

基础概念回顾

那么在开始讲解前，我们先回顾一下数学上的一些基本概念。

0.1 函数的近似

假设 $y = f(x)$ 在 x_0 处可导且导数 $y' = f'(x_0) \neq 0$ ，且 $|\Delta x|$ 很小时（ Δ 表示增量），有：

$$dy = f'(x) \Delta x$$

$$\text{即 } dy \approx \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x} \Delta x$$

$$\text{移项则得 } f(x_0 + \Delta x) \approx f(x_0) + f'(x_0) \Delta x$$

$$\text{此时令 } x = x_0 + \Delta x, \text{ 则有: } f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

0.2 泰勒展开

函数 $f(x)$ 在 x_0 处展开的带有拉格朗日型余项的 n 次泰勒公式如下：

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \cdots + \frac{f^{(n)}(x_0)}{n!}(x - x_0)^n + R_n(x), \text{ 其中 } R_n(x) \text{ 为拉格朗日型余项。}$$

0.3 L1\2正则

- L1正则化是根据权重的绝对值的总和来惩罚权重。

$$\text{公式: } l_1 : \Omega(\omega) = ||\omega||_1 = \sum_i^N |w_i|$$

- L2正则化是根据权重的平方和来惩罚权重。

$$\text{公式: } l_2 : \Omega(\omega) = ||\omega||_2 = \sum_i^N w_i^2$$

不同的最优化方法

1. 梯度下降法（Gradient Descent）

梯度下降法是最早最简单，也是最为常用的最优化方法。梯度下降法实现简单，当目标函数是凸函数时，梯度下降法的解是全局解。一般情况下，其解不保证是全局最优解，梯度下降法的速度也未必是最快的。梯度下降法的优化思想是用当前位置负梯度方向作为搜索方向，因为该方向为当前位置的最快下降方向，所以也被称为是“最速下降法”。最速下降法越接近目标值，步长越小，前进越慢。

1.1 梯度（Gradient）

导数

【导数】：导数我们都非常熟悉，既可以表示某点的切线斜率，也可以表示某点变化率，公式如下表示：

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x} = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x}$$

当函数是多元的，导数就变成了偏导数：如 $f_x(x, y)$ 表示当 y 不变时， $f(x, y)$ 沿着 x 轴的变化率；那么同理， $f_y(x, y)$ 表示当 x 不变时， $f(x, y)$ 沿着 y 轴的变化率；其公式表达如下：

$$f_x(x, y) = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x, y) - f(x, y)}{\Delta x}$$

$$f_y(x, y) = \lim_{\Delta y \rightarrow 0} \frac{f(x, y + \Delta y) - f(x, y)}{\Delta y}$$

但是多元函数是一个平面，方向有很多， x 轴、 y 轴只是其中两个方向而已，假如我们需要其他方向的变化率怎么办呢？这时候方向导数就有用了，顾名思义，方向导数可以表示任意方向的导数。

方向导数

【方向导数】：

设：二次函数 $f(x, y)$ ，方向（单位向量） $u = \cos \theta i + \sin \theta j$ ，那么该二次函数在这个方向下的导数为：

$$D_u f = \lim_{t \rightarrow 0} \frac{f(x+t \cos \theta, y+t \sin \theta) - f(x, y)}{t} = f_x(x, y) \cos \theta + f_y(x, y) \sin \theta = [f_x(x, y), f_y(x, y)] \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}$$

我们记为 $D_u f = A \times l = \text{abs}(A) \text{abs}(l) \cos \alpha$ ，其中 α 是两向量的夹角。我们可以推断出当 α 等于0时，方向导数 $D_u f$ 达到最大值，此时的方向导数即为梯度。

梯度公式

【梯度的定义 ∇ 】

梯度的定义：某一函数沿着某点处的方向导数可以以最快的速度到达极大值，该方向导数我们定义为该函数的梯度。

$$\nabla = \frac{d(f(\theta))}{d\theta}$$

其中 θ 是自变量，根据泰勒展开式我们可以得出梯度下降公式：

$$\theta = \theta_0 - \eta \cdot \nabla f(\theta_0)$$

其中 η 是步长。[1]

