

# 基本优化方法（现学现卖）

李显求

VIPL Group

April 10, 2016

# 提纲

- 主要内容
  - 凸问题
  - 泰勒展开式
  - 梯度下降法
  - 牛顿迭代法
  - Levenberg-Marquardt 算法
  - 拟牛顿法
    - BFGS 算法
    - L-BFGS 算法
  - 次梯度算法
  - Bregman Projection Problem 及解法
  - Lagrange 对偶问题
  - 几个例子
- Q&A

## 符号编辑页 1 (相当于草稿)

- 凸集合  $\mathbf{x}_1, \mathbf{x}_2 \in S$ , then  $\mathbf{x} = \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in S$
- 凸函数  $f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2) = f(\mathbf{x}_2 + \lambda(\mathbf{x}_1 - \mathbf{x}_2)) < \lambda f(\mathbf{x}_1) + (1 - \lambda)f(\mathbf{x}_2), \lambda \in (0, 1)$
- Bregman Divergence:  $D_F(\mathbf{p}, \mathbf{q}) = F(\mathbf{p}) - F(\mathbf{q}) - \langle \nabla F(\mathbf{q}), \mathbf{p} - \mathbf{q} \rangle$
- $f$  是一个线性函数和一个常数的和, 也就是  $f(\mathbf{x}) = A\mathbf{x} + \mathbf{b}$  的形式, 其中  $A \in R^{m \times n}, \mathbf{b} \in R^m$ , 那么  $f: R^n \rightarrow R^m$  是仿射的
- 逐点上确界  $g(\mathbf{x}) = \sup\{f_1(\mathbf{x}), f_2(\mathbf{x}), \dots\}$
- $\mathcal{L}(x, \lambda) = f_1(x) + \lambda f_0(x)$
- $s(x) = \frac{\exp(x)}{1 + \exp(x)} = \frac{1}{1 + \exp(-x)}, s'(x) = s(x)(1 - s(x))$
- $\min |x|, t_i = \frac{1}{i+1}$
- $\text{dist}_{\mathbf{x}' \in S}^2(\mathbf{x}', \mathbf{x}), S$

$$\begin{cases} \min f_1(x) \\ f_0(x) < 0 \end{cases} \quad (1)$$

## 符号编辑页 2 (相当于草稿)

- $D_{ld}(A, A_0) = \text{tr}(AA_0^{-1}) - \log \det(AA_0^{-1}) - n$
- $F(X) = -\log |\det(X)|, \nabla_X F(X) = -(X^T)^{-1}$

$$\begin{aligned} & \min_{A \succeq 0} D_{ld}(A, A_0) + \gamma D_{ld}(\text{diag}(\xi), \text{diag}(\xi_0)) \\ s.t. \quad & \text{tr}(A(x_i - x_j)(x_i - x_j)^T) \leq \xi_{ij}, (i, j) \in S \\ & \text{tr}(A(x_i - x_j)(x_i - x_j)^T) \geq \xi_{ij}, (i, j) \in D \end{aligned} \quad (2)$$

$$\begin{aligned} & \min_{A \succeq 0} D_{ld}(A_{k+1}, A_k) + \gamma D_{ld}(\text{diag}(\xi_{k+1}), \text{diag}(\xi_k)) \\ s.t. \quad & \delta_{ij} \text{tr}(A_{k+1}(x_i - x_j)(x_i - x_j)^T) \leq \delta_{ij} \xi_{ij} \\ & \delta_{ij} = \begin{cases} 1, & \text{if } (i, j) \in S \\ -1, & \text{if } (i, j) \in D \end{cases} \end{aligned} \quad (3)$$

## 符号编辑页 3 (相当于草稿)

- $D_{ld}(A, A_0) = \text{tr}(AA_0^{-1}) - \log \det(AA_0^{-1}) - n$
- $F(X) = -\log |\det(X)|, \nabla_X F(X) = -(X^T)^{-1}$

$$\begin{aligned}\mathcal{L}(A, \beta) &= D_{ld}(A, A_k) + \gamma D_{ld}(\text{diag}(\xi_{k+1}), \text{diag}(\xi_k)) \\ &\quad + \alpha(\delta_{ij} \text{tr}(A(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T) - \delta_{ij} \xi_{ij}) \\ \nabla_A \mathcal{L}(A, \beta) &= (A_k^{-1})^T - (A^{-1})^T + \alpha \delta_{ij} C_{ij}^T \\ A &= (A_k^{-1} + \alpha \delta_{ij} C_{ij})^{-1} \\ &= A_k - \alpha \delta_{ij} \frac{A_k(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T A_k}{1 + \alpha \delta_{ij}(\mathbf{x}_i - \mathbf{x}_j)^T A_k(\mathbf{x}_i - \mathbf{x}_j)}\end{aligned}\tag{4}$$

$$\begin{aligned}D \log \det(AB)(A)[H] &= D \log \det(AB)(AB)[D(AB)(A)[H]] \\ &= D \log \det(AB)(AB)[HB] \\ &= \langle \nabla_{AB} \log \det(AB), (HB)^T \rangle \\ &= \langle ((AB)^{-1})^T, (HB)^T \rangle \\ &= \text{tr}((A^{-1})^T (B^{-1})^T B^T H^T)\end{aligned}\tag{5}$$

