

梯度下降算法

Gradient Descent Algorithm

1、操作系统相关环境

1) 硬件环境：

➤ 电脑

2) 软件环境：

➤ Matlab2018a(程序设计软件)

3) 操作系统(2 选 1)：

➤ Windows10

2、梯度下降

梯度下降是迭代法的一种，可以用于求解最小二乘问题(线性和非线性的问题都可以求解)。在求解机器学习算法的模型参数，即无约束优化问题时，梯度下降（Gradient Descent）是最常采用的方法之一，另一种常用的方法是最小二乘法。在求解损失函数的最小值时，可以通过梯度下降法来一步步的迭代求解，得到最小化的损失函数和模型参数值。反过来，如果我们需要求解损失函数的最大值，这时就需要用梯度上升法来迭代了。在机器学习中，基于基本的梯度下降法发展了两种梯度下降方法，分别为随机梯度下降法和批量梯度下降

法。

1) 应用场景

梯度下降法是机器学习中一种常用到的算法，但其本身不是机器学习算法，而是一种求解的最优化算法。主要解决求最小值问题，其基本思想在于不断地逼近最优点，每一步优化方向就是梯度的方向。

机器学习的本质就是“喂”给模型数据，让模型不断地去学习，而这个学习的过程就是利用梯度下降法不断去优化的过程，目前最为常见的深度神经网络便是利用梯度的反向传播，反复更新模型参数直至收敛，从而达到优化模型的目的。

3、算法原理

对于最简单线性模型，如（1）：

$$h(\theta) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_i x_i \quad (1)$$

我们假设其损失函数为（2）：

$$J(\theta) = \frac{1}{2} [h_t(x) - y]^2 \quad (2)$$

那么梯度下降的基本形式为（3）：

$$\theta_{n+1} = \theta_n - \alpha J'(\theta) \quad (3)$$

其中， α 为学习率（learning_rate）。下一步便是要将损失函数最小化，需要对函数求导：

$$J'(\theta) = \frac{\partial J(\theta)}{\partial \theta} = [h_\theta(x) - y] * h'_\theta \quad (4)$$

求解可知：

$$J'(\theta) = [h_\theta(x) - y] * x \quad (5)$$

即：

$$\theta_{n+1} = \theta_n - \alpha [h_\theta(x^{(i)}) - y^{(i)}] * x^{(i)} \quad (6)$$

常见的批量梯度下降法(Batch Gradient Descent, BGD)、随机梯度下降法(Stochastic Gradient Descent, SGD)、小批量梯度下降法

(Mini-batch Gradient Descent, MBGD)等，被广泛的应用在线性回归（非线性）、神经网络等场景。

在使用梯度下降法时，需要选择合适的步长以及迭代次数，如果选择的步长过短，不仅会使得迭代的时间过长，还有可能让模型陷入局部最小值，无法跳出去。如果设置步长过长，模型无法找到最优解，损失函数无法下降。下面对各种梯度下降法进行介绍。

1) 批量梯度下降法

使用所有样本在当前点的梯度值来对变量参数进行更新操作，其更新公式为：

$$\theta^{k+1} = \theta^k - \alpha \sum_{i=1}^m \nabla f_{\theta^k}(x^i) \quad (7)$$

批量梯度下降，顾名思义就是将所有样本一次性喂入公式进行更新操作，迭代速度快。

2) 随机梯度下降法

在更新变量参数的时候，不再像批量梯度下降一样，将所有的样本喂入公式迭代，而是随机选取一个样本的梯度值来更新参数，更新公式为：

$$\theta^{k+1} = \theta^k - \alpha \nabla f_{\theta^k}(x^i) \quad (8)$$

3) 小批量梯度下降法

小批量梯度下降法结合了 BGD 和 SGD 的特性，从原始数据中，每次选择 n 个样本来更新参数值，更新公式为：

$$\theta^{k+1} = \theta^k - \alpha \sum_{i=t}^{t+n-1} \nabla f_{\theta^k}(x^i) \quad (9)$$

4、例题

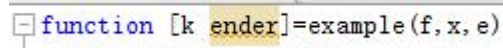
1) 拟合函数 $f = (x_1 - 2)^2 + 2 * (x_2 - 1)^2$; x_1, x_2 的初始位置为[1,3]。

2) 解题步骤

- 初始化样本集;
- 设置步长以及迭代次数;
- 根据公式 (6) 的推导更新权值;
- 计算损失函数;
- 得出结果。

3) 程序设计

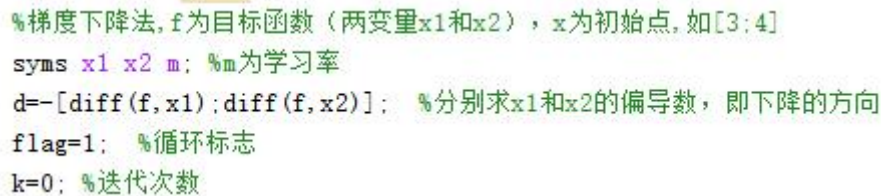
编写函数 example.m, 如图 1 所示:



```
function [k, ender]=example(f,x,e)
```

图 1.导入相关依赖
fig1.write a function

初始化样本集合, 即 x_1, x_2 的样本矩阵、权值矩阵[3,4], 函数 y, x_1, x_2 为一组等差数列, 与权值矩阵相乘后得到函数 y , 作为梯度下降的真实值, 如图 2 所示:



```
%梯度下降法,f为目标函数(两变量x1和x2),x为初始点,如[3;4]
syms x1 x2 m; %m为学习率
d=-[diff(f,x1);diff(f,x2)]; %分别求x1和x2的偏导数,即下降的方向
flag=1; %循环标志
k=0; %迭代次数
```

图 2.初始化参数
fig2.init parameter

根据公式 (7) 的推导, 撰写代码, 具体代码例程运行打开 example.m 文件, 在命令行输入:

```
syms x1 x2;
```

```
f=(x1-2)^2+2*(x2-1)^2;
```

```
x=[1;3];
```

```
e=10^(-20);
```

```
[k ender]=steepest(f,x,e)
```

然后调用函数[k ender]=example1(f,x,e)

核心程序如图 3 所示：

```
]while(flag)
    d_temp=subs(d,x1,x(1)); %将起始点代入，求得当次下降x1梯度值
    d_temp=subs(d_temp,x2,x(2)); %将起始点代入，求得当次下降x2梯度值
    nor=norm(d_temp); %范数
    if(nor>=e)
        x_temp=x+m*d_temp; %改变初始点x的值
        f_temp=subs(f,x1,x_temp(1)); %将改变后的x1和x2代入目标函数
        f_temp=subs(f_temp,x2,x_temp(2));
        h=diff(f_temp,m); %对m求导，找出最佳学习率
        m_temp=solve(h); %求方程，得到当次m
        x=x+m_temp*d_temp; %更新起始点x
        k=k+1;
    else
        flag=0;
    end
end
ender=double(x); %终点
end
```

图 3.梯度下降
fig3. gradient descent

运行结果如图 4 所示：

```
k =
    27

ender =
     2
     1
```

图 3.梯度下降迭代 27 次结果
fig3. gradient descent results of 27iterations