2} \|x - y\|_2^2$

and now we need a bound of the

form $D_\psi(x, y) \geq \alpha \|x - y\|_q^2$

For the domain \mathcal{B}^q we have

$$D_\psi(x, y) \geq \frac{q-1}{2} \|x - y\|_q^2 \text{ and}$$

$$D_\psi(y^*, 0) \leq \max_{y \in \mathcal{B}^q} D_\psi(y, 0) \leq 1/2 \quad (\text{for } q=2)$$

$$\text{so } D_\psi(x, y) \geq \alpha \|x - y\|_q^2, \quad \alpha = \frac{q-1}{2} =$$

$$\frac{\frac{1}{q} + \frac{1}{p} = 1}{\frac{p}{p-1} - 1} = \frac{p - p+1}{2(p-1)} =$$

$$\text{So to get } \varepsilon\text{-close we need } = \frac{1}{2(p-1)}$$

$$\left(\frac{4 \cdot 2^2 D_\psi(y^*, 0)}{\alpha \varepsilon^2} \right) \leq \frac{4 \cdot 2^2 \cdot 1/2}{\frac{1}{2(p-1)} \varepsilon^2} = O(p/\varepsilon^2) \text{ iterations}$$

which is also the # of v_i s that we need to use