

机器学习day_3

逻辑回归 (Logistic Regression) 是机器学习的一种分类模型，虽然名字中有“回归”，但它只是与回归有一定的联系，属于算法的混淆概念，在实战中应理解其本质。

定义

$$h(w) = w_1x_1 + w_2x_2 + w_3x_3 + \dots + b$$

原理

$$g(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

源函数

回归的结果输入到sigmoid函数当中
输出结果：0.1 0.5 0.9 的一个概率值，默认为0.5为阈值

解释

假设有两个类别A、B，并且假设我们的概率阈值为属于A(1)这个类别的概率值。现在有一个样本的输入到逻辑回归输出结果为0.6，那么这个概率值0.6 > 0.5，最终输出为1(属于类别A)的结果就是=1(1类别)。那么反之，如果输出结果为0.3，那么，该样本属于类别A的概率为0.3的类别。

损失及优化

$$cost(h(x), y) = \begin{cases} -\log(h(x)) & \text{if } y=1 \\ -\log(1 - h(x)) & \text{if } y=0 \end{cases}$$
$$cost(h(x), y) = -\sum_{i=1}^n y_i \log(h(x_i)) - (1 - y_i) \log(1 - h(x_i))$$

综合完整的损失函数

同时使用梯度下降优化算法，来减少损失函数的值。这样主要新逻辑回归前面对应算法的权重参数，提升原本属于1类别的概率，降低原本属于0类别的概率。

优化

```
sklearn.linear_model.LogisticRegression(solver='lbfgs', penalty='l2', C=1.0)
```

• solver 优化求解方式 (默认采用lbfgs实现，内部使用了坐标轴下降法来迭代优化损失函数)
• lsg: 梯度下降法求解
• penalty: 正则化的种类
• C: 正则化力量

API

数据以两列数据少的列做正例

1. 导入数据
2. 读取数据并处理缺失值
3. 去除特征目标值，分割数据集
4. 进行标准化
5. 逻辑回归进行训练与预测

案例：肿瘤病人预测

```
log_fit = train_x, y_train
```

混淆矩阵

	实际正例	实际负例
预测正例	TP (True Positive)	FP (False Positive)
预测负例	FN (False Negative)	TN (True Negative)

准确率

预测结果为正例样本中真实为正例的比例 (一般不关注)

精确率

真实为正例的样本中预测结果为正例的比例 (直白的，对样本的区分能力)

召回率

所有真实类别为1的样本中，预测类别为1的比例

TPR = TP / (TP + FN)
FPR = FP / (FP + TN)
所有真实类别为0的样本中，预测类别为1的比例

ROC曲线

AUC的概率意义是随机取一对正负样本，正样本得分大于负样本的概率

- AUC的最小值为0.5，最大值为1，取值越高越好
- AUC=1，完美分类器，采用这个预测模型时，不管设定什么阈值都能得出完美预测。绝大多数预测的场合，不存在完美分类器
- 0.5 < AUC < 1，优于随机猜测。这个分类器 (模型) 假设设定阈值的值，能有预测价值。
- AUC=0.5，随机猜测一样 (例：丢硬币)，模型没有预测价值。
- AUC < 0.5，比随机猜测还差；但只要总是反预测而行，就优于随机猜测。因此不存在 AUC < 0.5 的情况。

最终AUC的范围在0.5, 1之间，并且越接近1越好

AUC指标

```
from sklearn.metrics import roc_auc_score
```

```
sklearn.metrics.roc_auc_score(y_true, y_score)
```

