

## 简单线性回归算法公式推导

求损失函数  $\sum_{i=1}^m (y^{(i)} - ax^{(i)} - b)^2$  的最小值

最小二乘法求解:

本质是试图找到一条直线，使得样本上的点到直线的欧式距离之和最小

$$J(a, b) = \sum_{i=1}^m (y^{(i)} - ax^{(i)} - b)^2$$

对J(a,b)求导并且令  $\frac{\partial J(a,b)}{\partial a} = 0$  ,  $\frac{\partial J(a,b)}{\partial b} = 0$

(一)

$$\frac{\partial J(a, b)}{\partial b} = \sum_{i=1}^m 2(y^{(i)} - ax^{(i)} - b)(-1) = 0$$

化简

$$\sum_{i=1}^m y^{(i)} - a \sum_{i=1}^m x^{(i)} - \sum_{i=1}^m b = 0$$

$$\sum_{i=1}^m y^{(i)} - a \sum_{i=1}^m x^{(i)} - mb = 0$$

$$\sum_{i=1}^m y^{(i)} - a \sum_{i=1}^m x^{(i)} = mb$$

两边同除 m

$$b = \bar{y} - a\bar{x}$$

(二)

$$\frac{\partial J(a, b)}{\partial a} = \sum_{i=1}^m 2(y^{(i)} - ax^{(i)} - b)(-x^{(i)}) = 0$$

化简并将 b 带入

$$\sum_{i=1}^m (y^{(i)} - ax^{(i)} - \bar{y} + a\bar{x}) x^{(i)} = 0$$

$$\sum_{i=1}^m (x^{(i)} y^{(i)} - ax^{(i)} x^{(i)} - x^{(i)} \bar{y} + a\bar{x} x^{(i)}) = 0$$

$$\sum_{i=1}^m (x^{(i)} y^{(i)} - x^{(i)} \bar{y} + a\bar{x} x^{(i)} - ax^{(i)} x^{(i)}) = 0$$

$$\sum_{i=1}^m (x^{(i)}y^{(i)} - x^{(i)}\bar{y}) = a \sum_{i=1}^m (x^{(i)}x^{(i)} - \bar{x}x^{(i)})$$

$$a = \frac{\sum_{i=1}^m (x^{(i)}y^{(i)} - x^{(i)}\bar{y})}{\sum_{i=1}^m (x^{(i)}x^{(i)} - \bar{x}x^{(i)})}$$

因为

$$\sum_{i=1}^m x^{(i)}\bar{y} = \bar{y} \sum_{i=1}^m x^{(i)} = m\bar{y} \cdot \bar{x} = \bar{x} \sum_{i=1}^m y^{(i)} = \sum_{i=1}^m x^{(i)}\bar{y} = \sum_{i=1}^m \bar{y}\bar{x}$$

所以

$$a = \frac{\sum_{i=1}^m (x^{(i)}y^{(i)} - x^{(i)}\bar{y} - \bar{x}y^{(i)} + \bar{x}\bar{y})}{\sum_{i=1}^m (x^{(i)}x^{(i)} - \bar{x}x^{(i)} - \bar{x}x^{(i)} + \bar{x}^2)}$$

$$a = \frac{\sum_{i=1}^m (x^{(i)} - \bar{x})(y^{(i)} - \bar{y})}{\sum_{i=1}^m (x^{(i)} - \bar{x})^2}$$

模型输出  $y = ax_i + b$

## 多元线性回归公式推导

$$\text{样本数据集} \begin{pmatrix} x_1^1, x_2^1, x_3^1, \dots, x_n^1, y^1 \\ x_1^2, x_2^2, x_3^2, \dots, x_n^2, y^2 \\ \vdots \\ x_1^n, x_2^n, x_3^n, \dots, x_n^n, y^n \end{pmatrix} = (X_1^{(i)}, X_2^{(i)}, X_3^{(i)}, \dots, X_n^{(i)}, Y^{(i)})$$

$$\text{求损失函数} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 \text{ 的最小值}$$

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n \quad \text{截距为} \theta_0$$

$$\hat{y}^{(i)} = \theta_0 + \theta_1 X_1^{(i)} + \theta_2 X_2^{(i)} + \theta_3 X_3^{(i)} + \dots + \theta_n X_n^{(i)}$$

找出  $\theta_0, \theta_1, \theta_2, \dots, \theta_n$ , 使得  $\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$  最小

$$\text{令 } \boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)^T$$

$$\hat{y}^{(i)} = \theta_0 X_0^{(i)} + \theta_1 X_1^{(i)} + \theta_2 X_2^{(i)} + \theta_3 X_3^{(i)} + \dots + \theta_n X_n^{(i)}, X_0^{(i)} \equiv 1$$

$$X_b = X_0^{(i)}, X_1^{(i)}, X_2^{(i)}, X_3^{(i)}, \dots, X_n^{(i)}$$

$$X_b = \begin{pmatrix} 1, x_1^1, x_2^1, x_3^1, \dots, x_n^1, y^1 \\ 1, x_1^2, x_2^2, x_3^2, \dots, x_n^2, y^2 \\ \vdots \\ 1, x_1^n, x_2^n, x_3^n, \dots, x_n^n, y^n \end{pmatrix} \quad \boldsymbol{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{pmatrix}$$

$$\sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 = (\mathbf{y} - X_b \cdot \boldsymbol{\theta})^T (\mathbf{y} - X_b \cdot \boldsymbol{\theta})$$

最终多元线性回归方程的正规方程解  $\theta = (X_b^T X_b)^{-1} X_b^T y$

模型输出  $\hat{y}^{(i)} = X_b \theta$

参数

超参数	解释	数值类型(默认值)	optional
fit_intercept	是否计算此模型的截距。如果设置为 <code>False</code> ，则不会在计算中使用截距（例如，预计数据已经居中	Boolean（ <code>True</code> ）	True/False
normalize	<code>fit_intercept</code> 设置为 <code>False</code> 时，将忽略此参数。如果为真，则回归量 <code>X</code> 将在回归之前通过减去平均值并除以 <code>12</code> 范数来归一化。如果您希望标准化，请在使用估算器 <code>sklearn.preprocessing.StandardScaler</code> 之前 <code>fit</code> 使用 <code>normalize=False</code> 。	Boolean（ <code>False</code> ）	True/False
copy_X	如果为 <code>True</code> ，则将复制 <code>X</code> ；否则，它可能会被覆盖。（不太懂）	Boolean（ <code>True</code> ）	True/False
n_jobs	确定 <code>cpu</code> 的核数， <code>-1</code> 表示使用所有处理器	Int（ <code>None</code> ）	

属性	解释	类型	Shape
coef_	回归系数（斜率）	array	( <code>n_features</code> , ) ( <code>n_targets</code> , <code>n_features</code> )
intercept_	截距	Float	

方法	解释	类型
<code>fit(X, y, sample_weight=None)</code>	训练线性模型	<code>X</code> : array-like or 稀疏矩阵, <code>y</code> : array_like, <code>sample_weight</code> : numpy array
<code>get_params(deep=True)</code>	如果为 <code>True</code> ，将返回此估计器的参数并包含作为模型的子对象。	<b>boolean</b>
<code>predict(X)</code>	使用线性模型预测	<code>X</code> : array-like or 稀疏矩阵
<code>score(X, y[, sample_weight])</code>	评分标准 $R^2 = 1 - \frac{(\sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2)/m}{(\sum_{i=1}^m (y^{(i)} - \bar{y})^2)/m}$	float
<code>set_params(**params)</code>	设置参数	