

作业 1:

【逻辑回归算法】

• 输入:

样本集 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}, x_i \in R^n, y_i \in \{0,1\};$

• 输出过程:

1: 初始化模型参数: $w \in R^n, b \in R;$

2: 建立逻辑回归模型: $P(x) = p(y = 1|x) = \frac{e^{w^T x + b}}{1 + e^{w^T x + b}}; P'(x) = p(y = 0|x) = \frac{1}{1 + e^{w^T x + b}}$

3: 令 $\beta = (w, b), x_i = (x_i, 1)$, 则 $P(x) = p(y = 1|x; \beta) = \frac{e^{\beta^T x}}{1 + e^{\beta^T x}}; P'(x) = p(y = 0|x; \beta) = \frac{1}{1 + e^{\beta^T x}}$

4: 计算负对数似然函数:

$$l(\beta) = \sum_{i=1}^m (-y_i \beta^T x_i + \ln(1 + e^{\beta^T x_i}))$$

5: 用某种优化算法计算参数 β : $\beta^* = \operatorname{argmin} l(\beta)$

6: 计算 w^*, b^* : $(w^*, b^*) = \beta^*$

• 输出:

$$P(x) = p(y = 1|x) = \frac{e^{w^* T x + b^*}}{1 + e^{w^* T x + b^*}}$$
$$P'(x) = p(y = 0|x) = \frac{1}{1 + e^{w^* T x + b^*}}$$

作业 3:

对于 LogisticRegression 类 的 参数的解释网上有很多, 选取其中一个博客给大家参考, 而我主要用通俗的不怎么严谨的语言来对参数进行补充说明, 特别是大家有可能困惑的地方。

<https://blog.csdn.net/jagbiam1000/article/details/79764012>

(注: 下表中出现的 目标函数 = 需要被优化的函数)

参数名称	注释	备注
penalty	用于选择正则化项。 参数值为'l1'时, 表示正则化项为 l1 正则化; 参数值为'l2'时, 表示正则化项为 l2 正则化。	新目标函数=目标函数+正则化项 正则化项是防止模型过拟合的最为常用的手段之一。
dual	选择目标函数为原始形式还是其对偶形式。	何为对偶函数? 将原始函数等价转换为一个新函数, 这个新函数我们称为对偶函数。对偶函数比原始函数更易于优化。
tol	优化算法停止的条件。	一般优化算法都是迭代算法, 举个例子, 比如牛顿法, 假设我们现在要最小化一个函数, 每迭代一次, 更新一次变量值, 函数的值都要减

		少一点, 就这样一直迭代下去。那么应该什么时候停止呢, 比如我们设 $\text{tol}=0.001$, 就是当迭代前后的函数差值 ≤ 0.001 时就停止。
C	用来控制正则化项的强弱。 C 越小, 正则化越强。	可以简单把 C 理解成正则化项的系数的倒数。
fit_intercept	用来选择逻辑回归模型中是否含有 b。	b 即线性模型的常数项。如果不含有 b, 即等价于 $b=0$
Intercept_scaling		在西瓜书算法中会有一个步骤, 令 $x = (x, 1)$, 但是在具体的代码实现上是: $x = (x, \text{intercept_scaling})$ ($(x, 1)$ 意思就是在向量 x 后加一个数值 1, 形成一个新的向量)
class_weight	设置每个类别的权重。	在西瓜书 3.6 中介绍了类别不平衡的问题, 这个参数就是为了解决这个问题的。这个权重值我们可以事先自己计算好, 然后再赋值。也可以设置 class_weight 为 balanced, 即让程序自动根据数据集计算出每个类别对应的权重值。
random_state	随机数种子。	在程序中, 有很多变量的初始值都是随机产生的值, 那么这个种子就是控制产生什么值。比如, 为种子值为 20, 那么每次随机产生的值都是 20 这个种子对应的值。而且很多时候, 数据集中每个样本的顺序需要进行打乱, 那么这个种子就是控制打乱后的顺序的。
solver	选择使用哪个优化算法进行优化。	对于一个目标函数, 我们可以有很多优化算法供我们进行选择, 因为不同的优化算法所擅长解决的问题是不同。
max_iter	优化算法的迭代次数	前面参数中介绍了, 我们可以用 tol 参数来控制优化算法是否停止。还有就是, 我们也可以用迭代次数来控制停止与否。
multi_class	用于选择多分类的策略	由西瓜书 3.5 可知, 用二分类器去构造出一个多分类器, 有很多可供选择的策略, 比如 ovo, ovr, mvm。注意, 从数学理论上讲, 我们可以构造出一个多分类函数, 但是在实践过程中, 我们并不这样做, 更一般的做法是用多个二分类器构造出一个多分类器。

verbose	主要用来控制是否 print 训练过程	
warm_start	是否热启动	什么叫热启动呢？一般而言,当我们定义一个机器学习模型时,就是定义其参数变量,所以要在一开始阶段,对变量进行随机地初始化。但是热启动的意思是,我们不在随机地初始化,而是把之前训练好的参数值赋值给变量进行初始化。
n_jobs	用 cpu 的几个核来跑程序	