

最优化第三次作业

162050127 颜劭铭

2022 年 5 月 6 日

1 基础题

1.1 凸函数

1.1.1 证明 $f(x) = \log \sum_{k=1}^n \exp(x_k)$ 是凸函数

由题目可知: $f(x) = \log \sum_{k=1}^n \exp(x_k) = \log(e^{x_1} + e^{x_2} + \cdots + e^{x_n}), x \in \mathbf{R}^n$

因为对数的底数取多少并不影响函数 $f(X)$ 的凸性, 所以为了方便求导, 我们取底数为 e 。

因此对 $f(x)$ 求一阶偏导:

$$\frac{\partial f}{\partial x_i} = \frac{e^{x_i}}{e^{x_1} + e^{x_2} + \cdots + e^{x_n}}$$

接着对 $f(x)$ 求二阶偏导:

当 $i \neq j$ 的时候:

$$\frac{\partial^2 f}{\partial x_i \partial x_j} = \frac{-e^{x_i} e^{x_j}}{(e^{x_1} + e^{x_2} + \cdots + e^{x_n})^2}$$

当 $i = j$ 的时候:

$$\frac{\partial^2 f}{\partial x_i \partial x_i} = \frac{e^{x_i}(e^{x_1} + e^{x_2} + \cdots + e^{x_n}) - e^{x_i} e^{x_i}}{(e^{x_1} + e^{x_2} + \cdots + e^{x_n})^2}$$

此时定义一个列向量: $z = [e^{x_1}, e^{x_2}, \cdots, e^{x_n}]^T$

因此, Hessian 矩阵为:

$$H = \begin{vmatrix} \frac{e^{x_1}(e^{x_1}+e^{x_2}+\dots+e^{x_n})-e^{x_1}e^{x_1}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} & \frac{-e^{x_1}e^{x_2}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} & \dots & \frac{-e^{x_1}e^{x_n}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} \\ \frac{-e^{x_2}e^{x_1}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} & \frac{e^{x_2}(e^{x_1}+e^{x_2}+\dots+e^{x_n})-e^{x_2}e^{x_2}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} & \dots & \frac{-e^{x_2}e^{x_n}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{-e^{x_n}e^{x_1}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} & \frac{-e^{x_n}e^{x_2}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} & \dots & \frac{e^{x_n}(e^{x_1}+e^{x_2}+\dots+e^{x_n})-e^{x_n}e^{x_n}}{(e^{x_1}+e^{x_2}+\dots+e^{x_n})^2} \end{vmatrix}$$

化简可得：

$$H = \frac{1}{(e^{x_1} + e^{x_2} + \dots + e^{x_n})^2} \begin{bmatrix} e^{x_1}(e^{x_1} + e^{x_2} + \dots + e^{x_n}) & 0 & \dots & 0 \\ 0 & e^{x_2}(e^{x_1} + e^{x_2} + \dots + e^{x_n}) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & e^{x_n}(e^{x_1} + e^{x_2} + \dots + e^{x_n}) \end{bmatrix} - \frac{1}{(e^{x_1} + e^{x_2} + \dots + e^{x_n})^2} \begin{bmatrix} e^{x_1} \\ e^{x_2} \\ \vdots \\ e^{x_n} \end{bmatrix} \begin{bmatrix} e^{x_1} & e^{x_2} & \dots & e^{x_n} \end{bmatrix}$$

用 z 表示为：

$$H = \frac{1}{(1^T z)^2} (1^T z (\text{diag}(z)) - z z^T).$$

由二阶条件可得：当 $H \geq 0$ 时，即 H 为半正定矩阵时， $f(x)$ 为凸函数。

又 $\frac{1}{1^T z}$ 一定 > 0 ，所以我们定义矩阵 $K \in R^{n \times n}$ 为 $k = 1^T z (\text{diag}(z)) - z z^T$ 。

因此当且仅当对 $\forall v \in R^n, v^T k v \geq 0$ 时， k 为半正定矩阵。

所以进行计算：

$$\begin{aligned} v^T k v &= v^T ((1^T z) (\text{diag}(z)) - z z^T) v \\ &= (1^T z) v^T \text{diag} z v - v^T z z^T k v \\ &= \left(\sum_{i=1}^n z_i \right) \left(\sum_{i=1}^n v_i^2 z_i \right) - \left(\sum_{i=1}^n v_i z_i \right)^2 \end{aligned}$$

因此令 $a = v_i \sqrt{z_i}, b = \sqrt{z_i}$ ，上式可以化成：

$$v^T k v = (b^T b) (a^T a) - (a^T b)^2$$

由 Cauchy-Schwarz 不等式可得： $v^T k v \geq 0$

因此可得 log-sum-up 函数为凸函数。

参考资料：中国科学技术大学研究生院制作课程：《最优化理论》

<https://www.bilibili.com/video/BV1es411u7pK?p=12>

1.2 梯度下降法

1.2.1 算法的收敛速率

已知 $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T Q \mathbf{x} + b^T \mathbf{x}$ 是正定二次函数，因此 $\nabla f(\mathbf{x}_k) = Q\mathbf{x}_k + b$ 。

极小值点 \mathbf{x}^* 应满足 $\nabla f(\mathbf{x}^*) = Q\mathbf{x}^* + b = 0 \Rightarrow \mathbf{x}^* = -Q^{-1}b$

假设新的迭代点为 $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha d_k$ ，其中由于使用最速下降法， $d_k = -\nabla f(\mathbf{x}_k)$

因此问题转化为考虑一维最优化问题 $\min_{\alpha} h(\alpha) = f(\mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k))$

将上式展开，可得：

$$\begin{aligned}
 \min_{\alpha} h(\alpha) &= \frac{1}{2}(x_k - \alpha \nabla f(\mathbf{x}_k))^T Q (x_k - \alpha \nabla f(\mathbf{x}_k)) + b^T (x_k - \alpha \nabla f(\mathbf{x}_k)) \\
 &= \frac{1}{2}(x_k^T - \alpha \nabla f(\mathbf{x}_k)^T)(Q x_k - \alpha \nabla Q f(\mathbf{x}_k)) + b^T x_k - \alpha b^T \nabla f(\mathbf{x}_k) \\
 &= \frac{1}{2}x_k^T Q x_k - \frac{1}{2}\alpha x_k^T Q \nabla f(\mathbf{x}_k) - \frac{1}{2}\alpha \nabla f(\mathbf{x}_k)^T Q x_k + \frac{1}{2}\alpha^2 \nabla f(\mathbf{x}_k)^T Q \nabla f(\mathbf{x}_k) + b^T x_k - \alpha b^T \nabla f(\mathbf{x}_k) \\
 &= \frac{1}{2}\alpha^2 \nabla f(\mathbf{x}_k)^T Q \nabla f(\mathbf{x}_k) - \alpha \nabla f(\mathbf{x}_k)^T (Q x_k + b) + \frac{1}{2}x_k^T Q x_k + b^T x_k \\
 &= \frac{1}{2}\alpha^2 \nabla f(\mathbf{x}_k)^T Q \nabla f(\mathbf{x}_k) - \alpha \nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k) + f(x_k) \\
 &= \frac{1}{2}\|\nabla f(x_k)\|_Q^2 \alpha^2 - \|\nabla f(x_k)\|^2 \alpha + f(x_k)
 \end{aligned}$$