## 符号编辑页 3（相当于草稿）



$$\begin{aligned}\nabla_{A_{k+1}} D_{ld}(A_{k+1}, A_0) &= \nabla_{A_k} D_{ld}(A_k, A_0) + \alpha \delta_{ij} C_{ij} \\ (A_0^{-1})^T - (A_{k+1}^{-1})^T &= (A_0^{-1})^T - (A_k^{-1})^T + \alpha \delta_{ij} C_{ij} \\ A_{k+1} &= (A_k^{-1} - \alpha \delta_{ij} C_{ij})^{-1} \\ &= A_k + \alpha \delta_{ij} \frac{A_k(\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^T A_k}{1 - \alpha \delta_{ij} (\mathbf{x}_i - \mathbf{x}_j) A_k (\mathbf{x}_i - \mathbf{x}_j)}\end{aligned}\tag{6}$$

# 从泰勒展开说起

## 定义

设  $f(x)$  是实数域上的无穷可微的函数, 那么它在任意一点  $x_0$  泰勒展开式表示为<sup>a</sup>:

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \cdots \quad (7)$$

同样的对于无穷可微向量函数  $f(\mathbf{x})$  以及初始点  $\mathbf{x}_0$  有<sup>b</sup>:

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \left( \frac{\partial f(\mathbf{x}_0)}{\partial \mathbf{x}} \right)^T (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2!} (\mathbf{x} - \mathbf{x}_0)^T \mathbf{H} (\mathbf{x} - \mathbf{x}_0) + \cdots \quad (8)$$

<sup>a</sup>这里只讨论优化问题故省略了关于收敛域的讨论; 所以只需要知道它是  $x_0$  周围的一个小领域就可以了, 优化时只要步长不太大一般不会违反这个准则

<sup>b</sup>其中  $\mathbf{H}$  矩阵就是通常意义的 Hessian 矩阵

# 梯度下降优化方法

首先，梯度下降涉及两个参数：方向  $\mathbf{g}$  和步长  $\alpha$ ，为了方便起见这里不妨假设  $\|\mathbf{g}\| = 1$ （虽然在实际中我们并不一定做归一化）。最优化的目的就是使得  $f(\mathbf{x}_0 + \alpha\mathbf{g})$  最小，于是利用一阶泰勒展开我们得到：

$$f(\mathbf{x}_0 + \alpha\mathbf{g}) \approx f(\mathbf{x}_0) + \alpha \left( \frac{\partial f(\mathbf{x}_0)}{\partial \mathbf{x}} \right)^T \mathbf{g} \quad (9)$$

上式右边当  $\mathbf{g} = - \left( \frac{\partial f(\mathbf{x}_0)}{\partial \mathbf{x}} \right) / \left\| \frac{\partial f(\mathbf{x}_0)}{\partial \mathbf{x}} \right\|$  的时候是  $f(\mathbf{x}_0 + \alpha\mathbf{g})$  较  $f(\mathbf{x}_0)$  下降最快的方向，所以梯度下法也叫最速下降法。

至于步长  $\alpha$  的选择则是一个一维的优化问题，这个问题比较简单（二分法，0.618 法等都可以解决），当然也可以是预先定义的。

对于凸函数由于有： $f(\mathbf{x}) \geq f(\mathbf{x}_0) + \left( \frac{\partial f(\mathbf{x}_0)}{\partial \mathbf{x}} \right)^T (\mathbf{x} - \mathbf{x}_0)$ ，所以梯度下降算法可以保证目标函数值是不增加的



# 牛顿迭代法 (1/3)

下面的动画显示了牛顿法在一维时的情况，这也相当于阐述了牛顿法的几何意义

在下一页中将使用数学来解释这种过程了，也就是牛顿法（牛顿迭代）的数学形式。

## 牛顿迭代法 (2/3)

类似于梯度下降算法，首先给出  $f(x)$  的二阶泰勒展开式<sup>1</sup>

$$\begin{aligned} f(x) &= f(x_0) + \left( \frac{\partial f(x_0)}{\partial x} \right)^T (x - x_0) + \frac{1}{2!} (x - x_0)^T H (x - x_0) + \cdots \\ \frac{\partial f(x)}{\partial x} &= \frac{\partial f(x_0)}{\partial x} + H(x - x_0) + \cdots \end{aligned} \quad (10)$$