从几何意义上来说，梯度向量就是函数变化增加最快的地方。

具体来说，对于函数 $f(x, y)$ ，在给定点 (x_0, y_0) 处，沿着梯度向量 $\left[\frac{\partial f}{\partial x_0}, \frac{\partial f}{\partial y_0} \right]$ 的方向就是该函数增加最快的地方。

换个理解方式：沿着梯度向量的方向，更加容易找到函数的最大值。那么反过来说，沿着梯度向量相反的方向（去负号），则就是更加容易找到函数的最小值。

1.2 梯度下降的直观解释

比如我们在一座大山上的某处位置，由于我们不知道怎么下山，于是决定走一步算一步，也就是在每走到一个位置的时候，求解当前位置的梯度，沿着梯度的负方向，也就是当前最陡峭的位置向下走一步，然后继续求解当前位置梯度，向这一步所在位置沿着最陡峭最易下山的位置走一步。这样一步步的走下去，一直走到觉得我们已经到了山脚。当然这样走下去，有可能我们不能走到山脚，而是到了某一个局部的山峰低处（局部最小值）。

从上面的解释可以看出，梯度下降不一定能够找到全局的最优解，有可能是一个局部最优解。当然，如果损失函数是凸函数，梯度下降法得到的解就一定是全局最优解。

1.3 梯度下降法迭代公式

了解了梯度下降的基本原理，我们来继续讲一下梯度下降法的具体实现。

设 $f(x)$ 是 R^n 上具有一阶连续偏导数的函数，我们的目标就是求解一个无约束最优化问题： $f^* = \min_{x \in R^n} f(x)$ ，其中 f^* 就是目标函数要求的极小值。

而梯度下降法是一种迭代算法。首先，我们需要选取适当的初值 $x^{(0)}$ ，然后不断迭代，更新 x 的值，进行目标函数的极小化，直到收敛。根据1.2可知，负梯度方向是使函数下降最快的方向，所以在迭代的每一步，我们需要以负梯度方向更新 x 的值，从而达到减少函数值的目的。

举个例子，在梯度下降法的第 $k+1$ 次迭代值为： $x^{(k+1)} = x^{(k)} + \lambda_k p_k$

其中 $p_k = -\nabla f(x^{(k)})$ ，是搜索方向，即负梯度方向；

其中 λ_k 是搜索步长(即stride即学习率lr)，需要满足 $f(x^{(k)} + \lambda_k p_k) = \min_{\lambda \geq 0} f(x^{(k)} + \lambda_k p_k)$

梯度下降法步骤

【输入】：目标函数 $f(x)$ ，其梯度函数 $g(x^{(k)}) = \nabla f(x^{(k)})$ ，计算精度 ϵ

【输出】： $f(x)$ 的极小点 x^*

【步骤】：

1. 取初始值 $x^{(0)} \in R^n$ ，置 $k=0$
2. 计算 $f(x^{(k)})$
3. 计算梯度 $g_k = g(x^{(k)})$
当 $|g_k| \leq \epsilon$ 时，停止迭代，且令 $x^* = x^{(k)}$
否则，令搜索方向 $p_k = -g(x^{(k)})$ 等于负梯度方向，然后求 λ_k ，使： $f(x^{(k)} + \lambda_k p_k) = \min_{\lambda \geq 0} f(x^{(k)} + \lambda_k p_k)$
4. 下一步，置 $x^{(k+1)} = x^{(k)} + \lambda_k p_k$ ，然后继续计算 $f(x^{(k+1)})$
当 $|f(x^{(k+1)}) - f(x^{(k)})| \leq \epsilon$ 或 $|x^{(k+1)} - x^{(k)}| \leq \epsilon$ 时，停止迭代，且令 $x^* = x^{(k+1)}$
否则 $k=k+1$ ，转步骤3。