这是一个关于变量 α 的二次函数。

接下来用精确线搜索确定步长：

$$h'(\alpha) = \frac{\partial h}{\partial \alpha} = \alpha \|\nabla f(x_k)\|_Q^2 - \|\nabla f(x_k)\|^2$$

令 $h'(\alpha) = 0$ ，可得：

$$\alpha_k = \frac{\|\nabla f(x_k)\|^2}{\|\nabla f(x_k)\|_Q^2}$$

因此下一个迭代点为： $\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha d_k = \mathbf{x}_k - \frac{\|\nabla f(x_k)\|^2}{\|\nabla f(x_k)\|_Q^2} \nabla f(x_k)$

令 $g(k) = \nabla f(x_k)$ ，所以 $\alpha_k = \frac{g_k^T Q g_k}{g_k^T g_k}$

因此下一个迭代点的函数值为 $f(x_{k+1}) = f(x_k) - \alpha_k g_k^T g_k + \frac{1}{2}\alpha_k^2 g_k^T Q g_k = f(x_k) - \frac{1}{2} \frac{(g_k^T g_k)^2}{g_k^T Q g_k}$

因为 $f(x^*) = \frac{1}{2}b^T Q^{-1}b$

所以可得：

$$\begin{aligned}
\frac{f(x_{k+1}) - f(x^*)}{f(x_k) - f(x^*)} &= \frac{f(x_k) - \frac{1}{2} \frac{(g_k^T g_k)^2}{g_k^T Q g_k} + \frac{1}{2} b^T Q^{-1} b}{f(x_k) + \frac{1}{2} b^T Q^{-1} b} \\
&= 1 + \frac{-\frac{1}{2} \frac{(g_k^T g_k)^2}{g_k^T Q g_k}}{\frac{1}{2} x_k^T Q x_k + b^T x_k + \frac{1}{2} b^T Q^{-1} b} \\
&= 1 - \frac{\frac{1}{2} \frac{(g_k^T g_k)^2}{g_k^T Q g_k}}{\frac{1}{2} (Q x_k + b)^T Q^{-1} (Q x_k + b)} \\
&= 1 - \frac{\frac{1}{2} \frac{(g_k^T g_k)^2}{g_k^T Q g_k}}{\frac{1}{2} g_k^T Q^{-1} g_k} \\
&= 1 - \frac{(g_k^T g_k)^2}{(g_k^T Q g_k)(g_k^T Q^{-1} g_k)}
\end{aligned}$$

因为 Q 为正定矩阵，所以由 Kantorovich 不等式可得：

$$\frac{(g_k^T g_k)^2}{(g_k^T Q g_k)(g_k^T Q^{-1} g_k)} \geq \frac{4\lambda_1 \lambda_d}{(\lambda_1 + \lambda_d)^2}, \text{ 其中 } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0, \lambda_i \text{ 为 } Q \text{ 的特征值}$$

所以代入原式可得：

$$\frac{f(x_{k+1}) - f(x^*)}{f(x_k) - f(x^*)} \leq 1 - \frac{4\lambda_1 \lambda_d}{(\lambda_1 + \lambda_d)^2} = \left(\frac{\lambda_1 - \lambda_d}{\lambda_1 + \lambda_d} \right)^2$$

由于 $Q = \text{diag}(1, 2, 10)$ ，可得： $\lambda_1 = 10, \lambda_d = 1$ ，所以

$$\left(\frac{\lambda_1 - \lambda_d}{\lambda_1 + \lambda_d} \right)^2 = \frac{81}{121}$$

因此收敛速率

$$\frac{f(x_{k+1}) - f(x^*)}{f(x_k) - f(x^*)} \leq \frac{81}{121}$$

1.2.2 要使 $f(x_T) - f(x^*) < 10^{-10}$ ，大约需要多少次迭代

因为 $b = (1, 0, 0)^T$ ，初始点 $x_0 = (1, 1, 1)^T, Q = \text{diag}(1, 2, 10)$ ，而且

$$\frac{f(x_{k+1}) - f(x^*)}{f(x_k) - f(x^*)} \leq \left(1 - \frac{2}{1 + \kappa} \right)^2$$

可得： $\kappa = \frac{\lambda_1}{\lambda_d} = 10$

所以，对于式子进行化简可得：

$$f(x_{k+1}) - f(x^*) \leq \left(1 - \frac{2}{1+\kappa}\right)^2 (f(x_k) - f(x^*)) = \left(1 - \frac{2}{1+\kappa}\right)^{2k} (f(x_0) - f(x^*))$$

令 $R = f(x_0) - f(x^*)$

要计算 $f(x_T) - f(x^*) < 10^{-10}$ 的迭代次数，即计算 $\left(1 - \frac{2}{1+\kappa}\right)^{2(T-1)} R < 10^{-10}$

接下来计算 $f(x_0) - f(x^*)$ ：

$$\begin{aligned} f(x_0) &= \frac{1}{2}(1, 1, 1) \text{diag}(1, 2, 10)(1, 1, 1)^T + (1, 0, 0)(1, 1, 1)^T = \frac{15}{2} \\ f(x^*) &= -\frac{1}{2}(1, 0, 0) \text{diag}(1, 2, 10)(1, 0, 0)^T = -\frac{1}{2} \end{aligned}$$

所以 $R=8$

所以转化为计算：

$$\left(1 - \frac{2}{11}\right)^{2t} 8 < 10^{-10}$$

将其转化为对数形式进行计算可得： $t \geq 63$

所以大约需要 63 次迭代。

1.3 强凸与光滑问题判定

1.3.1 计算 $f(x)$ 相对于 x 的梯度

由题意可得：

$$\begin{aligned} f(x) &= \frac{1}{2}\|Ax - y\|^2 = \frac{1}{2}(Ax - y)^T(Ax - y) = \frac{1}{2}(x^T A^T - y^T)(Ax - y) \\ &= \frac{1}{2}(x^T A^T Ax - x^T A^T y - y^T Ax + y^T y) \end{aligned}$$

因此要求 $f(x)$ 相对于 x 的梯度可转化为求 $\frac{\partial f}{\partial x}$

因为

$$\frac{\partial x^T Ax}{\partial x} = (A + A^T)x, \frac{\partial A^T x}{\partial x} = A, \frac{\partial x^T A}{\partial x} = A$$