不难发现的是过  $x_0$  点  $\frac{\partial f(x)}{\partial x}$  的“切平面”的就是：

$$y = \frac{\partial f(x_0)}{\partial x} + H(x - x_0) \quad (11)$$

---

<sup>1</sup>需要注意的是：牛顿法计算的是零点（或者说方程的根，因此，其实牛顿法是通过极值点处的梯度为零将两者联系起来

## 牛顿迭代法 (3/3)

在公式  $y = \frac{\partial f(x_0)}{\partial x} + H(x - x_0)$  中, 令  $y = 0$  可以得到:

$$x = x_0 - H^{-1} \frac{\partial f(x_0)}{\partial x}$$

or

(12)

$$x_{k+1} = x_k - H^{-1} \frac{\partial f(x_k)}{\partial x}$$

在最后, 这里在举一个例子用来说明牛顿法的应用 (求一个实数  $a$  的  $n$  次方根):

- 首先写出方程:  $f(x) = x^n - a$ , 则  $f(x)$  的零点即为所求。
- 根据牛顿迭代法得到更新公式:  $x_{k+1} = x_k - \frac{x_k^n - a}{nx_k^{n-1}}$ 。
- 最后由于算法是二阶的所以很快就会收敛了。

# Levenberg-Marquardt 算法

假设  $f$  是一个从  $\mathbb{R}^m \rightarrow \mathbb{R}^n$  的非线性映射，也就是说  $\mathbf{x} \in \mathbb{R}^m$ ，经过  $f$  变换  $\mathbf{y} = f(\mathbf{w}|\mathbf{x})$  有  $\mathbf{y} \in \mathbb{R}^n$ 。

Levenberg-Marquardt 算法的目的就是希望任意给定一个  $\mathbf{y}$  以及合理的初始值  $\mathbf{w}_0$ ，能找到一个  $\mathbf{w}^*$ ，使得  $\epsilon^T \epsilon$  尽量小（局部极小），其中

$$\epsilon = f(\mathbf{w}^*|\mathbf{x}) - \mathbf{y}$$

接着还是请出  $f(\mathbf{w}|\mathbf{x})$  的泰勒展开（一阶）： $f(\mathbf{w} + \mathbf{g}) \approx f(\mathbf{w}|\mathbf{x}) + \mathbf{J}\mathbf{g}$  对于第  $k+1$  次迭代：

$$\begin{aligned}\min \|\epsilon_{k+1}\|^2 &= \min \|\mathbf{y} - f(\mathbf{w} + \mathbf{g}_k)\| \\ &\approx \min \|\mathbf{y} - f(\mathbf{w}|\mathbf{x}) - \mathbf{J}\mathbf{g}_k\|^2 \\ &= \min \|\epsilon_k - \mathbf{J}\mathbf{g}_k\|\end{aligned}\tag{13}$$

该最小二乘问题有解析解： $(\mathbf{J}^T \mathbf{J})\mathbf{g}_k = \mathbf{J}^T \epsilon_k$

对该式稍作修改可以得到（类似于岭回归）： $(\mu \mathbf{I} + \mathbf{J}^T \mathbf{J})\mathbf{g}_k = \mathbf{J}^T \epsilon_k$

最后再说一句：当  $\mu$  较大的时候上式接近梯度下降，反之接近牛顿法

## BFGS 算法 (1/2)

BFGS(**B**royden,**F**letcher,**G**oldfarb,**S**hanno); 本节内容参考:

<http://blog.csdn.net/itplus/article/details/21896453>

牛顿法是二阶收敛的, 速度上是非常理想的; 但是 Hessian 矩阵非常难以计算, 在实际计算中往往得不偿失甚至计算不出来; 拟牛顿法应运而生, 其核心就是如何估计 Hessian 矩阵 (或其逆)。

另一方面为了判断估计的好坏, 于是有了拟牛顿条件, 首先回顾牛顿法中的公式:

$$\frac{\partial f(x)}{\partial x} \approx \frac{\partial f(x_0)}{\partial x} + H(x - x_0)$$

令  $x = x_{k+1}$ ,  $x_0 = x_k$  并令  $g_{k+1} = \frac{\partial f(x_{k+1})}{\partial x}$ ,  $g_k = \frac{\partial f(x_k)}{\partial x}$  带入上式得到拟牛顿条件:

$$g_{k+1} - g_k \approx H_k(x_{k+1} - x_k)$$

最后  $s_k \triangleq x_{k+1} - x_k$ ,  $y_k \triangleq g_{k+1} - g_k$  带入得到:

$$s_k \approx H^{-1}y_k \text{ or } y_k \approx Hs_k$$

## BFGS 算法 (2/2)

与 DFP<sup>2</sup>不同 BFGS 算法估计的是 Hessian 矩阵而前者估计的是 Hessian 矩阵的逆

BFGS: 利用迭代式  $B_{k+1} = B_k + \Delta B_k$  估计  $H_k: B_{k+1} \approx H_k$ , 其主要有以下几步:

- 假设  $\Delta B_k = \alpha u u^T + \beta v v^T$ , 其中  $\alpha, \beta, u, v$  待定
- 带入拟牛顿条件:  
件:  $y_k \approx B_k s_k + (\alpha u u^T + \beta v v^T) s_k = B_k s_k + (\alpha u^T s_k) u + (\beta v^T s_k) v$ .
- 令  $(\alpha u^T s_k) = 1, (\beta v^T s_k) = -1, u = y_k, v = B_k s_k$ .
- 解得  $\alpha = \frac{1}{y_k^T s_k}, \beta = \frac{1}{s_k^T B_k s_k}$ , 于是  $\Delta B_k = \frac{y_k y_k^T}{y_k^T s_k} - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k}$ .
- 更新公式  $B_{k+1} = B_k + \frac{y_k y_k^T}{y_k^T s_k} - \frac{B_k s_k s_k^T B_k}{s_k^T B_k s_k}$ .

最后利用 Sherman-Morrison 公式<sup>3</sup>可以方便的计算  $B_{k+1}^{-1}$ :

$$B_{k+1}^{-1} = \left( I - \frac{s_k y_k^T}{y_k^T s_k} \right) B_k^{-1} \left( I - \frac{y_k s_k^T}{y_k^T s_k} \right) + \frac{s_k s_k^T}{y_k^T s_k}$$

<sup>2</sup><http://blog.csdn.net/itplus/article/details/21896981>

# Sherman-Morrison 公式

## 定义

设  $A \in \mathbb{R}^{n \times n}$  为非奇异方阵,  $u, v \in \mathbb{R}^n$ , 若  $1 + v^T A^{-1} u \neq 0$ , 则有:

$$(A + uv^T)^{-1} = A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u}$$

借此空间, show 一下 BFGS 算法的作者:

Broyden, Fletcher, Goldfarb, Shanno



## L-BFGS 算法 (1/3)

L-BFGS 是 Limit-Storage/Limit-Memory BFGS 的缩写，其设计的目的就是为了解决 BFGS 的内存问题（因为  $B_k$  在数据量较大时候其存储量是非常巨大的），解决思路就是从下式出发：

$$B_{k+1}^{-1} = \left( I - \frac{s_k y_k^T}{y_k^T s_k} \right) B_k^{-1} \left( I - \frac{y_k s_k^T}{y_k^T s_k} \right) + \frac{s_k s_k^T}{y_k^T s_k}$$

- 记  $\rho_k = \frac{1}{y_k^T s_k}$ ,  $V_k = I - \rho_k y_k s_k$ ;
- 于是有:  $B_{k+1}^{-1} = V_k^T B_k V_k + \rho_k s_k s_k^T$ ; 递推得到:

$$B_{k+1}^{-1} = \left( \Pi_{i=0}^k V_i^T \right) B_0 \left( \Pi_{i=0}^k V_i \right) + \sum_{j=1}^k \left( \left( \Pi_{i=j}^k V_i^T \right) (\rho_{j-1} s_{j-1} s_{j-1}^T) \left( \Pi_{i=j}^k V_i \right) \right) + \rho_k s_k s_k^T$$



## L-BFGS 算法 (2/3)

接上一页内容:

- 计算  $\mathbf{B}_{k+1}$  需要  $\{\mathbf{s}_i, \mathbf{y}_i\}_{i=0}^k$ , 这很耗空间, 当  $k+1 > m$  考虑  $m$  阶截断:

$$\begin{aligned}\mathbf{B}_{k+1}^{-1} = & \left( \Pi_{i=k-m+1}^k \mathbf{V}_i^T \right) \mathbf{B}_0 \left( \Pi_{i=k-m+1}^k \mathbf{V}_i \right) \\ & + \sum_{j=k-m+2}^k \left( \left( \Pi_{i=j}^k \mathbf{V}_i^T \right) (\rho_{j-1} \mathbf{s}_{j-1} \mathbf{s}_{j-1}^T) \left( \Pi_{i=j}^k \mathbf{V}_i \right) \right) + \rho_k \mathbf{s}_k \mathbf{s}_k^T\end{aligned}$$

- 通常  $m < 10$ ; 至此, L-BFGS 算法的理论部分就基本完成了, 下一节介绍其算法流程。

# L-BFGS 算法 (3/3): 算法流程<sup>4</sup>

算法 ( $B_k \cdot g_k$  的快速算法)

Step 1 初始化.

$$\delta = \begin{cases} 0, & \text{若 } k \leq m \\ k - m, & \text{若 } k > m \end{cases}; \quad L = \begin{cases} k, & \text{若 } k \leq m \\ m, & \text{若 } k > m \end{cases}; \quad q_L = g_k.$$

Step 2 后向循环.

