

人工智能之机器学习

回归算法

上海育创网络科技有限公司

主讲人：刘老师(GerryLiu)

课程要求

- 课上课下“九字”真言
 - 认真听，善摘录，勤思考
 - 多温故，乐实践，再发散
- 四不原则
 - 不懒散惰性，不迟到早退
 - 不请假旷课，不拖延作业
- 一点注意事项
 - 违反“四不原则”，不包就业和推荐就业

课程内容

- 线性回归算法
- 多项式回归算法
- 正则化
- Logistic回归算法
- Softmax回归算法
- 梯度下降
- 特征抽取
- 线性回归案例

什么是回归算法

- 回归算法是一种有监督算法
- 回归算法是一种比较常用的机器学习算法，用来建立“解释”变量(自变量 X)和观测值(因变量 Y)之间的关系；从机器学习的角度来讲，用于构建一个算法模型(函数)来做属性(X)与标签(Y)之间的映射关系，在算法的学习过程中，试图寻找一个函数 $h: R^d \rightarrow R$ 使得参数之间的关系拟合性最好。
- 回归算法中算法(函数)的最终结果是一个**连续**的数据值，输入值(属性值)是一个 d 维度的属性/数值向量

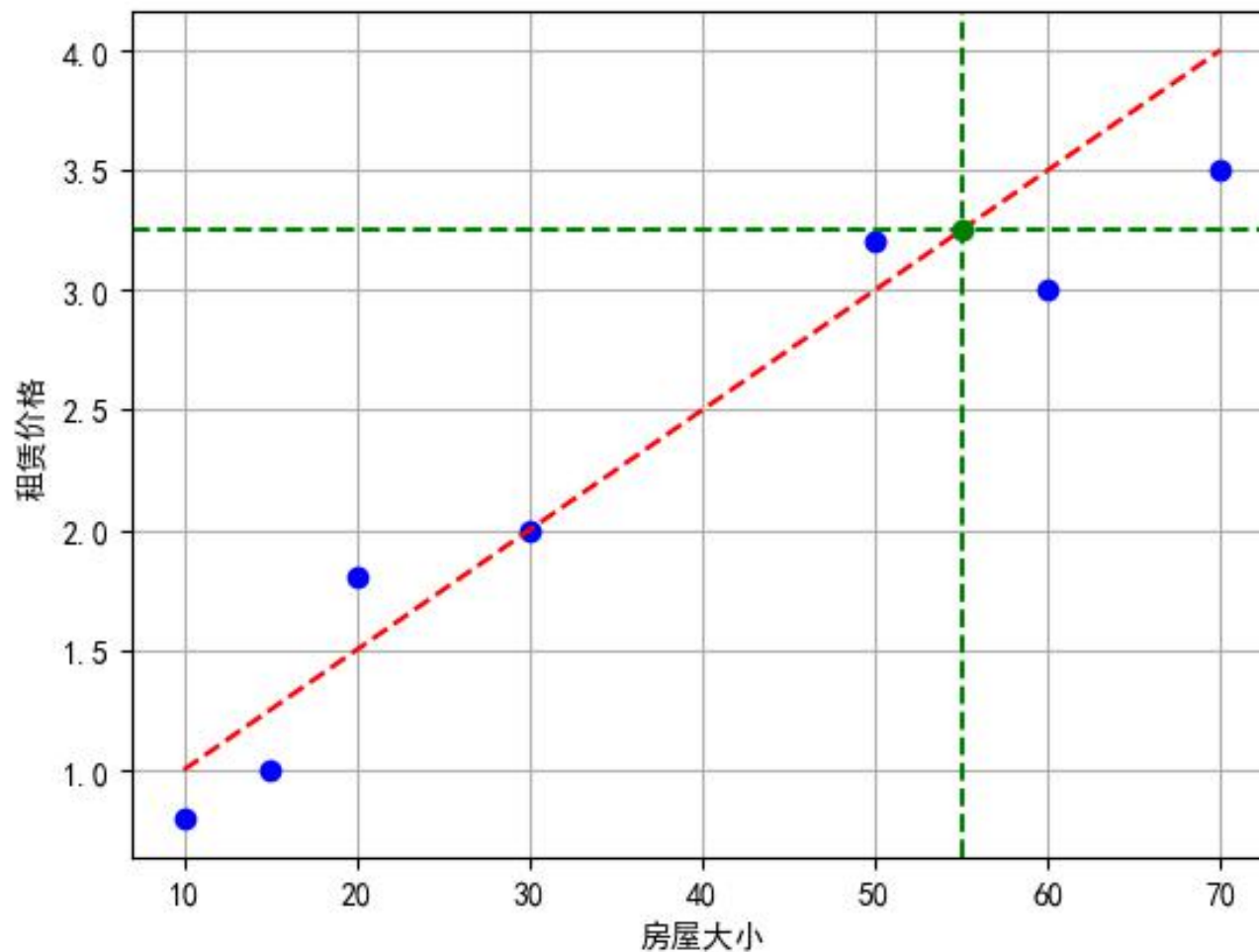
线性回归

- 认为数据中存在线性关系，也就是特征属性 X 和目标属性 Y 之间的关系是满足线性关系。
- 在线性回归算法中，找出的模型对象是期望所有训练数据比较均匀的分布在直线或者平面的两侧。
- 在线性回归中，最优模型也就是所有样本(训练数据)离模型的直线或者平面距离最小。