所以对上式进行化简：

$$\frac{\partial f}{\partial x} = \frac{1}{2} \left((A^T A + (A^T A)^T) x - A^T y - A^T y \right) = \frac{1}{2} (2A^T Ax - 2A^T y) = A^T(Ax - y)$$

所以 $\nabla f(x) = A^T(Ax - y)$

1.3.2 设 $A = \text{diag}(-1, 5)$, 证明 $f(x)$ 是强凸函数, 计算参数 μ

要证 $f(x)$ 是强凸函数, 只需构造 $g(x) = f(x) - \frac{1}{2}\|x\|^2$ 并证明其为凸函数证明:

$$\begin{aligned} g(x) &= f(x) - \frac{1}{2}\|x\|^2 \\ \nabla g(x) &= A^T(Ax - y) - x \\ \nabla g(x)^2 &= A^T A - I = \begin{bmatrix} 0 & 0 \\ 0 & 24 \end{bmatrix} \geq 0 \end{aligned}$$

因此, $f(x)$ 为强凸函数

若 $f(x)$ 是强凸函数且二阶连续可微, 那么 $\nabla^2 f(x) \geq \mu I$

因为 $\nabla^2 f(x) = \frac{\partial \nabla f(x)}{\partial x} = AA^T$

因此可以转化为: $AA^T \geq \mu I$

因为 $A = \text{diag}(-1, 5)$, 所以:

$$AA^T = \begin{bmatrix} 1 & 0 \\ 0 & 25 \end{bmatrix} \geq \mu I$$

因此可得: $\mu = 1$

1.3.3 设 $A = \text{diag}(-1, 5)$, 证明 $f(x)$ 是光滑函数, 计算光滑函数的参数 L

要证 $f(x)$ 是光滑函数, 只需构造 $g(x) = \frac{25}{2}\|x\|^2 - f(x)$ 并证明其为凸函数证明:

$$\begin{aligned} g(x) &= \frac{25}{2}\|x\|^2 - f(x) \\ \nabla g(x) &= 25x - A^T(Ax - y) \\ \nabla g(x)^2 &= 25I - A^T A = \begin{bmatrix} 24 & 0 \\ 0 & 0 \end{bmatrix} \geq 0 \end{aligned}$$

因此, $f(x)$ 为光滑函数

若 $f(x)$ 是光滑函数且二阶连续可微, 那么 $\nabla^2 f(x) \leq LI$

因为 $\nabla^2 f(x) = \frac{\partial \nabla f(x)}{\partial x} = AA^T$

因此可以转化为: $AA^T \leq LI$

因为 $A = \text{diag}(-1, 5)$, 所以:

$$AA^T = \begin{bmatrix} 1 & 0 \\ 0 & 25 \end{bmatrix} \leq LI$$

因此可得： $L = 25$

1.3.4 设 $A = \text{diag}(1, 2)$ ，并且步长为 $\alpha_k = \frac{1}{4}$ ，求算法的单步下降幅度 $f(x_{k+1}) - f(x_k)$

由光滑函数的定义：

$$f(y) \leq f(x) + \nabla f^T(x)(y - x) + \frac{L}{2} \|y - x\|^2$$

因为步长： $x_{k+1} = x_k - \frac{1}{4} \nabla f(x_k) = x_k - \frac{1}{4} [A^T (Ax - y)]$

那么对于上式进行化简可得：

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \nabla^T f(x_k)(x_{k+1} - x_k) + \frac{L}{2} \|x_{k+1} - x_k\|^2 \\ &= f(x_k) - \frac{1}{4} \|\nabla f(x_k)\|^2 + \frac{1}{2L} \|\nabla f(x_k)\|^2 \end{aligned}$$

由 1.3.3 可得：因为 $f(x)$ 是光滑函数，所以 $\nabla^2 f(x) \leq LI$

因为 $A = \text{diag}(1, 2)$ ，所以 $\nabla^2 f(x) = AA^T = \text{diag}(1, 4) \leq LI \Rightarrow L = 4$

所以可得：

$$f(x_{k+1}) - f(x_k) \leq -\frac{1}{4} \|\nabla f(x_k)\|^2 + \frac{1}{8} \|\nabla f(x_k)\|^2 = -\frac{1}{8} \|\nabla f(x_k)\|^2$$

因此算法的单步下降幅度 $f(x_{k+1}) - f(x_k) = \frac{1}{8} \|\nabla f(x_k)\|^2$

1.3.5 设 $A = \text{diag}(1, 2)$ ，算法使用固定步长，问如果要使迭代序列是下降序列，步长应满足什么条件

设要使用固定步长为 α ，且 $A = \text{diag}(1, 2)$ ，由 1.3.4 可得：

$$f(x_{K+1}) - f(x_k) \leq -(\alpha - \frac{1}{8}) \|\nabla f(x_k)\|^2$$

要保证为下降序列，可转化为： $\alpha - \frac{1}{8} > 0$

因此步长 $\alpha > \frac{1}{8}$

1.3.6 设 $A = \text{diag}(1, 2)$, 如果 $f(x_0) - f(x^*) = 1$, 要使 $f(x_k) - f(x^*) < 10^{-3}$, 试估计需要多少次迭代

由 3.1 ~ 3.5 的证明可得: $f(x)$ 是光滑且强凸的, 并且它的 Hessian 矩阵满足 $\mu I \leq \nabla^2 f(x) \leq LI$

因此梯度下降法的收敛速率为 $(1 - \frac{\mu}{L})$

$$\begin{aligned} A = \text{diag}(1, 2) &\Rightarrow \mu I \leq \nabla^2 f(x) = AA^T = \text{diag}(1, 4) \leq LI \Rightarrow \mu = 1, L = 4 \\ &\Rightarrow 1 - \frac{\mu}{L} = \frac{3}{4} \end{aligned}$$

又因为 $f(x_0) - f(x^*) = 1$

所以要求迭代次数, 可转化为:

$$\begin{aligned} f(x_k) - f(x^*) &= \left(\frac{3}{4}\right)^k f(x_0) - f(x^*) < 10^{-3} \\ &\Rightarrow k(\ln 3 - \ln 4) < -3 \ln 10 \\ &\Rightarrow k > 24.0118 \end{aligned}$$

因此大约需要 25 次迭代.