计算ROC曲线面积，即AUC值

y_true: 每个样本的真实类别，必须为0(反例) 1(正例)标记

y_score: 每个样本预测的概率值

API

- AUC只能用来评价二分类
- AUC非常适合评价样本不平衡中的分类器性能
- AUC会比较高预测出来的概率，而不仅仅提供分类

总结

分类算法-逻辑回归与二分类

线性回归

定义

线性回归(linear regression)是利用“回归方程(函数)”对“一个或”多个“自变量(特征值和因变量(目标值)之间”关系进行建模的一种分析方法。

特点

特点：只有一个自变量的情况称为“单变量回归”，大于一个自变量情况称为“多变量回归”。

损失函数

每个样本的真实值与预测值的差的平方求和

$$J(\theta) = (\theta_0(x_1) - y_1)^2 + (\theta_0(x_2) - y_2)^2 + \dots + (\theta_0(x_n) - y_n)^2$$
$$= \sum_{i=1}^n (\theta_0(x_i) - y_i)^2$$

这种方法也称为最小二乘法

其实就相当与数学中的求导问题

优化算法

法一：正规方程

$$w = (X^T X)^{-1} X^T y$$

缺点：当特征经过多次迭代时，求解速度太慢并且得不到结果

法二：梯度下降 (Gradient Descent)

$$w_1 := w_1 - \alpha \frac{\partial cost(w_1, w_2)}{\partial w_1}$$
$$w_0 := w_0 - \alpha \frac{\partial cost(w_1, w_2)}{\partial w_0}$$

α为学习速率，需要手动指定 (超参数)。α旁边的梯度表示方向沿着这个函数下降的方向，最后到达函数山脚的最低点，附近更新V值使用：因为V值高说明下降十分巨大的任务，能够找到较好的结果

正规方程

```
sklearn.linear_model.LinearRegression(fit_intercept=True)
```

通过正规方程优化

fit_intercept: 是否计算截距

LinearRegression.coef_: 回归系数

LinearRegression.intercept_: 截距

API

```
sklearn.linear_model.SGDRegressor(loss='squared_loss', fit_intercept=True, learning_rate='invscaling', eta=0.01)
```

SGDRegressor实现了随机梯度下降学习，它支持不同的loss函数和正则化惩罚项来拟合线性回归模型

- loss: 损失类型
- loss: 'squared_loss': 普通最小二乘法
- fit_intercept: 是否计算截距
- learning_rate: string, optional
- 学习策略
- 'constant': eta = eta0
- 'logistic': eta = 1.0 / (alpha * (1 + exp(p * x))) [default]
- 'invscaling': eta = eta0 / pow(t, power)
- power: b>0.25时又被称为

对于每一个参数都学习速率，可以使用learning_rate='constant'，并使用eta来控制学习速率。

SGDRegressor.coef_: 回归系数

SGDRegressor.intercept_: 截距

梯度下降

回归当中的数据大小不一致，会导致结果影响较大，所以需要做标准化处理 (前面K近邻算法也进行了标准化，贝叶斯和随机森林没有)。

分析

均方误差(Mean Squared Error) (MSE)的评价机制

$$MSE = \frac{1}{n} \sum_{i=1}^n (y^i - \hat{y}^i)^2$$

• y^i为预测值，y^i为真实值

回归性能评估 (重要)

```
sklearn.metrics.mean_squared_error(y_true, y_pred)
```

均方误差回归损失

y_true: 真实值

y_pred: 预测值

return: 均方误差

案例：波士顿房价

过程

1. 获取数据进行数据分割
2. 数据标准化处理
3. 使用线性回归模型进行预测
4. 调用predict方法进行预测
5. 利用均方误差来评估回归性能

正规方程 (不能解决拟合，一般不用)

梯度下降

可以试试指定学习率 (一般都很小)

小结

小数据集：正规方程

大数据集：梯度下降

岭回归

岭回归

岭回归，其实也是一种线性回归，只不过在算法建立回归方程时候，加上正则化的限制，从而达到解决拟合的效果

```
sklearn.linear_model.Ridge(alpha=1.0, fit_intercept=True, solver='choi', normalize=False)
```

