

1, 有人说当批量大小为 1 时基于随机梯度下降法 (Stochastic Gradient Descent, SGD) 的逻辑斯蒂回归 (Logistic Regression) 算法可以被看

作为“软性”的感知器算法 (PLA), 你认同这个说法吗? 请给出你的

理由。  
 答: 同意 ①采用批量为1的SGD优化逻辑斯蒂回归, 其权重更新方式与PLA  
 理由。 ②权重更新方式类似 ③PLA与逻辑斯蒂回归均适用于二分类问题 ④逻辑斯  
 蒂回归之所以被称为“软性”, 是因为它输出为0,1间的连续值, 表征属于某类别的概率

2, 在 Logistic regression 中当标签  $y=\{+1,-1\}$  时常用交叉熵作为损失函

数:  $L_{in}(\mathbf{w}) = \frac{1}{N} \sum_1^N \ln(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n))$ , 请推导出该函数的梯

度表达式解:  $\frac{\partial L_{in}(\mathbf{w})}{\partial \mathbf{w}} = \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{w}^T \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)} \cdot (-y_n \mathbf{x}_n) = \frac{1}{N} \sum_{n=1}^N \theta(-y_n \mathbf{w}^T \mathbf{x}_n) (-y_n \mathbf{x}_n)$

3, 为什么在 Logistic Regression 中不用  $L_{in}(\mathbf{w}) = (\theta(y \mathbf{w}^T \mathbf{x}) - 1)^2$  作

为损失函数, 这里假设  $\theta(\cdot)$  是 Sigmoid 函数, 标签  $y=\{+1,-1\}$ 。

解: 使用MSE作损失函数时, 其梯度

$$\frac{\partial L_{in}(\mathbf{w})}{\partial \mathbf{w}} = 2(\theta(y \mathbf{w}^T \mathbf{x}) - 1) \cdot \frac{\partial \theta(y \mathbf{w}^T \mathbf{x})}{\partial \mathbf{w}}$$

而其中  $\frac{\partial \theta(y \mathbf{w}^T \mathbf{x})}{\partial \mathbf{w}}$  项在  $y \mathbf{w}^T \mathbf{x} \rightarrow -\infty$  与  $y \mathbf{w}^T \mathbf{x} \rightarrow +\infty$  时均趋于0, 使梯  
 度极平稳, 非常不利于使用梯度下降计算