1.4 共轭梯度法性质

1.4.1 求 $x^{(1)}$ 处的搜索方向

因为在共轭梯度法中, 梯度等于残差, 并且由题意得: $f: \mathbf{R}^3 \rightarrow \mathbf{R}$,

$$\left. \frac{\partial f(x)}{\partial x_1} \right|_{x^{(1)}} = \left. \frac{\partial f(x)}{\partial x_2} \right|_{x^{(1)}} = -2$$

因此可得: (其中 a 是由于我们不知道残差 r_1 的第三个分量, 所以进行假设)

$$r_1 = \begin{bmatrix} -2 \\ -2 \\ a \end{bmatrix}, d_0 = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix}, r_0 = -d_0 = \begin{bmatrix} -1 \\ 1 \\ -2 \end{bmatrix}$$

因为如下公式:

$$\alpha_0 = \frac{r_0^T r_0}{d_0^T G d_0}, r_1 - r_0 = \alpha_0 G d_0, \beta_0 = \frac{r_1^T r_1}{r_0^T r_0}, d_1 = -r_1 + \beta_0 d_0$$

所以我们可以设:

$$\alpha_0 Gd_0 = \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

进行如下计算：

$$d_0^T \cdot \alpha_0 Gd_0 = r_0^T r_0 = \|r_0\|^2 = 6 \Rightarrow a_1 - a_2 + 2a_3 = 6$$

结合 $r_1 - r_0 = \alpha_0 Gd_0$ ，我们可以得到如下式子：

$$\begin{cases} -2 + 1 = a_1 \\ -2 - 1 = a_2 \\ a + 2 = a_3 \\ a_1 - a_2 + 2a_3 = 6 \end{cases}$$

解得：

$$a_1 = -1, a_2 = -3, a_3 = 2, a = 0, \alpha_0 Gd_0 = \begin{bmatrix} -1 \\ -3 \\ 2 \end{bmatrix}, r_1 = \begin{bmatrix} -2 \\ -2 \\ 0 \end{bmatrix}$$

接着进行计算可得：

$$\beta_0 = \frac{4}{3}, d_1 = -r_1 + \beta_0 d_0 = \begin{bmatrix} \frac{10}{3} \\ \frac{2}{3} \\ \frac{8}{3} \end{bmatrix}$$

因此可得到 $x^{(1)}$ 处的搜索方向 d_1 。

1.4.2 如果在 $x^{(2)}$ 处梯度长度为 $\|\nabla f(x^{(2)})\| = 1$ ，求 $\nabla f(x^{(2)})$

已知：

$$r_1 = \begin{bmatrix} -2 \\ -2 \\ 0 \end{bmatrix} \quad d_1 = \begin{bmatrix} \frac{10}{3} \\ \frac{2}{3} \\ \frac{8}{3} \end{bmatrix}$$

由共轭梯度法可得：残差等与梯度，因此要求 $\nabla f(x^{(2)})$ 可转化为求 r_2 ，设 $r_2 = [b_1, b_2, b_3]^T$

因为 $\|\nabla f(x^{(2)})\| = 1 \Rightarrow \|r_2\| = 1 \Rightarrow b_1^2 + b_2^2 + b_3^2 = 1$

因为 r_2 与已有共轭方向 d_0, d_1 正交，所以对以上式子进行联立可得：

$$\begin{cases} b_1 - b_2 + 2b_3 = 0 \\ 5b_1 + b_2 + 4b_3 = 0 \\ b_1^2 + b_2^2 + b_3^2 = 1 \end{cases}$$

可得：

$$b_1 = -\frac{1}{\sqrt{3}}, b_2 = \frac{1}{\sqrt{3}}, b_3 = \frac{1}{\sqrt{3}}$$

所以可得：

$$\nabla f(x^{(2)}) = \left[-\frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{3}} \right]^T$$

1.5 共轭梯度法的等价性

1.5.1 给定下降方向 d ，计算精确线搜索确定的步长

已知： $f(x) = \frac{1}{2}x^T Qx + b^T x + c$ ，下降方向为 d ，因此精确线搜索确定的步长可以转化为： $\alpha^* = \arg \min_{\alpha} f(x_k + \alpha d)$ ，令 $g = -d$

$$\begin{aligned} \alpha^* &= \arg \min_{\alpha} f(x_k - \alpha g) \\ &= \arg \min_{\alpha} \left(\frac{1}{2}(x_k - \alpha g)^T Q(x_k - \alpha g) + b^T(x_k - \alpha g) + c \right) \\ &= \arg \min_{\alpha} \left(\frac{1}{2}(x_k^T - \alpha g^T)(Qx_k - \alpha Qg) + b^T x_k - \alpha b^T g + c \right) \\ &= \arg \min_{\alpha} \left(\frac{1}{2}x_k^T Qx_k - \frac{1}{2}\alpha x_k^T Qg - \frac{1}{2}\alpha g^T Qx_k + \frac{1}{2}\alpha^2 g^T Qg + b^T x_k - \alpha b^T g + c \right) \\ &= \arg \min_{\alpha} \left(\frac{1}{2}\alpha^2 g^T Qg - \alpha g^T(Qx_k + b) + \frac{1}{2}x_k^T Qx_k + b^T x_k + c \right) \\ &= \arg \min_{\alpha} \left(\frac{1}{2}g^T Qg\alpha^2 - g^T g\alpha + f(x_k) \right) \end{aligned}$$

可以看出上述函数是关于 α 的二次函数

所以精确线搜索确定的步长为：

$$\alpha^* = \frac{g^T g}{g^T Qg} = \frac{d^T d}{d^T Qd}$$

1.5.2 证明, 在上述问题上, 使用精确线搜索确定步长时, 三种共轭梯度法所确定的点列 $\{x_i\}$ 是相同的

由共轭梯度法的计算过程可以看出三种方法的差别主要在为 β 的迭代公式不同, 因此要证明三种共轭梯度法所确定的点列是否相同, 即证明 β 迭代公式的等价性。

这边给出 PPT 上的六种 β 迭代公式, 其中前三种已经在 PPT 上给出证明, 因此这里主要对于后三种进行证明:

$$\beta_{k-1}^{HS} = \frac{g_k^T G d_{k-1}}{d_{k-1}^T G d_{k-1}} \quad (\text{Hestenes - Stiefel 公式})$$

$$\beta_{k-1}^{CW} = \frac{g_k^T (g_k - g_{k-1})}{d_{k-1}^T (g_k - g_{k-1})} \quad (\text{Crowder - Wolfe 公式})$$

$$\beta_{k-1}^{FR} = \frac{g_k^T g_k}{g_{k-1}^T g_{k-1}} \quad (\text{Fletcher - Reeves 公式})$$

$$\beta_{k-1}^{PRP} = \frac{g_k^T (g_k - g_{k-1})}{g_{k-1}^T g_{k-1}} \quad (\text{Polak - Ribiere - Polyak 公式})$$

$$\beta_{k-1}^D = -\frac{g_k^T g_k}{d_{k-1}^T g_{k-1}} \quad (\text{Dixon 公式})$$

$$\beta_{k-1}^{DY} = \frac{g_k^T g_k}{d_{k-1}^T (g_k - g_{k-1})} \quad (\text{Dai - Yuan 公式})$$