FOR  $i = L - 1, L - 2, \dots, 1, 0$  DO

{

$$j = i + \delta;$$

$$\alpha_i = \rho_j s_j^T q_{i+1}; \quad // \alpha_i \text{ 需要存下来, 前向循环要用!}$$

$$q_i = q_{i+1} - \alpha_i y_j.$$

}

Step 3 前向循环.

$$r_0 = B_0 \cdot q_0;$$

FOR  $i = 0, 1, \dots, L - 2, L - 1$  DO

{

$$j = i + \delta;$$

$$\beta_j = \rho_j y_j^T r_i;$$

$$r_{i+1} = r_i + (\alpha_i - \beta_j) s_j.$$

}

最后算出的  $r_L$  即为  $H_k \cdot g_k$  的值.

<sup>4</sup> Jorge Nocedal. updating quasi-newton matrices with limited storage. Mathematics of Computation. 1980.

## 次梯度算法 (1/3)

关于次梯度这里: <http://www.hanlongfei.com/convex/2015/10/02/cmu-10725-subgradidient/>做了比较详细的介绍(这里开始就完全是现学现卖了); 首先, 来回顾一下凸函数的一阶判别条件<sup>5</sup>

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x})$$

两个条件: 1)  $f(x)$  可微, 2) 上式在**整个定义域**都满足  
点  $x$  的**次梯度**指在函数  $f$  上的点  $x$  满足以下条件的  $\mathbf{g} \in \mathbb{R}^n$ :

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x})$$

几点注意:

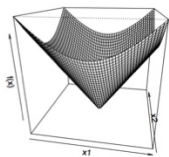
- $f(x)$  可以不是凸函数也不要求它是可微的
- 该定义是在一个点 ( $x$ ) 上定义的<sup>6</sup>
- 次梯度可以有多个(一个集合), 可以只有有一个, 也可以是空集

<sup>5</sup>顺带说一下它的二阶判别条件就是, 其 Hessian 阵是半正定的(如果存在的话)

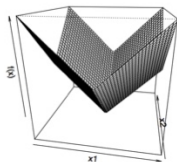
<sup>6</sup>个人认为应该还有一定的范围限制, 如在  $x$  的领域内

## 次梯度算法 (2/3)

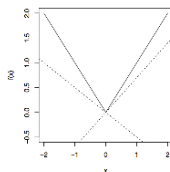
几个需要计算次梯度的例子（各个图的函数形式位于图篇的底部）：



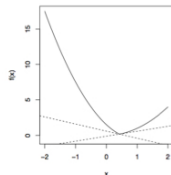
$\|x\|_2$



$\|x\|_1$



$|x|$



$\max\{f_1(x), f_2(x)\}$

有了次梯度接下来定义次微分（记为  $\partial f(x)$ ）的概念：

$$\partial f(x) = \{g \in \mathbb{R}^n : g \text{ is a subgradient of } f \text{ at } x\}$$

在第  $k$  次迭代时有了次梯度  $g_k \in \partial f(x)$  就可以类似于梯度下降的算法定于次梯度下降算法（其中  $t_k$  表示步长）：

$$x_{k+1} = x_k - t_k g_k, k = 1, 2, 3, \dots$$

## 次梯度算法 (3/3)

关于次梯度算法，有几点需要注意：

- 如果目标函数是凸的那么它能保证算法是收敛的，但是次梯度算法并不能保证目标函数每次都是递减的（或者说是单调下降的），为此会有

$$f(\mathbf{x}_{best}^{(k)}) = \min_{i=0,1,\dots,k} \{f(\mathbf{x}_i)\} = \min\{f(\mathbf{x}_k), f(\mathbf{x})_{k-1}\}$$

- 当次梯度方向不唯一的时候理论上是选取任意一个都可以的
- 与梯度下降不同次梯度下降的步长是预先设定好的且满足：

$$\sum_{i=0}^{\infty} t_i = \infty \text{ and } \sum_{i=0}^{\infty} t_i^2 < \infty \quad (14)$$

关于公式14最后再说两句：1）公式前半部分保证了搜索范围；2）后半部分保证了步长减小（收敛）

# Bregman Projection<sup>7</sup>(1/2)

首先，来回顾一下 Bregman Divergence:

$$D_F(p, q) = F(p) - F(q) - \langle p - q, \nabla F(q) \rangle$$

其中  $F(x)$  是实值可微的凸函数:  $F(x) : \Omega \rightarrow R$ , 而所谓 Projection 就是指的是将原始空间中的数据 (例如下面的  $W_1$ ) 投影到某个子空间 (集合) 的过程。

其物理意义就是原数据到这个子空间 (集合) 的距离最短; 所以 Bregman Projection 描述的是类似于如下形式的问题:

$$\begin{aligned} \mathbf{W}^* &= \min_{\mathbf{W}} D_F(\mathbf{W}, \mathbf{W}_1) \\ \text{s.t. } \mathbf{W} &= \mathbf{W}^T, \text{tr}(\mathbf{W}) = 1 \\ \text{tr}(\mathbf{W}\mathbf{C}_j) &< 0, \text{for } j = 1, 2, \dots, n \end{aligned}$$

---

<sup>7</sup>Tsuda, K., Raetsch, G., & Warmuth, M. K. (2005). Matrix exponentiated gradient updates of online learning and bregman projection. Journal of Machine Learning Research, 6.

