# 机器学习里的数学第7课凸优化进阶

管老师

七月在线

June, 2016

### 主要内容

- 共轭函数
  - 共轭函数
- 对偶问题
  - 拉格朗日对偶函数
  - 拉格朗日对偶问题
  - 共轭函数与拉格朗日对偶函数
- 对偶性
  - 弱对偶性与强对偶性
  - 强对偶性成立的几种情况
  - 凸优化问题求解 (KKT)
- 应用举例
  - 支持向量机 (SVM) 的最简单形式
- 总结寄语 (数学部分)

### 记号

• 本节课常用数学记号

V, W 向量空间

v, w 向量

 $\mathbb{R}^n$ ,  $\mathbb{R}^m$  实坐标空间

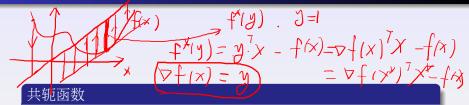
 $\alpha,\beta$  V 和 W 的基

 $T: V \to W$  向量空间 V 到 W 的线性映射

 $A_{\alpha,\beta}(T)$  线性映射 T 在  $\alpha$  和  $\beta$  这两组基下的矩阵

 $G(v_1, v_2)$  内积空间 V 上的内积

 $H_{\alpha}$  G 在基  $\alpha$  下的矩阵形式



如果  $f: \mathbb{R}^n \to \mathbb{R}$  是一个函数, 那么 f 的共轭函数

$$f^*(y) = \sup_{x \in \mathsf{dom} f} (y^T x - f(x))$$

其中  $f^*(y)$  的定义域是使得等式右边有上界的那些 y.

$$f_{2}$$
 =  $\begin{cases} xy - f^{*}(y) \\ y \in f_{0} \\ y = f_{1} \end{cases}$   $\begin{cases} z^{T}y - f^{*}(y) \\ y^{T}x - f_{1}x \end{cases}$   $\begin{cases} x^{T}y - f^{T}(y) \\ y^{T}x - f_{1}x \end{cases}$ 



#### 共轭函数的基本性质

- 共轭函数  $f^*$  是一个凸函数
- 如果 g 是 f 的凸闭包,那么  $g^* = f^*$
- 对一般的函数  $f, f^{**} \leq f$
- 如果 f 是一个凸函数,那么  $f^{**} = f$

#### 共轭函数的进一步性质

- $f(x) + f^*(y) \ge x^T y$
- 如果 f 是凸函数而且可微,那么  $f^*(y) = x^{\mathsf{T}} \mathcal{T} \nabla f(x^*) f(x^*)$ ,其中  $x^*$  满足  $\nabla f(x^*) = y$ .
- 如果 g(x) = f(Ax + b), 则  $g^*(y) = f^*(A^{-T}y) b^T A^{-T}y$ .
- 如果  $f(u,v) = f_1(u) + f_2(v)$ , 那么  $f^*(w,z) = f_1^*(w) + f_2^*(z)$

$$f(x) = x \ln x$$

$$f(x) = y \ln x - f(x) = y - \ln x - \frac{x}{x} = y - l - l - x$$

$$f(y) = y - l - x - \frac{x}{x} = y - l - l - x$$

$$f(y) = y - e^{y - l} (e^{y - y})$$

$$= y e^{y - l} - e^{y - l} (y - l)$$

$$= (e^{y - l})$$

考虑  $\mathbb{R}^n$  上的优化问题:

可约数三面介绍的

### 优化问题

最小化:
$$f_0(x)$$
  
不等条件: $f_i(x) \leq \sum_{i=1}^m i = 1, \dots, m$   
等式条件: $h_i(x) = 0, i = 1, \dots, p$   
定义域: $\mathcal{D} = \bigcap_{i=0}^m \mathsf{dom} f_i \cap \bigcap_{i=1}^p \mathsf{dom} h_i$ .

请注意定义域  $\mathcal{D}$  指的是使得所有函数  $f_i, h_i$  有定义的区域。而可行域指的是定义域中满足不等条件与等式条件的那些点. 本课中把这个优化问题称为原问题,优化点称为  $x^*$ , 最优化值为  $p^*$ .

根据原函数与限制条件我们定义拉格朗日量  $L(x,\lambda,\nu):\mathbb{R}^{n+m+p}\to\mathbb{R}$ 

#### 拉格朗日量

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} \nu_i h_i(x),$$

根据拉格朗日函数我们定义拉格朗日对偶函数  $g(\lambda, \nu): \mathbb{R}^{m+p} \to \mathbb{R}$ 

#### 拉格朗日对偶函数

$$g(\lambda, \nu) = \inf_{x \in \mathcal{D}} L(x, \lambda, \nu)$$
$$= \inf_{x \in \mathcal{D}} \left( f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x) \right)$$

拉格朗日量在数学与物理中有极为广泛的应用,感兴趣的同学可以在"信息几何学","理论力学"中找到它其他的应用.