已知: $d_k^T g_{k+1} = -g_k^T g_{k+1} + \beta_{k-1} d_{k-1}^T g_{k+1}$, 由子空间扩展定理: $d_k^T g_{k+1} = 0, d_{k-1}^T g_{k+1} = 0$, 可得: $g_{k+1}^T g_k = 0$

由梯度的正交性 $g_k^T d_{k-1} = 0$, 因此可由 Crowder-Wolfe 公式推出 PRP 公式。

因为 $d_{k-1}^T g_{k-1} = -g_{k-1}^T g_{k-1}$, 代入 Crowder-Wolfe 公式可以推出 Dixon 公式。

由梯度的正交性: $g_k^T g_{k-1} = 0$, 可以由 Crowder-Wolfe 公式推出 Dai-Yuan 公式。

1.6 共轭梯度法的子空间性质与收敛性

1.6.1 每次迭代的搜索方向 d_i 都可以写成下面向量组的线性组合

$$\{d, Qd, Q^2d, \dots, Q^i d\}$$

由子空间扩展定理可得: 每次迭代 g_k 与整个子空间 $\text{span}\{d_0, d_1, \dots, d_{k-1}\}$ 正交, 与之前搜索过的所有方向均相互独立, 并且正交。

因为 $r_i = g_i = Qx_i - b$, 可得:

$$\begin{aligned}
r_{i+1} &= g_{i+1} = Qx_{i+1} - b = Q(x_i + \alpha_i d_i) - b \\
&= Qx_i - b + \alpha_i Qd_i \\
&= r_i + \alpha_i Qd_i
\end{aligned}$$

因此要证明每次搜索方向 d_i 可以写成 $\{d, Qd, Q^2d, \dots, Q^i d\}$ 的线性组合, 可用数学归纳法证明如下:

当 $i = 0$ 时, $d_0 = g_0$, 易得: 显然 d_0 可以写成 d_0 自己的线性组合。

假设 $r_i, d_i \in \text{span}\{d, Qd, Q^2d, \dots, Q^i d\}$

则由上式证明可得:

$$r_{i+1} = r_i + \alpha_i Qd_i \in \text{span}\{d, Qd, Q^2d, \dots, Q^i d, Q^{i+1}d\}$$

$$d_{i+1} = r_{i+1} + \beta d_i \in \text{span}\{d, Qd, Q^2d, \dots, Q^i d, Q^{i+1}d\}$$

因此可得每次搜索方向 d_i 可以写成 $\{d, Qd, Q^2d, \dots, Q^i d\}$ 的线性组合。

1.6.2 对于第 $i + 1$ 次迭代的序列 x_{i+1} , 证明, 对于任意由上述向量线性组合的向量 x' , 都有 $f(x') \geq f(x_{i+1})$

由共轭梯度法可得:

$$\begin{aligned}
x_{i+1} &= x_i + \alpha_i d_i \\
&= x_0 + \sum_{k=0}^i \alpha_k d_k = x_0 + \begin{bmatrix} \alpha_0 & \alpha_1 & \dots & \alpha_i \end{bmatrix} \begin{bmatrix} d_0 \\ d_1 \\ \vdots \\ d_i \end{bmatrix} \\
&= x_0 + A^T D
\end{aligned}$$

令子空间:

$$\begin{aligned}
V &= x_0 + \text{span}\{d, Qd, Q^2d, \dots, Q^i d\} \\
&= x_0 + \text{span}\{d, d_1, d_2, \dots, d_i\}
\end{aligned}$$

所以对于任意迭代点 x' , 可转化为:

$$\begin{aligned}
x' &= x_0 + \begin{bmatrix} c_0 & c_1 & \cdots & c_i \end{bmatrix} \begin{bmatrix} d_0 \\ d_1 \\ \vdots \\ d_i \end{bmatrix} \\
&= x_0 + c^T D \in V
\end{aligned}$$

因此，要证明对于任意由上述向量线性组合的向量 x' ，都有 $f(x') \geq f(x_{i+1})$ ，可以转化为求 $\arg \min_{x' \in V} f(x')$
即 $\nabla f(x') = \nabla f(x_0 + c^T D)^T D = 0$
由共轭的条件可知： $\nabla f(x_{i+1})^T D = 0$
即 $x' = x_{i+1}$ 的时候有最小值
因此 $f(x') \geq f(x_{i+1})$

1.6.3 使用上述结论证明，共轭梯度法在二次函数 $f(x)$ 上迭代 K 次即可收敛到最优点，其中 K 是矩阵 Q 的不同的特征值的数目

通过翻阅参考资料《Numerical Optimization》可知：

$$\|x_{k+1} - x^*\|_Q^2 \leq \min_{P_k} \max_{1 \leq i \leq n} [1 + \lambda_i P_k(\lambda_i)]^2 \|x_0 - x^*\|_Q^2$$

假设 Q 的 N 个特征值只能取 k 个互不相同的值且满足 $k_1 < k_2 < \cdots < k_k$
那么令：

$$Q_k(\lambda) = \frac{(-1)^k}{k_1 k_2 \cdots k_k} (\lambda - k_1)(\lambda - k_2) \cdots (\lambda - k_k)$$

对上式进行进一步化简变换可得：

$$Q_k(\lambda_i) = 0 \quad i = 1, 2, \cdots, k$$

$$Q_k(0) = 1$$

$$Q_k(\lambda) - 1 \text{ 为 } k \text{ 阶多项式}$$

因此可得：

$$\bar{P}_{k-1}(\lambda) = \frac{Q_k(\lambda) - 1}{\lambda} \text{ 为 } k-1 \text{ 阶多项式}$$

所以可得出以下结论：

$$0 \leq \min_{P_{k-1}} \max_{1 \leq i \leq n} [1 + \lambda_i P_{k-1}(\lambda_i)]^2 \leq \max_{1 \leq i \leq n} [1 + \lambda_i \bar{P}_{k-1}(\lambda_i)]^2 = \max_{1 \leq i \leq n} Q_k^2(\lambda_i) = 0$$

因此可证明: k 次迭代后 $\|x_k - x^*\|_Q^2 = 0$

所以, $x_k = x^*$, 即共轭梯度法在二次函数 $f(x)$ 上迭代 K 次即可收敛到最优点.