## Bregman Projection<sup>8</sup>(2/2)

然而全部约束一起考虑太难了，因此 Bregman 说了：

*It is well known that the Bregman projection into the intersection of convex regions can be solved by sequential projections to each region (Bregman, 1967; Censor and Lent, 1981)*

于是前述问题可以改成每次只考虑一个约束条件的形式（通常是不满足的那个），循环迭代：

$$\begin{aligned} \mathbf{W}^* &= \min_{\mathbf{W}} D_F(\mathbf{W}, \mathbf{W}_1) \\ s.t \quad &\mathbf{W} = \mathbf{W}^T, \text{tr}(\mathbf{W}) = 1 \\ &\text{tr}(\mathbf{W}\mathbf{C}_j) \end{aligned}$$

---

<sup>8</sup>Tsuda, K., Raetsch, G., & Warmuth, M. K. (2005). Matrix exponentiated gradient updates of online learning and bregman projection. Journal of Machine Learning Research, 6.

## Lagrange 对偶问题 (1/3)

前面介绍的内容都是目标函数已有的优化方法，这一部分将会介绍如何将一个一般的有约束的问题转化为 Lagrange 对偶问题。

原问题：

$$\begin{aligned} \min f_0(\mathbf{x}) \\ \text{s.t. } f_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m \\ h_i(\mathbf{x}) = 0, i = 1, 2, \dots, p \end{aligned} \quad (15)$$

Lagrange 函数：

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v}) = f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p v_i h_i(\mathbf{x}) \quad (16)$$

Lagrange 对偶函数：

$$g(\boldsymbol{\lambda}, \mathbf{v}) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{v}) = \inf_{\mathbf{x}} \left( f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i f_i(\mathbf{x}) + \sum_{i=1}^p v_i h_i(\mathbf{x}) \right) \quad (17)$$



## Lagrange 对偶问题 (2/3)

对偶问题17对任意的  $\lambda \geq 0$  以及  $v$  都是原问题15的下界，此外对于  $\lambda < 0$  的情形这将导致  $g(\lambda, v) = -\infty$  并没有实际意义。

既然  $g(\lambda, v)$  对于任意的  $\lambda \geq 0$  以及  $v$  是原问题15的下界，那么什么样的  $\lambda, v$  才是好的，对偶问题考虑：

$$\max g(\lambda, v), s.t \lambda \geq 0 \quad (18)$$

所以实际上的 Lagrange 对偶问题优化的是：

$$\max_{\lambda, v} \left( \min_x \mathcal{L}(x, \lambda, v) \right)$$

关于 SVM 的博客：

[http://blog.csdn.net/v\\_july\\_v/article/details/7624837](http://blog.csdn.net/v_july_v/article/details/7624837)

中把 SVM 作为一个实例对 Lagrange 对偶问题进行了解说

## Lagrange 对偶问题 (3/3)

当原问题是凸问题的时候, 满足 **KKT** 条件<sup>9</sup>的点也是对偶问题的解。  
**KKT** 条件:

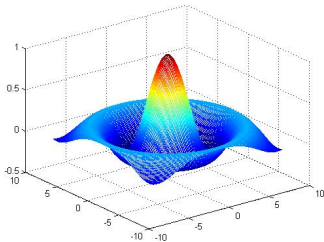
$$\left\{ \begin{array}{l} \nabla f_0(\mathbf{x}) + \sum_{i=1}^m \lambda_i \nabla f_i(\mathbf{x}) + \sum_{i=1}^p v_i \nabla h_i(\mathbf{x}) = \mathbf{0} \\ f_i(\mathbf{x}) \leq 0, i = 1, 2, \dots, m \\ h_i(\mathbf{x}) = 0, i = 1, 2, \dots, p \\ \lambda_i \geq 0, i = 1, 2, \dots, m \\ \lambda_i f_i(\mathbf{x}) = 0, i = 1, 2, \dots, m \end{array} \right.$$

<sup>9</sup>KKT: Karush-Kuhn-Tucker



## 几个例子——之一 (1/4)

首先，第一个例子来介绍一下 Mexico 草帽：



函数的形式：
$$f(x, y) = \frac{\sin(\sqrt{x^2 + y^2})}{\sqrt{x^2 + y^2}}$$



## 几个例子——之一 (2/4)

接上一页:

- 函数的梯度形式:

$$\partial f(x, y) = \begin{cases} \cos(\sqrt{x^2 + y^2}) \frac{x}{x^2 + y^2} - \sin(\sqrt{x^2 + y^2}) \frac{x}{\sqrt{x^2 + y^2}^3} \\ \cos(\sqrt{x^2 + y^2}) \frac{y}{x^2 + y^2} - \sin(\sqrt{x^2 + y^2}) \frac{y}{\sqrt{x^2 + y^2}^3} \end{cases}$$