具体例子

单变量示例

现有函数 $J(x) = x^2$

那么导数（单变量情况下也就是他的梯度）为 $J'(x) = 2x$

所以梯度函数就是 $g(x^{(k)}) = \nabla f(x^{(k)}) = 2x^{(k)}$

设初始点 $x^0 = 1$ ，设置搜索步长（即学习率） $\alpha = 0.4$

根据梯度下降公式，迭代过程如下：

$$x^0 = 1$$

$$x^1 = x^0 - \alpha J'(x^0) = 1 - 0.4 * 2 = 0.2$$

$$x^2 = x^1 - \alpha J'(x^1) = 0.2 - 0.4 * 0.4 = 0.04$$

$$x^3 = x^2 - \alpha J'(x^2) = 0.04 - 0.4 * 0.08 = 0.008$$

$$x^4 = x^3 - \alpha J'(x^3) = 0.008 - 0.4 * 0.016 = 0.0016$$

$$J'(x^4) = 2 * 0.0016 = 0.0032 \approx 0 \text{ 基本找到最低点}$$

多变量示例

现有函数 $J(x) = x^2 + y^2$

那么他的梯度函数为 $\nabla J(x, y) = \nabla J(\theta) = \langle 2x, 2y \rangle$

设初始点 $\theta^0 = (1, 3)$ ，设置搜索步长（即学习率） $\alpha = 0.1$

根据梯度下降公式，迭代过程如下：

$$\theta^0 = (1, 3)$$

$$\theta^1 = \theta^0 - \alpha \nabla J(\theta^0) = [0.8 \ 2.4]$$

$$\theta^2 = \theta^1 - \alpha \nabla J(\theta^1) = [0.64 \ 1.92]$$

$$\theta^3 = \theta^2 - \alpha \nabla J(\theta^2) = [0.512 \ 1.536]$$

.....

$\nabla J(\theta^{100}) \approx 0$ 此时梯度接近0，因此在 $k=100$ 时达到了局部最小值。

1.4 梯度下降算法的调优

在使用梯度下降时，需要进行调优。哪些地方需要调优呢？

1. 算法的步长选择：步长取值取决于数据样本，可以多取一些值，从大到小，分别运行算法，看看迭代效果，如果损失函数在变小，说明取值有效，否则要增大步长。步长太大，会导致迭代过快，甚至有可能错过最优解。步长太小，迭代速度太慢，很长时间算法都不能结束。所以算法的步长需要多次运行后才能得到一个较为优的值。

2. 算法参数的初始值选择：初始值不同，获得的最小值也有可能不同，因此梯度下降求得的只是局部最小（当然如果损失函数是凸函数则一定是最优解）。由于有局部最优解的风险，需要多次用不同初始值运行算法，关键损失函数的最小值，选择损失函数最小化的初值。

3. 归一化：由于样本不同特征的取值范围不一样，可能导致迭代很慢，为了减少特征取值的影响，可以对特征数据归一化。

1.5 梯度下降法的缺点

- \1. 靠近极小值时收敛速度减慢，如下图所示；
- \2. 直线搜索时可能会产生一些问题；
- \3. 寻找的是局部最优，可能会“之字形”地下降。

1.6 梯度下降法大家族（BGD，SGD，MBGD）

批量梯度下降法（Batch Gradient Descent）

批量梯度下降法，是梯度下降法最常用的形式，具体做法也就是在更新参数时使用所有的样本来进行更新。这样一来每迭代一步，都要用到训练集所有的数据，如果数据量很大，那么可想而知这种方法的迭代速度会很慢。所以，这就引入了另外一种方法，随机梯度下降。

随机梯度下降法（Stochastic Gradient Descent）

随机梯度下降（Stochastic Gradient Descent）每次迭代只用到了一个样本，在样本量很大的情况下，常见的情况是只用到了其中一部分样本数据即可迭代到最优解。因此随机梯度下降比批量梯度下降在计算量上会大大减少。SGD有一个缺点是，其噪音较BGD要多，使得SGD并不是每次迭代都向着整体最优优化方向。而且SGD因为每次都是使用一个样本进行迭代，因此最终求得的最优解往往不是全局最优解，而只是局部最优解。但是大的整体的方向是向全局最优解的，最终的结果往往是在全局最优解附近。