1.7 约束优化问题的 KKT 条件

1.7.1 写出其 KKT 条件, 并分析 KKT 点

题目可进行以下转化:

$$\begin{cases} \arg \max_v & tr(v^T S v) \\ s.t. & v^T v = I_k \end{cases}$$

构造拉格朗日函数 $L(v, \lambda) = tr(v^T S v) + \lambda(I_k - v^T v)$

所以求一阶导可得:

$$\nabla_v L(v, \lambda) = (S + S^T)v - 2\lambda v, \quad \nabla_\lambda = I_k - v^T v$$

因此令 $\nabla_v L(v, \lambda) = 0$, 可得:

$$(S + S^T)v - 2\lambda v = 0 \Rightarrow \lambda v^* = \frac{S + S^T}{2} v^*$$

即 v^* 是矩阵 $\frac{S+S^T}{2}$ 的对应于特征值 λ 的特征向量。

由 $\lambda v^* = \frac{S+S^T}{2} v^*$, 等式两边左乘 v^{*T} , 并有 $v^T v = I_k$ 可得:

$$\lambda^* = v^{*T} \left(\frac{S + S^T}{2} \right) v^*$$

因此 λ^* 为 $f(v) = v^T \left(\frac{S+S^T}{2} \right) v$ 的最大值。

1.7.2 如果使用投影梯度法进行求解上述问题, 写出投影公式

投影公式如下:

- 如果 $v_t + \alpha_t d_t \in V$, 则 $v_{t+1} = v_t + \alpha_t d_t$
- 如果 $v_t + \alpha_t d_t \notin V$, 则 $v_{t+1} = \arg \min_{x \in V} \|x - (v_t + \alpha_t d_t)\|$

2 思考题

2.1 共轭方向与特征向量

2.1.1 证明：方向组 $1, \dots, n$ 是关于 Q 相互共轭的

依题意得： Q 是正定对称矩阵，并且其特征分解为 $Q = U\Lambda U^T$ ，且特征值满足 $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0$ ，其对应的特征向量分别为 $1, 2, \dots, n$ 。

由对称矩阵特征分解的定理可得：矩阵 $U = [u_1, u_2, \dots, u_n]$ 为正交矩阵，且满足 $u_i^T u_j = \delta_{ij}$ ，其中当 $i=j$ 的时候为 1，否则为 0。

因此要证明方向组 $1, \dots, n$ 是关于 Q 相互共轭的，可将方向组用 u_1, u_2, \dots, u_n 表示，对于 $i \neq j$ ：

$$\begin{aligned} u_i^T Q u_j &= u_i^T \begin{pmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \dots & \\ & & & \lambda_n \end{pmatrix} u_j \\ &= u_i^T \begin{pmatrix} 0 \\ \dots \\ \lambda_j \\ \dots \\ 0 \end{pmatrix} \\ &= 0 \end{aligned}$$

这满足共轭的定义，因此得证：方向组 $1, \dots, n$ 是关于 Q 相互共轭的。

2.1.2 对于任意 $k = 1, \dots, n-1$ ，求单步改进幅度 $f(x_{k+1}) - f(x_k)$

由题意得： $x_{k+1} = x_k - \beta_k k$ ，其中 k 指的是矩阵 Q 的特征向量 u_i ，且由第一题得到特征向量是共轭方向组，因此可得： $k = -\nabla f(x_k)$ ，所以直接展开进行计算：

其中涉及到的有 2.1.1 的共轭方向组的性质以及对称矩阵的特征分解一些性质，例如当 $i=j$ 的时候， $u_i^T u_j = 1$ 。

$$\begin{aligned}
f(x_{k+1}) - f(x_k) &= f(x_k - \beta_k k) - f(x_k) \\
&= \frac{1}{2}(x_k - \beta_k k)^T Q(x_k - \beta_k k) + b^T(x_k - \beta_k k) - \frac{1}{2}x_k^T Q x_k - b^T x_k \\
&= \frac{1}{2}(x_k^T Q - \beta_k k^T Q)(x_k - \beta_k k) + b^T x_k - b^T \beta_k k - \frac{1}{2}x_k^T Q x_k - b^T x_k \\
&= \frac{1}{2}(x_k^T Q x_k) - \frac{1}{2}\beta_k x_k^T Q k - \frac{1}{2}\beta_k k^T Q x_k + \frac{1}{2}\beta_k^2 k^T Q k - \frac{1}{2}x_k Q^T x_k - \beta_k b^T k \\
&= \frac{1}{2}\beta_k^2 k^T Q k - \beta_k k^T (Q x_k + b) \\
&= \frac{1}{2}\beta_k^2 k^T \lambda_k k + \beta_k \nabla f^T(x_k) \nabla f(x_k) \\
&= \frac{1}{2}\beta_k^2 \lambda_k + \beta_k \|\nabla f(x_k)\|^2
\end{aligned}$$

因此单步改进幅度为： $\frac{1}{2}\beta_k^2 \lambda_k + \beta_k \|\nabla f(x_k)\|^2$

2.1.3 若对于 $k = 1, \dots, n-1$ ，每一次迭代的改进幅度 $\delta_k = |f(x_{k+1}) - f(x_k)|$ 都相等，则步长 β_k 满足什么条件

结合 2.1.2 可得，因为每一次迭代的改进幅度 $\delta_k = |f(x_{k+1}) - f(x_k)|$ 都相等

$$\frac{1}{2}\beta_k^2 \lambda_k + \beta_k \|\nabla f(x_k)\|^2 = \frac{1}{2}\beta_{k+1}^2 \lambda_{k+1} + \beta_{k+1} \|\nabla f(x_{k+1})\|^2$$

写到这里写不下去了，这个式子不知道该怎么化简。

3 编程题

3.1 使用 Newton-CG 与 L-BFGS 法求解 LASSO 问题

3.1.1 Newton-CG 法

用 python 实现的 Newton-CG 法代码：

```

1  # 导入相关包
2
3  import numpy as np
4  import matplotlib.pyplot as plt
5
6  # 写入目标函数
7
8  def testFun(x, y):

```



```

9         t = 4.0*x - 1.0*y
10        t1 = x - 1.0
11        z = np.power(t, 2) + np.power(t1, 4)
12        return z
13
14    # 求函数的梯度
15
16    def gradTestFun(x, y):
17        delta_x = 1e-6    #x方向差分小量
18        delta_y = 1e-6    #y方向差分小量
19        grad_x = (testFun(x+delta_x, y)-testFun(x-delta_x, y))/(2.0*delta_x)
20        grad_y = (testFun(x, y+delta_y)-testFun(x, y-delta_y))/(2.0*delta_y)
21        grad_xy = np.array([grad_x, grad_y])
22        return grad_xy
23
24    # 采用牛顿法, 精确线性搜索确定移动步长
25
26    def getStepLengthByNewton(array_xy, array_d):
27        a0 = 1.0          #初始猜测值
28        e0 = 1e-6          #退出搜索循环的条件
29        delta_a = 1e-6     #对a作差分的小量
30        while(1):
31            new_a = array_xy + a0*array_d
32            new_a_l = array_xy + (a0-delta_a)*array_d
33            new_a_h = array_xy + (a0+delta_a)*array_d
34            diff_a0 = (testFun(new_a_h[0], new_a_h[1]) - testFun(new_a_l[0], new_a_l[1]))
35            / (2.0*delta_a)
36            if np.abs(diff_a0) < e0:
37                break
38            ddiff_a0 = (testFun(new_a_h[0], new_a_h[1]) + testFun(new_a_l[0], new_a_l[1]) -
39                2.0*testFun(new_a[0], new_a[1]))/(delta_a*delta_a)
40            a0 = a0 - diff_a0/ddiff_a0
41        return a0
42
43    # 可视化
44
45    def plotResult(array_xy_history):
46        x = np.linspace(-1.0, 4.0, 100)
47        y = np.linspace(-4.0, 8.0, 100)
48        X, Y = np.meshgrid(x, y)
49        Z = testFun(X, Y)
50        plt.figure(dpi=300)
51        plt.xlim(-1.0, 4.0)
52        plt.ylim(-4.0, 8.0)
53        plt.xlabel("x")
54        plt.ylabel("y")
55        plt.contour(X, Y, Z, 40)