对偶函数有如下重要性质

#### 对偶函数为原问题提供下界

如果限制  $\lambda_i \geq 0, \forall i = 1, \dots, m$ , 则

$$g(\lambda, \nu) \le p^*$$

#### 证明

对任意一个  $x \in \mathcal{D}$ , 如果 x 在可行域中, 那么

$$g(\lambda, \nu) \le f_0(x) + \sum_{i=1}^m \lambda_i f_i(x) + \sum_{i=1}^p \nu_i h_i(x)$$
$$= f_0(x) + \sum_{i=1}^m \lambda_i f_i(x)$$
$$\le f_0(x)$$

### 对偶问题:对偶问题

根据对偶函数, 定义对偶问题的一般形式

### 对偶问题

最大化: $g(\lambda, \nu)$ 不等条件: $\lambda_i \geq 0, i = 1, \cdots, m$ 

我们把对偶问题的最大值点称为  $(\lambda^*, \nu^*)$ , 相应的最大值称为  $d^*$ , 这里面的对偶函数 g 定义域为  $\mathsf{dom} g = \{(\lambda, \nu) : g(\lambda, \nu) > -\infty\}$ . 在 g 的定义域中满足  $\lambda_i \geq 0$  的那些  $(\lambda, \nu)$  全体,叫做对偶可行域. 也就是对偶问题的可行域.

#### 对偶问题举例

原问题, 线性规划

最小化: $c^T x$ 

条件:Ax = b, 而且 $x_i \ge 0$ ,  $\forall i = 1, \dots, n$ 

$$f_{o} = C^{T} \times,$$

$$f_{o}(x) = -X \le 0$$

$$f_{o}(x) = 0, \quad 0$$

$$A = \begin{bmatrix} 0, & 0 \\ h_{o}(x) \end{bmatrix}$$

$$A^{T} \cdot X = \begin{bmatrix} 0 \\ h_{o}(x) \\ h_{o}(x) \end{bmatrix}$$

$$A^{T} \cdot X = \begin{bmatrix} 0 \\ h_{o}(x) \\ h_{o}(x) \\ h_{o}(x) \end{bmatrix}$$

$$\frac{L(x,\lambda,\mu)}{=c^{T}x+\lambda^{T}(x)+\nu^{T}(A^{T}x-b)}$$

$$=\frac{1}{2}(\lambda,\nu)=\inf\left(\frac{1}{2}(c^{T}x+\lambda^{T}x+\nu^{T}A^{T}x-\nu^{T}b)\right)$$

$$=\inf\left(\frac{1}{2}(\lambda^{T}+\nu^{T}A^{T}x-\nu^{T}b)\right)$$

$$=\frac{1}{2}(-\nu^{T}b^{T},\nu^{T}\lambda^{T}+\nu^{T}\lambda^{T}a^{T}a^{T}b)$$

#### 对偶问题举例

线性规划的对偶问题

最小化: $b^T \nu$ 

条件: $A^T \nu - \lambda + c = 0$ ,而且为 $\geq 0, \forall i = 1, \dots, n$ .

当优化问题的限制条件是线性条件时,可以利用共轭函数的一些 性质方便的得到对偶问题

### 线性约束优化问题的对偶问题

最小化: $f_0(x)$ 不等条件: $AX \le b$ 等式条件:Cx = d

这里面向量的比较 u < v 指的是 u 里边每一个分量都小于 v 里面对应的分量.

$$\inf\left(\mathcal{Y}^{T}x + f_{0}(x)\right) = - \sup\left(-\mathcal{Y}^{T}x - f_{0}(x)\right)$$

$$= - \int_{x}^{x} \left(-\mathcal{Y}^{T}x - f_{0}(x)\right)$$

$$= - \int_{x}^{$$

而这个函数的定义域就是

$$\operatorname{dom} g = \{(\lambda, \nu) : -A^T \lambda - C^T \nu \in \operatorname{dom} f_0^* \}.$$

$$|X| = |X|_{L^{1}} = \int_{X_{1}}^{2} \frac{1}{1 \cdot 1 \cdot 1} \frac{1}{1 \cdot 1$$

最小化向量范数

对偶问题举例\_\_\_\_\_

最小化:
$$|x|$$
  $A^T \mathcal{V} \leq A^T$ 

若 
$$f(x) = |x|$$
 那么共轭函数是  $f^*(y) = 0$ ,  $\forall |y| \le 1$ , 否则  $f^*(y) = +\infty$ . 所以对偶问题就是: ②  $|y| > 1$ ,  $x_t = t$  ②  $|y| > 1$ ,  $x_t = t$  ②  $|y| > 1$  ③  $|y| > 1$  ④  $|y| > 1$ 