小批量梯度下降法（Mini-batch Gradient Descent）

小批量梯度下降法是批量梯度下降法和随机梯度下降法的折衷，也就是对于 m 个样本，我们采用 x 个样子来迭代， $1 < x < m$ 。

以多元线性回归为例子

设有线性函数 $f(x) = \omega_0 + \omega_1 x_1 + \dots + \omega_n x_n = \omega^T x$

那么损失函数就是真实值跟预测值的MSE: $J(\omega) = \sum_{i=1}^n (\omega^T x^i - y^i)^2$

那么梯度就是 $\nabla J(\omega) = \frac{\partial J(\omega)}{\partial \omega} = \sum_{i=1}^n 2 (\omega^T x^i - y^i) * x_j^i$

那么批量梯度下降法 (Batch Gradient Descent) 的更新规则就是:

$$\omega_j = \omega_j - 2\alpha \sum_{i=1}^n (\omega^T x^i - y^i) * x_j^i$$

小批量梯度下降法 (Mini-batch Gradient Descent) 的更新规则就是:

$$\omega_j = \omega_j - 2\alpha \sum_{i=1}^m (\omega^T x^i - y^i) * x_j^i$$

而随机梯度下降法 (Stochastic Gradient Descent) 的更新规则则为:

$$\omega_j = \omega_j - 2\alpha (\omega^T x^i - y^i) * x_j^i$$

根据loss图像我们可以看到, BGD的loss是一直下降的, 而SGD的loss虽然整体是下降的, 但是中途有升有降 [2]

1.7 批量梯度下降法python实现

```
1 def train(X, y, W, B, alpha, max_iters):
2     '''
3     使用了所有的样本进行梯度下降
4     x: 训练集,
5     y: 标签,
6     w: 权重向量,
7     B: bias,
8     alpha: 学习率,
9     max_iters: 最大迭代次数.
10    '''
11    dw = 0 # 权重梯度收集器
12    dB = 0 # Bias梯度的收集器
13    m = X.shape[0] # 样本数
14
15    for i in range(max_iters):
16        dw = 0 # 每次迭代重置
17        dB = 0
18        for j in range(m):
19            # 1. 迭代所有的样本
20            # 2. 计算权重和bias的梯度保存在w_grad和b_grad,
21            # 3. 通过增加w_grad和b_grad来更新dw和dB
22            w = w - alpha * (dw / m) # 更新权重
23            B = B - alpha * (dB / m) # 更新bias
24    return W, B
```

2. 牛顿法 (Newton's method)

牛顿法是一种在实数域和复数域上近似求解方程的方法。方法使用函数 $f(x)$ 的泰勒级数的前面几项来寻找方程 $f(x)=0$ 的根。牛顿法最大的特点就在于它的收敛速度很快。

2.1 牛顿法的求解方程

并不是所有的方程都有求根公式，或者求根公式很复杂，导致求解困难。利用牛顿法，可以迭代求解。原理是利用泰勒公式 [3]，在 x_0 处展开到一阶，即 $f(x) = f(x_0) + f'(x_0)(x - x_0)$ 【式2.1】

然后求解 $f(x) = 0$ ，解得 $x = x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$ 。

这是利用泰勒公式一阶展开求得的近似解，所以 x_1 并不能让 $f(x) = 0$

只能说 $f(x_1)$ 的值比 $f(x_0)$ 的值更接近于0，同理， $f(x_2)$ 的值比 $f(x_1)$ 的值更接近于0，etc...

所以迭代的想法因此而生，进而推出了公式： $x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$

通过迭代，该式在 $f(x^*) = 0$ 的时候收敛，此时我们得方程根。

2.2 牛顿法的迭代公式

同样地，上述求解方程的过程也可以看作是一个无约束最优化问题： $\min_{x \in R^n} f(x)$ ，我们要找的也就是目标函数的极小点 x^* 。

设， $f(x)$ 具有二阶连续偏导数，那么我们就可以对【式2.1】两边取导数，得：

$$f'(x) = f'(x_0) + f''(x_0)(x - x_0), \text{ 设 } \Delta x = x - x_0$$

$$\Leftrightarrow f'(x) = f'(x_0) + f''(x_0) \Delta x$$

根据微积分的性质， $f(x)$ 取最小值时，有 $f'(x)=0$ ，我们把这个性质代入上面的式子，有：

$$0 = f'(x_0) + f''(x_0) \Delta x$$

$$\Leftrightarrow \Delta x = -\frac{f'(x_0)}{f''(x_0)}$$

这样我们就得到了牛顿法的参数迭代更新公式如下：

$$x_{n+1} = x_n - \frac{f'(x_n)}{f''(x_n)} = x_n - \frac{\nabla f(x_n)}{H(f(x_n))} \quad \text{【式2.2】}$$