```

```

54     plt.plot(array_xy_history[:,0], array_xy_history[:,1], marker='.', ms=10)
55     xy_count = array_xy_history.shape[0]
56     for i in range(xy_count):
57         if i == xy_count-1:
58             break
59         dx = (array_xy_history[i+1][0] - array_xy_history[i][0])*0.6
60         dy = (array_xy_history[i+1][1] - array_xy_history[i][1])*0.6
61         plt.arrow(array_xy_history[i][0], array_xy_history[i][1], dx, dy, width=0.1)
62
63     # 使用CG算法优化，用FR公式计算组合系数
64
65     def mainFRCG():
66         ls_xy_history = [] #存储初始坐标的迭代结果
67         xy0 = np.array([4.0, -2.0]) #初始点
68         grad_xy = gradTestFun(xy0[0], xy0[1])
69         d = -1.0*grad_xy #初始搜索方向
70         e0 = 1e-6 #迭代退出条件
71         xy = xy0
72         while(1):
73             ls_xy_history.append(xy)
74             grad_xy = gradTestFun(xy[0], xy[1])
75             tag_reach = np.abs(grad_xy) < e0
76             if tag_reach.all():
77                 break
78             step_length = getStepLengthByNewton(xy, d)
79             xy_new = xy + step_length*d
80             grad_xy_new = gradTestFun(xy_new[0], xy_new[1])
81             b = np.dot(grad_xy_new, grad_xy_new)/np.dot(grad_xy, grad_xy) #根据FR公式
            # 计算组合系数
82             d = b*d - grad_xy_new
83             xy = xy_new
84             array_xy_history = np.array(ls_xy_history)
85             plotResult(array_xy_history)
86             return array_xy_history
87
88     # 主函数
89
90     mainFRCG()
91

```

选择目标函数后的可视化结果：

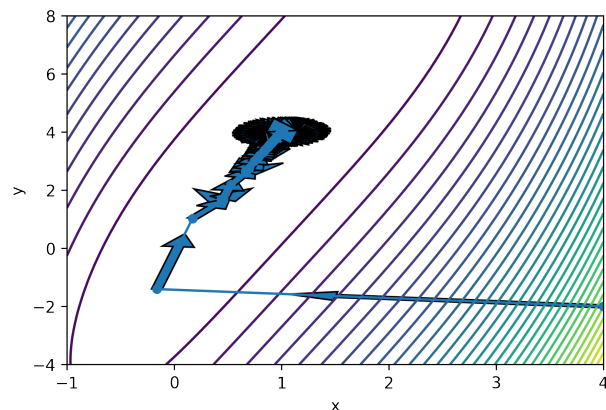


图 1: Newton-CG 可视化结果

3.1.2 应用牛顿-共轭梯度法解逻辑回归问题

看了那个题目提供的界面的代码后，发现里面的数据集是下载不下来的，只好换了个数据集进行运行，然后题目中的可视化纵坐标使用的是迭代点的梯度 2 范数，这个需要使用自动微分，不知道咋写，所以只能将纵坐标设为损失函数：所有模型误差的平方和进行可视化。

代码如下：

```

1  import numpy as np
2
3  # 读取文件的内容转化成数组的形式
4
5  def file2matrix(filename,length=2):
6      fr=open(filename)
7      lines=fr.readlines()
8      numberoflines=len(lines)
9      returnmat=np.zeros((numberoflines,length))
10     index=0
11     for line in lines:
12         line=line.strip()
13         listofline=line.split(' ')
14         returnmat[index,:]=listofline
15         index=index+1
16     return returnmat
17
18 # 数据归一化
19
20 def normalizefeature(x):
21     x_norm=x
22     mu=np.mean(x,axis=0)

```

```

23     sigma=np.std(x,axis=0)
24     x_norm=(x-mu)/sigma
25     return x_norm,mu,sigma
26
27 # 数据的可视化
28
29 import matplotlib.pyplot as plt
30
31 def plot_data(*args):
32     X = file2matrix('./ex4x.dat')
33     y = file2matrix('./ex4y.dat', 1)
34
35     X,_=normalizefeature(X)
36     pos = list(np.where(y == 1.0)[0])
37     X_pos = X[pos]
38     neg = list(np.where(y == 0.0)[0])
39     X_neg = X[neg]
40     plt.plot(X_pos[:,0],X_pos[:,1],'+',label='admitted')
41     plt.plot(X_neg[:,0],X_neg[:,1],'o',label='Not admitted')
42     plt.xlabel("exam1 score")
43     plt.ylabel("exam2 score")
44     plt.legend()
45     plt.show()
46
47 # 假设函数为Sigmoid函数
48
49 def sigmoid(z):
50     res=1/(1+np.exp(-z))
51     return res
52
53 # 计算代价的文件
54
55 def compute_loss(x,y,theta):
56     loss=-np.sum((np.dot(y.T,np.log(sigmoid(np.dot(x,theta))))
57     +np.dot((1-y).T,np.log(1-sigmoid(np.dot(x,theta)))))/x.shape[0])
58     return loss
59
60 # 牛顿法进行参数的更新
61
62 a=np.diag(np.array([1,2]))
63 print(a)
64
65 def nt(x,y,theta,iterations):
66     n,m=x.shape
67     J_loss=[]
68     orig_loss=np.inf
69     for i in range(iterations):