### 最大熵问题

### 对偶问题举例

最小化:
$$f_0(x) = \sum_{i=1}^n x_i \ln x_i$$

不等条件: $Ax \leq b$ 

等式条件: $\mathbf{1}^T x = 1$ 

因为 
$$(f_1(x) + f_2(y))^* = f_1^*(w) + f_2^*(z)$$
, 我们有

$$f_0^*(y) = \sum_{i=1}^n e^{y_i - 1}$$

所以对偶函数为 
$$g(\lambda, \nu) = -b^T \lambda - \nu - e^{-\nu - 1} \sum_{i=1}^n e^{-a_i^T \lambda}$$

最大熵问题的对偶问题就是

### 最大熵问题的对偶问题

最大化: 
$$-b^T \lambda - \nu - e^{-\nu - 1} \sum_{i=1}^n e^{-a_i^T \lambda}$$

不等条件: $\lambda \ge 0$ 

### 对偶性

根据对偶函数的性质我们已经知道在对偶可行域中, $g(\lambda,\nu)$  总是不大于  $p^*$ . 所以就有

### 弱对偶性

$$d* \leq p*$$

若对偶性总是对的. 相对而言的强对偶性是指一部分优化问题来说, 有更好的结论.

#### 强对偶性

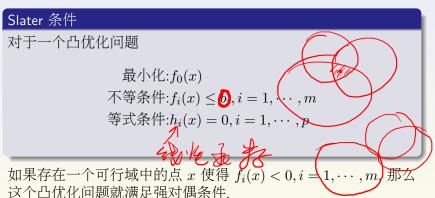
$$d* = p*$$

强对偶性不总成立.



### 强对偶性条件

第一个强对偶性的条件, 几乎所有的凸优化问题都满足强对偶性.



### 强对偶性条件

根据 Slater 条件,我们前述的几个例子都是满足强对偶性的.

### 满足强对偶性的例子

- 线性规划
- 最小二乘
- 最大熵问题

这种情况下我们如果发现对偶问题比原问题更容易解决,那么就可以使用对偶问题来解出  $d^* = p^*$ 

我们来看一下如果强对偶性满足的话,这些最优化点应该满足何种条件. 这一部分中我们假定所有的函数都是可微函数. 如果  $x^*$ ,  $(\lambda^*, \nu^*)$  分别是原问题与对偶问题的最优解,那么首先这些点应该满足可行域条件

- $f_i(x^*) \leq 0$
- $h_i(x^*) = 0$
- $\lambda_i^* \geq 0$

其次我们已经知道

$$d^* = g(\lambda^*, \nu^*)$$

$$\leq f_0(x^*) + \sum_{i=1}^m \lambda_i^* f_i(x^*) + \sum_{i=1}^p \nu_i^* h_i(x^*)$$

$$= f_0(x^*) + \sum_{i=1}^m \lambda_i^* f_i(x^*)$$

$$\leq f_0(x^*)$$

$$= p^*$$

于是  $d^* = p^*$  意味着上述不等式全都是等式.

所以我们有

$$\sum_{i=1}^{m} \lambda_i^* f_i(x^*) = 0$$

,以及

$$g(\lambda^*,\nu^*) = L(x^*,\lambda^*,\nu^*)$$

而因为

$$g(\lambda^*, \nu^*) = \inf(L(x^*, \lambda^*, \nu^*))$$

所以  $x^*$  是拉格朗日函数在 x 方向的驻点,所以有

$$\nabla_x L(x^*, \lambda^*, \nu^*) = 0$$

. 综上所述我们就的到了 KKT 条件.



#### KKT 条件

- $f_i(x^*) \leq 0, i = 1, \cdots, m$
- $h_i(x^*) = 0, i = 1, \cdots, p$
- $\lambda_i^* \geq 0$ ,  $i = 1, \dots, m$
- $\lambda_i^* f_i(x^*) = 0, i = 1, \dots, m$
- $\nabla_x L(x^*, \lambda^*, \nu^*) = 0$

其中第四个条件是由  $\sum_{i=1}^{m} \lambda_i^* f_i(x^*) = 0$  以及第一个和第三个条件 共同得到的.

所以我们有

$$\sum_{i=1}^{m} \lambda_i^* f_i(x^*) = 0$$

,以及

$$g(\lambda^*,\nu^*) = L(x^*,\lambda^*,\nu^*)$$

而因为

$$g(\lambda^*, \nu^*) = \inf(L(x^*, \lambda^*, \nu^*))$$

所以  $x^*$  是拉格朗日函数在 x 方向的驻点,所以有

$$\nabla_x L(x^*, \lambda^*, \nu^*) = 0$$

. 综上所述我们就的到了 KKT 条件.