- Hessian 矩阵:

$$H = \begin{bmatrix} a_{11} & a_{12} \\ a_{12} & a_{22} \end{bmatrix} \text{ where } \begin{cases} a_{11} = -\sin(\sqrt{x^2 + y^2}) \frac{x^2(x^2 + y^2) + y^2 - 2x^2}{\sqrt{x^2 + y^2}^5} \\ \quad + \cos(\sqrt{x^2 + y^2}) \frac{(y^2 - 2x^2)}{\sqrt{x^2 + y^2}^4} \\ a_{12} = -\sin(\sqrt{x^2 + y^2}) \frac{xy(x^2 + y^2) - 3xy}{\sqrt{x^2 + y^2}^5} \\ \quad - \cos(\sqrt{x^2 + y^2}) \frac{3xy}{\sqrt{x^2 + y^2}^4} \\ a_{22} = -\sin(\sqrt{x^2 + y^2}) \frac{y^2(x^2 + y^2) + x^2 - 2y^2}{\sqrt{x^2 + y^2}^5} \\ \quad + \cos(\sqrt{x^2 + y^2}) \frac{(x^2 - 2y^2)}{\sqrt{x^2 + y^2}^4} \end{cases}$$



## 几个例子——之一 (3/4)

首先注意到的是：即便是如此简单的一个函数函数其 Hessian 矩阵也是难以计算的。

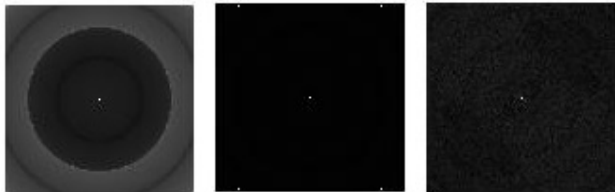
下面就来看一下墨西哥草帽上的优化问题的一些结果：

Mexico_hat(range=[-10,10,-10,10],step=0.2,scale=101×101)					
算法	总时间（单位：s）	maxiter	average_iter	存储	备注
grad	8.450431	500	93.97	$O(n)$	###
newton	3.720508	500	5.43	$O(n^2)$	###
bfgs	14.776714	500	42.27	$O(n^2)$	###

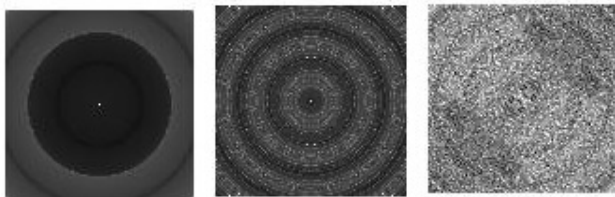


## 几个例子——之一 (4/4)

再来看几张图（从左到右：grad->newton->bfgs）：



把图片稍微锐化一下就是下面的样子了：



这些图片代表了在  $[-10,10,-10,10]$  的范围内以 0.2 为步长遍历所有点，并把它作为初始值赋给算法，然后各个点的灰度值代表了迭代次数大小。



## 几个例子——之二 (1/2)

这次做的是回归问题： $f(x) = p_1 e^{p_2 x}$ ，采样  $x = -1 : 0.02 : 1$ ，获得观测值  $y_n = f(x)$ ，其中真实参数  $(p_1, p_2) = (0.5, 0.1)$ ，目标函数：

$$\min \|\epsilon \epsilon^T\|, \epsilon = \hat{y} - y_n$$

这次直接上结果：

fit(range=[-1,1,-1,1],step=0.02,scale=101×101)					
算法	总时间（单位：s）	maxiter	average_iter	存储	备注
grad	90.776934	178	156.14	$O(ns)$	###
levmar	160.395254	384	351.70	$O(ns)$	$\mu = 1$
levmar	91.96887	217	200.64	$O(ns)$	$\mu = 0.1$
levmar	84.676262	200	185.18	$O(ns)$	$\mu = 0.01$
levmar	83.078666	198	182.28	$O(ns)$	$\mu = 0$
fit(range=[-10,10,-10,10],step=0.02,scale=1001×1001)					
grad	14965.815	500	332.04	$O(ns)$	###
levmar	31156.32704	500	486.58	$O(ns)$	$\mu = 0$