```

```

70         l=compute_loss(x,y,theta)
71         J_loss.append(l)
72         h=sigmoid(np.dot(x,theta))#(n,1)
73         j_first_order=1/n*np.dot(x.T,h-y)#(m,1)
74         j_second_order=1/n*np.dot(np.dot(np.dot(x.T,np.diag(h.reshape(n))),np.diag(1-h.
reshape(n))),x)
75         theta=theta-np.dot(np.linalg.inv(j_second_order),j_first_order)
76     return theta,J_loss
77
78     if __name__=="__main__":
79         X = file2matrix('./ex4x.dat')
80         y = file2matrix('./ex4y.dat', 1)
81
82         n, m = X.shape
83         X = np.column_stack((np.ones(n), X))
84         m = m + 1
85         theta = np.zeros((m, 1))
86
87         theta, J_his = nt(X, y, theta, 5)
88
89         print("theta", theta)
90         print("Loss", J_his)
91         plt.xlabel("iteration")
92         plt.ylabel("Loss")
93         plt.plot(np.arange(5), J_his)
94         plt.show()
95
96         pos = list(np.where(y == 1.0)[0])
97         X_pos = X[pos,1:3]
98         neg = list(np.where(y == 0.0)[0])
99         X_neg = X[neg,1:3]
100         plt.plot(X_pos[:, 0], X_pos[:, 1], '+', label='admitted')
101         plt.plot(X_neg[:, 0], X_neg[:, 1], 'o', label='Not admitted')
102         plt.xlabel("exam1 score")
103         plt.ylabel("exam2 score")
104         plt.legend()
105
106         xx = np.linspace(20, 70, 6)
107         yy = []
108         for i in xx:
109             res = (i * -(theta[1][0]) - (theta[0][0])) / (theta[2][0])
110             yy.append(res)
111         plt.plot(xx, yy)
112         plt.show()
113

```

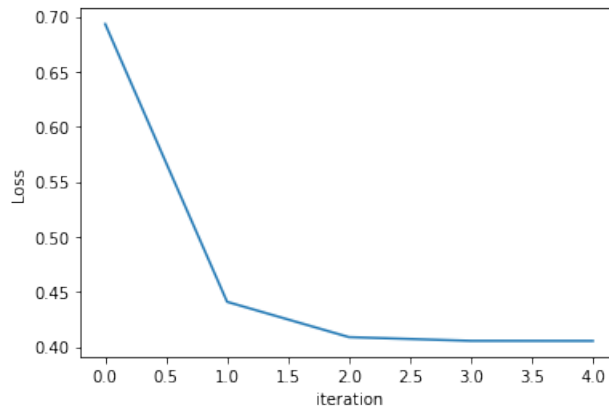


图 2: Newton-CG 的逻辑回归问题可视化结果

3.1.3 L-BFGS 算法的 MATLAB 实现

用 python 实现的 L-BFGS 算法代码:

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3  from numpy import empty
4  from numpy import dot
5
6  # 向量转化
7
8  def twoloop(s, y, rho, gk):
9
10     n = len(s) #向量序列的长度
11
12     if np.shape(s)[0] >= 1:
13         h0 = 1.0*np.dot(s[-1],y[-1])/np.dot(y[-1],y[-1])
14     else:
15         h0 = 1
16
17     a = empty((n,))
18
19     q = gk.copy()
20     for i in range(n - 1, -1, -1):
21         a[i] = rho[i] * dot(s[i], q)
22         q -= a[i] * y[i]
23     z = h0*q
24
25     for i in range(n):
26         b = rho[i] * dot(y[i], z)

```

```

27         z += s[i] * (a[i] - b)
28
29     return z
30
31     # LBFGS
32
33     def lbfgs(fun, gfun, x0, m=5):
34         # fun和gfun分别是目标函数及其一阶导数,x0是初值,m为储存的序列的大小
35         maxk = 5000
36         rou = 0.55
37         sigma = 0.4
38         epsilon = 1e-5
39         k = 0
40         W = np.zeros((2, 10 ** 3))
41         n = np.shape(x0)[0]
42         Bk = np.eye(n) # 初始对称正定矩阵, Bk=np.linalg.inv(hess(x0))
43         W[:, 0] = x0
44
45         s, y, rho = [], [], []
46
47         while k < maxk :
48             gk = gfun(x0)
49             if np.linalg.norm(gk) < epsilon:
50                 break
51
52             dk = -1.0*twoloop(s, y, rho, gk)
53
54             m0=0;
55             mk=0
56             while m0 < 20: # 用Armijo搜索求步长
57                 if fun(x0+rou**m0*dk) < fun(x0)+sigma*rou**m0*np.dot(gk,dk):
58                     mk = m0
59                     break
60                 m0 += 1
61
62             x = x0 + rou**mk*dk
63             sk = x - x0
64             yk = gfun(x) - gk
65
66             if np.dot(sk, yk) > 0: #增加新的向量
67                 rho.append(1.0/np.dot(sk, yk))
68                 s.append(sk)
69                 y.append(yk)
70             if np.shape(rho)[0] > m: #弃掉最旧向量
71                 rho.pop(0)
72                 s.pop(0)
73                 y.pop(0)

```

```

74
75         k += 1
76         x0 = x
77
78         W = W[:, 0:k] # 记录迭代点
79         return x0, fun(x0), k, W # 分别是最优点坐标, 最优值, 迭代次数
80
81 # 主函数及可视化
82
83 # 函数表达式 fun
84 fun = lambda x: 100*(x[0]**2-x[1])**2 + (x[0]-1)**2
85
86 # 梯度向量 gfun
87 gfun = lambda x: np.array([400*x[0]*(x[0]**2-x[1])+2*(x[0]-1), -200*(x[0]**2-x[1])])
88
89 # 海森矩阵 hess
90 hess = lambda x: np.array([[1200*x[0]**2-400*x[1]+2, -400*x[0]], [-400*x[0], 200]])
91
92 X0 = np.arange(-1.5, 1.5-0.05, 0.05)
93 X1 = np.arange(-3.5, 2+0.05, 0.05)
94 [x0, x1] = np.meshgrid(X0, X1)
95 f = 100*(x1-x0**2)**2+(1-x0)**2 # 给定的函数
96 plt.contour(x0, x1, f, 20)
97
98
99 x0, fun0, k, W = lbfgs(fun, gfun, np.array([0, 0])) # 此处x0是行向量, 计算时要转成列向量
100 print(x0, fun0, k)
101
102 plt.plot(W[0, :], W[1, :], 'g*', W[0, :], W[1, :]) # 画出迭代点轨迹
103 plt.show()
104

```

选择目标函数后的可视化结果:

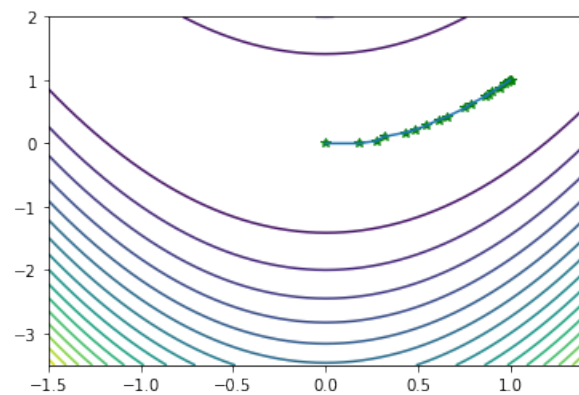


图 3: L-BFGS 算法可视化结果

3.1.4 利用 L-BFGS 方法求解逻辑回归问题

这道题不知道该咋解决了，因为之前的 L-BFGS 方法就实现得很奇怪。