#### KKT 条件叙述

- $f_i(x^*) \leq 0, i = 1, \cdots, m$
- $h_i(x^*) = 0, i = 1, \cdots, p$
- $\lambda_i^* \ge 0$ ,  $i = 1, \dots, m$
- $\lambda_i^* f_i(x^*) = 0, i = 1, \dots, m$
- $\nabla_x L(x^*, \lambda^*, \nu^*) = 0$

其中第四个条件是由  $\sum_{i=1}^{m} \lambda_i^* f_i(x^*) = 0$  以及第一个和第三个条件 共同得到的.

### KKT 条件使用

- 对于凸优化问题,KKT 条件是  $x^*$ ,  $(\lambda^*, \nu^*)$  分别作为原问题和对偶问题的最优解的充分必要条件.
- 对于非凸优化问题, KKT 条件仅仅是必要而非充分.

### 使用 KKT 条件解决优化问题例子

最小化:
$$(1/2)x^{P}x + q^{T}x + r$$
  
等式条件: $Ax = b$ 

KKT 条件是

$$Ax^* = b$$

以及

$$Px^* + q + A^Tv^* = 0$$

我们求解这个线性方程即可得到索要的结果.



### 应用举例: 支持向量机最简单形式



### 支持向量机的最简单形式

空间  $\mathbb{R}^n$  中有可分的两个点集 C,D. 我们希望找到一个最合适的 超平面  $a^Tx=b$  对他们进行区分. 也就是说对于  $p\in C, q\in D$ , 有

$$a^T p > b$$
  $a^T q < b$ 

也就是说存在一个正数 t 使得

$$2 a^T p - 2b \ge t, \quad a^T q - 2b \le 2t$$

这个正数 t 一定程度上描述了两个集合被分开的程度.

为了体现"最合适",一个比较好的想法就是希望当 a 是单位向量的时候,t 越大越好. 于是的到了一个优化问题

#### 应用举例: 支持向量机最简单形式

$$f_i \in O$$



### SVM 优化问题



最小化:

 $-a^T p_i + b \le -t$ 不等条件 1:

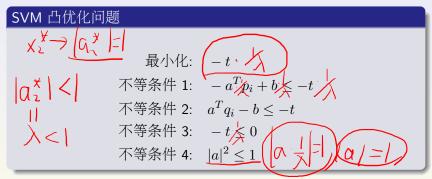
不等条件 2:  $a^T q_i - b < -t$ 

不等条件 3: -t < 0

但是这里面  $|a|^2 = 1$  并非凸条件,所以我们把它换成  $|a|^2 \le 1$ .

就得到一个跟原问题等价的凸优化问题.

### 应用举例: 支持向量机最简单形式



具体求解,因为这是一个凸优化问题,所以可以选择使用 KKT 条件来解方程求解.

# 总结寄语 (数学部分)

#### 总结

- 本次课程数学部分涵盖微积分,线性代数,概率统计,凸优化的内容
- 复习目标:
  - 分别对每一部分进行总结,得到清晰知识结构脉络
  - 最好对每一部分都能记住一些有代表性的例子,比如矩阵里边的旋转,拉伸变换
- 使用目标: 在工作和学习中遇到相关知识的时候
  - 可以回忆起最基本的想法(比如微分的本质是逼近,矩阵变换本质是换基)
  - 对于技术细节能很快找到参考资料

# 总结寄语 (数学部分)

#### 小故事

R.Thom 是法国人, 35 岁得的 Fields 奖。在一次采访当中, 作为数学家的 Thom 同两位古人类学家讨论问题。谈到远古的人们为什么要保存火种时

- 一个人类学家说, 因为保存火种可以取暖御寒;
- 另外一个人类学家说, 因为保存火种可以烧出鲜美的肉食。
- 而 Thom 说,因为夜幕来临之际,火光摇曳妩媚,灿烂多姿,是最美最美的。

除了上面所有目标,我还希望本次课程中讲述的数学各个分支之间产生的一些相互作用能或多或少给大家一点数学独特美学的享受. 比如微积分与线性代数之间,矩阵与几何之间,凸集合与凸函数之间那些意料之外又在情理之中的关系. 毕竟数学的终极目的还是美.

$$f(v,v) = f(v) + f(v)$$

$$f^*(w,z) = \sup_{x \in \mathbb{R}} \{ w^{\tau}u + z^{\tau}v - f(w - f(v)) \}$$

$$= \sup_{x \in \mathbb{R}} \{ w^{\tau}u + f(u) \} + \sup_{x \in \mathbb{R}} \{ z^{\tau}u + f(u) \}$$

$$= f^*(w) + f^*(z)$$