其中H就是 $f(x)$ 的Hessian海塞矩阵 [3]

牛顿法步骤

【输入】：目标函数 $f(x)$ ，其梯度函数 $g(x^{(k)}) = \nabla f(x^{(k)})$ ，计算精度 ϵ

【输出】： $f(x)$ 的极小点 x^*

【步骤】：

1. 取初始值 $x^{(0)} \in R^n$ ，置 $k=0$
2. 计算梯度 $g_k = g(x^{(k)})$
3. 若 $|g_k| \leq \epsilon$ ，停止迭代，求得近似解 $x^* = x^{(k)}$
4. 计算 $H_k = H(x^{(k)})$
5. 置 $x^{(k+1)} = x^{(k)} - \frac{g_k}{H_k}$
6. 置 $k=k+1$ ，转步骤2。

具体例子

用牛顿法求解 $\min f(x) = \frac{1}{2}x_1^2 + \frac{9}{2}x_2^2, x^0 = [9 \ 1]$

- 计算梯度 $\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} x_1 & 9x_2 \end{bmatrix}$
- 计算海塞矩阵 $Hf(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 9 \end{bmatrix}$
- 所以 $x^1 = x^0 - \frac{\nabla f(x^0)}{H(f(x^0))} = [9 \ 1] - [1 \ 0 \ 0 \ 9]^{-1} [9 \ 9] = [0 \ 0]$
- $|\nabla f(x^1)| = 0 \leq \varepsilon$, 所以迭代终止, 极小点为 $x^1 = [0 \ 0]$

2.3 牛顿法与梯度下降法的差异

总的来说, 梯度法和牛顿法有如下区别:

- 梯度下降法和牛顿法相比, 两者都是迭代求解, 不过梯度下降法是梯度求解, 而牛顿法是用二阶的海森矩阵的逆矩阵求解。相对而言, 使用牛顿法收敛更快(迭代更少次数)。但是每次迭代的时间比梯度下降法长。(至于为什么牛顿法收敛更快, 通俗来说梯度下降法每次只从你当前所处位置选一个坡度最大的方向走一步, 牛顿法在选择方向时, 不仅会考虑坡度是否够大, 还会考虑你走了一步之后, 坡度是否会变得更大。所以, 可以说牛顿法比梯度下降法看得更远一点, 能更快地走到最底部。)
- 牛顿法对初始值有一定要求, 在非凸优化问题中(如神经网络训练), 牛顿法很容易陷入鞍点(牛顿法步长会越来越小), 而梯度下降法则很容易逃离鞍点(因此在神经网络训练中一般使用梯度下降法, 高维空间的神经网络中存在大量鞍点)
- 梯度下降法在靠近最优值时会震荡, 因此步长调整在梯度下降法中是必要的, 具体有adagrad, adadelta, rmsprop, adam等一系列自适应学习率的方法。

2.4 牛顿法python实现

```
1 def newton(f, x, iters):
2     """
3     实现牛顿法
4     :param f: 原函数
5     :param x: 初始值
6     :param iters: 遍历的最大epoch
7     :return:
8     """
9     Hessian_T = np.linalg.inv(hessian(f, x))
10    H_G = np.matmul(Hessian_T, jacobian(f, x))
11    x_new = x - H_G
12    print("第1次迭代后的结果为:", x_new)
13    for i in range(1, iters):
14        Hessian_T = np.linalg.inv(hessian(f, x_new))
15        H_G = np.matmul(Hessian_T, jacobian(f, x_new))
16        x_new = x_new - H_G
17        print("第"+str(i+1)+"次迭代后的结果为:", x_new)
18    return x_new
```

3. 总结

最基本的最优化方法就介绍到这里，其余的还有拟牛顿法、DFP(Davidon-Fletcher-Powell)算法、BFGS(Broyden-Fletcher-Goldfard-Shano)算法等，有需求的朋友自己查。