

Optimización

Algoritmo BFGS

Tarea 12

20 de mayo de 2021

Erika Rivadeneira Pérez
erika.rivadeneira@cimat.mx
Matemáticas Aplicadas - CIMAT

y posteriormente se muestra la implementación del método en una función y se comparan los resultados con el método del máximo descenso.

1. Resumen

En el presente reporte se implementan los métodos del descenso por gradiente y del método Cuasi-Newton BFGS. Se comparan ambos métodos en cuanto al tiempo de cómputo, número de iteraciones y error. Se observa que el método BFGS converge mucho más rápido que el descenso por gradiente.

2. Introducción

Los métodos Cuasi-Newton son de gran importancia para la resolución de problemas de programación no lineal, dado que el costo computacional comparado con el método de Newton, el cual tiene una tasa de convergencia cuadrática, es mucho menor en problemas donde la matriz es muy grande. Dos de los métodos cuasi Newton son el método DFP y el BFGS. La fórmula de actualización DFP es bastante efectiva, pero pronto fue superada por la fórmula BFGS, la cual es actualmente considerada la más eficiente de todas las fórmulas de actualización Quasi-Newton.

En el presente trabajo se enfoca en el método BFGS ya que este método es más robusto que el DFP como se mencionó anteriormente y se lo compara con el método de máximo descenso (mencionado en reportes pasados). A continuación se menciona cómo se construye el método

3. Metodología

3.1. Algoritmo BFGS

En el BFGS la idea es calcular la aproximación de matriz inversa del Hessiano \mathbf{H}_{k+1} basado en \mathbf{B}_{k+1} . Para ello se considera la ecuación DFP, con $\rho_k = \frac{1}{s_k^T y_k}$

$$\mathbf{B}_{k+1} = (\mathbf{I} - \rho_k y_k s_k^T) \mathbf{B}_k (\mathbf{I} - \rho_k s_k y_k^T) + \rho_k y_k y_k^T$$

y las relaciones

$$\mathbf{B}_{k+1} s_k = y_k$$

y

$$\mathbf{H}_{k+1} y_k = s_k$$

Ahora $\rho_k = \frac{1}{y_k^T s_k}$ y por tanto:

$$\mathbf{H}_{k+1} = (\mathbf{I} - \rho_k s_k y_k^T) \mathbf{H}_k (\mathbf{I} - \rho_k y_k s_k^T) + \rho_k s_k s_k^T$$

Como Hessiano inicial H_0 se puede usar información específica acerca del problema, por ejemplo, configurándolo como el inverso de una aproximación del Hessiano calculado por diferencias finitas en x_0 . De otra manera, se puede simplemente configurarla para que sea la matriz identidad, o un múltiplo de la matriz identidad, donde el múltiplo es escogido para reflejar el escalamiento de las variables. En nuestro problema se escoge un múltiplo de la matriz identidad.

4. Implementación

A continuación se muestra el pseudo-código del método BFGS. Se requiere de input el vector x_0 de valores iniciales y un Hessiano inicial H_0 , el cual es un múltiplo de la matriz identidad como se mencionó en la anterior sección.

| BFGS |
|---|
| Input: x_0, H_0 |
| 1: Hacer $k = 0$ |
| 2: while $\ \nabla f_k\ \neq 0$ no converja do |
| 3: $d_k = -H_k \nabla f_k$ |
| 4: Calcular α_k usando búsqueda en línea. |
| 5: $x_{k+1} = x_k + \alpha_k d_k$ |
| 6: Calcular $\nabla f_{k+1}, \mathbf{y}_k, \mathbf{s}_k$ y actualizar $\mathbf{H}_{k+1} = (\mathbf{I} - \rho_k \mathbf{s}_k \mathbf{y}_k^T) \mathbf{H}_k (\mathbf{I} - \rho_k \mathbf{y}_k \mathbf{s}_k^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T$ |
| 7: $k = k + 1$ |
| 8: end while |
| Output: x^* |

Dados los datos, se utilizó una tolerancia de 10^{-4} para ambos métodos y se consideró un número máximo de iteraciones de 10000 para el método BFGS. Además, se consideró $m = 2$ con parámetros iniciales

$$a = [1, 1], b = 1, c = \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \end{bmatrix} \text{ y } d = [1, 1],$$

para los cuales se obtuvieron los resultados de la tabla (1)

| | Descenso por Gradiente | BFGS |
|--------------------|------------------------|---------|
| Iteraciones | 25 | 1 |
| Tiempo (s) | 899.23 | 37.81 |
| Error | 0.46484 | 0.46484 |

Tabla 1: Iteraciones, tiempo de compilación y error de los métodos de descenso por gradiente y BFGS

El error fue calculado de la siguiente manera

$$\text{error} = \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{T}} \left| \mathbf{1}_{h_{\hat{\theta}}(\mathbf{x}_i) > 0.5} (\mathbf{x}_i) - y_i \right|,$$

donde $|\mathcal{T}|$ representa al número de elementos del conjunto \mathcal{T} .

Adicionalmente, se calculó la norma del gradiente (aproximado) de la función (1) en cada método obteniendo

- Por el método de descenso por gradiente:

$$|\nabla f(\theta^*)| = 0.000561359$$

- Por el método BFGS:

$$|\nabla f(\theta^*)| = 0.0$$

5. Resultados

Se utilizaron los métodos de descenso por gradiente y BFGS para resolver el siguiente problema de optimización

$$F(\theta) = \frac{1}{N} \sum_{i=1}^N (h_{\theta}(\mathbf{x}_i) - y_i)^2 \quad (1)$$

donde $(\mathbf{x}_i, y_i), \mathbf{x}_i \in \mathbb{R}^n, y_i \in \{0, 1\}, i = 1, 2, \dots, N$ vienen dados y

$$\begin{aligned} h_{\theta}(\mathbf{x}) &= f_{\mathbf{a},b}(g_{\mathbf{c},d}(\mathbf{x})) \\ g_{\mathbf{c},d} &: \mathbb{R}^n \rightarrow \mathbb{R}^m \\ f_{\mathbf{a},b} &: \mathbb{R}^m \rightarrow \mathbb{R}, \end{aligned}$$

$m \in \mathbb{N}$ es conocida,

$$\begin{aligned} g_{\mathbf{c},d}(\mathbf{x}) &= [\sigma(\mathbf{c}_j^T \mathbf{x} + d_j)]_{j=1}^m \\ f_{\mathbf{a},b}(\mathbf{z}) &= \sigma(\mathbf{a}^T \mathbf{z} + b) \\ \sigma(t) &= \frac{1}{1 + e^{-t}}, t \in \mathbb{R} \end{aligned}$$

y θ corresponde al conjunto de parámetros $\mathbf{a}, b, \mathbf{c}, d$.

6. Conclusiones

A medida que se aumentó el valor de m el tiempo de cómputo aumenta, dado este hecho se decidió tomar $m = 2$. Se puede observar en la table (1) de resultados que el mejor método fue el de BFGS dado que su tiempo de cómputo fue mucho más bajo, al igual que sus iteraciones, que el método de descenso por gradiente. El error fue el mismo para ambos métodos. Adicionalmente, es necesario mencionar que se aproximó al gradiente de la función (1) por diferencias de Taylor hacia adelante.

Referencias

- [1] J. Nocedal and S. J. Wright. Numerical Optimization. Springer Series in Operation Research, 2000.