

Introduction to Machine Learning

Homework 2

181860066 牛铭杨

2020 年 11 月 16 日

1 [30pts] Multi-Label Logistic Regression

In multi-label problem, each instance \mathbf{x} has a label set $\mathbf{y} = \{y_1, y_2, \dots, y_L\}$ and each label $y_i \in \{0, 1\}, \forall 1 \leq i \leq L$. Assume the post probability $p(\mathbf{y} | \mathbf{x})$ follows the conditional independence:

$$p(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^L p(y_i | \mathbf{x}). \quad (1.1)$$

Please use the logistic regression method to handle the following questions.

(1) [15pts] Please give the log-likelihood function of your logistic regression model;

给定数据集 $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^m$, 对数似然函数为

$$l(\mathbf{w}, b) = \sum_{i=1}^m \ln p(\mathbf{y}_i | \mathbf{x}_i; \mathbf{w}, b) = \sum_{j=1}^L \sum_{i=1}^m \ln p(y_{ij} | \mathbf{x}_i; \mathbf{w}, b) \quad (1.2)$$

其中 y_{ij} 是第 i 个样本在第 j 个标签上的属性

所以采用书上的记号, 只需最小化

$$l(\beta) = \sum_{i=1}^m \left(- \sum_{j=1}^L y_{ij} \beta^T \hat{\mathbf{x}}_i + L \ln(1 + e^{\beta^T \hat{\mathbf{x}}_i}) \right) \quad (1.3)$$

(2) [15pts] Please calculate the gradient of your log-likelihood function and show the parameters updating step using gradient descent.

似然函数的梯度为

$$\nabla l(\beta) = \sum_{i=1}^m \hat{\mathbf{x}}_i \left(-\sum_{j=1}^L y_{ij} + \frac{L e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right) \quad (1.4)$$

第 $t + 1$ 轮迭代解的更新公式为

$$\beta^{t+1} = \beta^t - \gamma \nabla l(\beta) = \beta^t - \gamma \sum_{i=1}^m \hat{\mathbf{x}}_i \left(-\sum_{j=1}^L y_{ij} + \frac{L e^{\beta^T \hat{\mathbf{x}}_i}}{1 + e^{\beta^T \hat{\mathbf{x}}_i}} \right) \quad (1.5)$$

2 [70pts] Logistic Regression from scratch

2.1 实现细节

我将这次实验分为 3 个阶段，导入训练数据，梯度下降/上升法训练出 $\mathbf{w} = (\mathbf{w}; b)$ ，根据测试集预测结果。

导入数据时使用两个矩阵 X, Y 来存储输入的样本和样本的标签，其中 X 每个行向量（样本）都加上一列，其值为 1，以便计算。

然后使用 *Over*，训练 26 个分类器，在训练每个分类器时，作 Y 的深拷贝，并将当前识别的那类标注为正例，其他均为反例。

使用梯度下降/上升法训练出每个分类器的 \mathbf{w} ，测试时识别为概率最高的那一类。测试时，可以统计每个类别的混淆矩阵，之后总计即可算出查全率和查准率以及 $F1$ 。

梯度下降/上升法每一步沿着梯度的方向前进一个步长的距离，经过多次迭代就可以近似最值。我们要最大化

$$l(\mathbf{w}) = \sum_{i=1}^m (y_i \mathbf{w}^T \hat{\mathbf{x}}_i - \ln(1 + e^{\mathbf{w}^T \hat{\mathbf{x}}_i})) \quad (2.1)$$

对其求偏导，得到

$$\begin{aligned} \frac{\partial l(\mathbf{w})}{\partial w_j} &= \sum_{i=1}^m (y_i x_{ij} - x_{ij} \frac{e^{\mathbf{w}^T \hat{\mathbf{x}}_i}}{1 + e^{\mathbf{w}^T \hat{\mathbf{x}}_i}}) \\ &= \sum_{i=1}^m x_{ij} (y_i - \frac{e^{\mathbf{w}^T \hat{\mathbf{x}}_i}}{1 + e^{\mathbf{w}^T \hat{\mathbf{x}}_i}}) \\ &= \sum_{i=1}^m x_{ij} (y_i - \text{Sigmoid}(\mathbf{w}^T \hat{\mathbf{x}}_i)) \end{aligned} \quad (2.2)$$

所以有

$$\nabla l(\mathbf{w}) = X^T(Y - \text{Sigmoid}(X\mathbf{w})) \quad (2.3)$$

而第 $t + 1$ 轮迭代解的更新公式为

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \gamma \nabla l(\mathbf{w}) = \mathbf{w}_t + \gamma X^T(Y - \text{Sigmoid}(X\mathbf{w})) \quad (2.4)$$

其中 γ 为步长，上面 (??) 即为我的实现参照的等式

2.2 优化与参数设置

实验的步长参数 γ 和循环次数 $loops$ 是可以调节的，我通过调节这两个参数来使实验效果更好。

一开始我设置 $\gamma = 0.001$ $loops = 1000$ ，效果并不好，准确率只能达到百分四十几，然后我意识到可能并不需要循环这么多次，就已经收敛了，步长太大每次都越过了最值点。所以我根据循环次数来调节步长，经过不断实践，我发现， $loops < 400$ 时， $\gamma = 0.001$ ， $loops < 480$ 时， $\gamma = 0.0001$ ， $loops < 500$ 时， $\gamma = 0.00001$ ，这样是较好的选择，准确率可以达到百分之六十几，而仅仅增加循环次数效果并不是很理想。

下面是实现的结果。

表 1: Performance on test set.

Performance Metric	Value (%)
accuracy	66.95
micro Precision	65.22
micro Recall	66.15
micro F_1	65.68
macro Precision	65.23
macro Recall	68.55
macro F_1	66.85