具有L2正则化的线性回归

- alpha正则化力度，也叫L2
- 正则化力度越大，权重系数会越小
- 正则化力度越小，权重系数会越大
- *A取值: 0 < 1-10^4
- solver: 求解线性回归的优化方法
- *sag: 如果数据量、特征数比较大，选择该随机梯度下降优化

normalize: 数据是否进行标准化

normalize=False: 可以在X之前调用preprocessing.StandardScaler进行数据标准化

Ridge.coef_: 回归权重

Ridge.intercept_: 回归截距

交叉验证

```
sklearn.linear_model.RidgeCV(fit_intercept=True, scoring='neg_mean_squared_error')
```

正则化力度越大，权重系数会越小

正则化力度越小，权重系数会越大

使用带L2正则化的线性回归去预测

```
rd = Ridge(alpha=1.0)
```

```
rd.fit(X_train, y_train)
```

```
print("岭回归计算出的权重: ", rd.coef_)
```

```
print("岭回归计算出的截距: ", rd.intercept_)
```

再谈波士顿房价

过拟合与欠拟合

定义

一个假设在训练数据集上能够取得比其他假设更好的拟合，但是在测试数据集上却不能取得好的拟合数据，此时认为这个假设出现了过拟合的现象。(模型过于复杂)

解决方法

正则化 (尽量减小高次项特征的影响)

欠拟合

一个假设在训练数据集上不能取得更好的拟合，并且在测试数据集上也不能取得好的拟合数据，此时认为这个假设出现了欠拟合的现象。(模型过于简单)

解决方法

增加数据的特征数量

L2正则化

作用：可以使得模型中一些w的都很小，都接近于0 (不为0)，则忽略某些特征的影响

优点：减小模型参数使得模型更简单，越简单的模型则越不容易产生过拟合现象

- Ridge回归

L1正则化

作用：可以使得模型中一些w的值为直接为0，忽略这个特征的影响

- LASSO回归

原理

线性回归的损失函数用最小二乘法，等价于当预测值与真实值的误差满足正态分布时误差最大似然估计；岭回归的损失函数，是满足正态分布 (先验分布) 的最大似然估计；LASSO的损失函数，是满足正态分布 (先验分布) 的最大似然估计；LASSO的损失函数，是满足正态分布 (先验分布) 的最大似然估计。

总结

解决过拟合：决策树-剪枝、线性回归-L2正则化

解决欠拟合：增加数据、增大模型 (增大)

模型保存与加载

模型保存与加载

是什么

当训练模型计算好一个模型之后，那么如果别人需要我们提供结果预测，需要模型保存 (主要是模型训练好的参数)

API

```
from sklearn.externals import joblib
```

保存: joblib.dump(model, 'test.pkl')

加载: estimator = joblib.load('test.pkl')

无监督学习的K-means算法

无监督学习的K-means算法

是什么

1. 随机选取K个特征空间内的点作为初始的聚类中心 (K值一般是确定的)

2. 对于数据每个点计算到K个中心的距离，未知的点选择最近的一个聚类中心点作为标记位置，重新计算出每个聚类的新中心点 (迭代过程)

3. 查看看标记点距离聚类中心之后，重新计算出每个聚类的新中心点 (迭代过程)

4. 如果计算得出的新中心点与旧中心点一样，那么结束，否则重新进行第二步迭代

步骤

API

```
sklearn.cluster.KMeans(n_clusters=8, init='k-means++')
```

案例: (待更新)

$$SC_i = \frac{b_i - a_i}{\max(b_i, a_i)}$$

对于每个点 为已聚类数据中的样本，b_i 为到其它类的所有样本的距离最大值，a_i 为到本类簇的距离平均值。通过计算出的样本的相似度，相似度越高 (不患难越佳)

轮廓系数

计算所有样本的平均轮廓系数

X: 特征值

labels: 聚类类别的标记值

API

```
sklearn.metrics.silhouette_score(X, labels)
```

性能评价指标