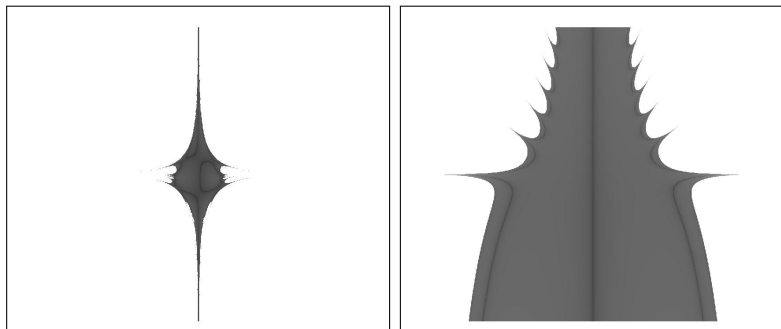


## 几个例子——之二 (2/2)

最后还是看几个图（从左到右：grad, levmar( $\mu = 0.01$ ), levmar( $\mu = 1$ )):



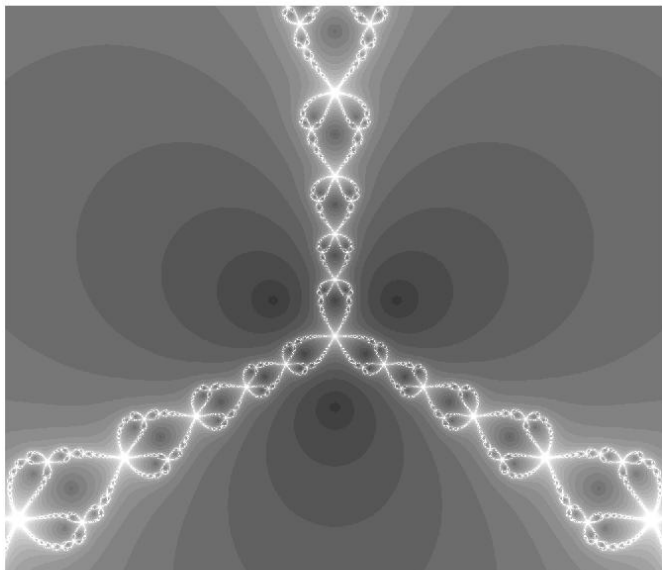
更大范围内的情况 (range=[-10,10,-10,10], 左到右: grad, levmar( $\mu = 0$ )):





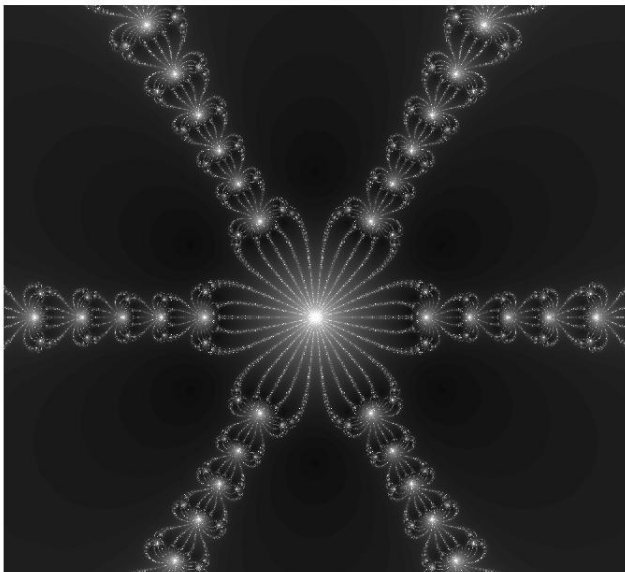


## 几个例子——之 Newton 法与多项式 ( $x^3 - 0.01$ )



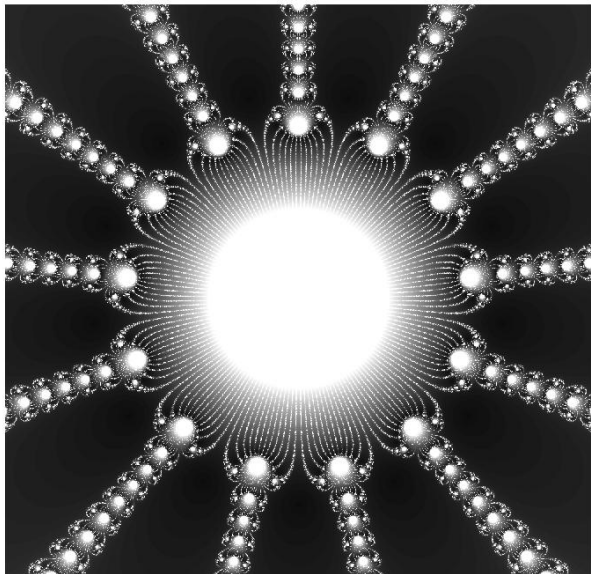


## 几个例子——之 Newton 法与多项式 ( $x^6 - 0.01$ )





## 几个例子——之 Newton 法与多项式 ( $x^{13} - 0.01$ )



# 下集预告

- 本集遗漏的内容
  - 共轭梯度算法
  - EM 算法
- 一维搜索（线搜），以下内容全凭记忆想到，欢迎指正和补充
  - 二分查找
  - 0.618 法
  - BFGS 一维的情况
  - 其它（暂时没想好）

谢谢！