线性回归

- $y = ax + b$

房屋面积(m ²)	租赁价格(1000¥)
10	0.8
15	1
20	1.8
30	2
50	3.2
60	3
60	3.1
70	3.5



回归算法理性认知

- 房价的预测

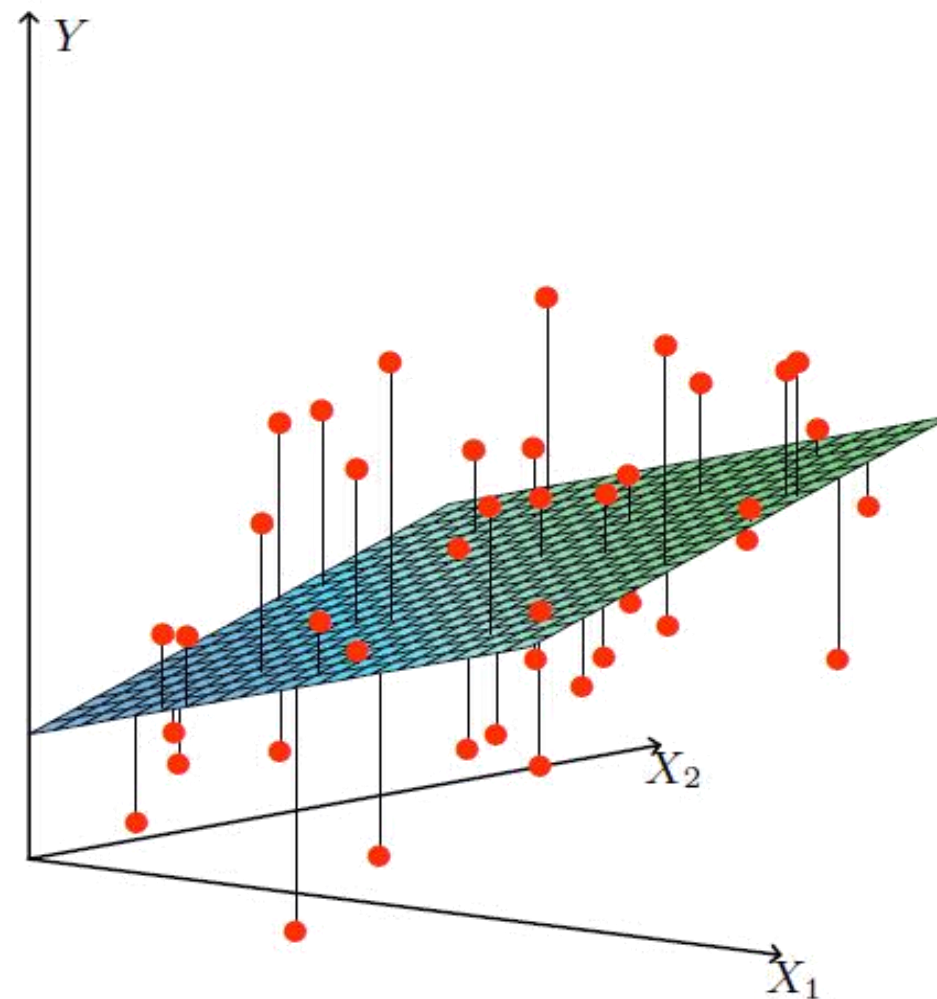
房屋面积(m ²)	租赁价格(1000¥)
10	0.8
15	1
20	1.8
30	2
50	3.2
60	3
60	3.1
70	3.5

请问，如果现在有一个房屋面积为55平，请问最终的租赁价格是多少比较合适？

线性回归

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

房屋面积	房间数量	租赁价格
10	1	0.8
20	1	1.8
30	1	2.2
30	2	2.5
70	3	5.5
70	2	5.2
.....



线性回归

$$\begin{aligned}h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n \\&= \theta_0 1 + \theta_1 x_1 + \cdots + \theta_n x_n \\&= \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n \\&= \sum_{i=0}^n \theta_i x_i = \theta^T x\end{aligned}$$

最终要求是计算出 θ 的值，并选择最优的 θ 值构成算法公式

线性回归、最大似然估计及二乘法

$$y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$$

- 误差 $\varepsilon^{(i)} (1 \leq i \leq n)$ 是独立同分布的，并且是服从均值为0，方差为某定值 σ^2 的 **高斯分布**。
 - 原因： **中心极限定理**
- 实际问题中，很多随机现象可以看做 **众多因素** 的独立影响的综合反应，往往服从正态分布

最小二乘

- 也就是说我们线性回归模型最优的时候是所有样本的预测值和实际值之间的差值最小化，由于预测值和实际值之间的差值存在正负性，所以要求平方后的值最小化。也就是可以得到如下的一个目标函数：

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (\varepsilon^{(i)})^2 = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

似然函数

$$y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)} \quad p(\varepsilon^{(i)}) = \frac{1}{\sigma\sqrt{2\pi}} e^{\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right)}$$

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$L(\theta) = \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^m \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

对数似然、目标函数及最小二乘

$$\begin{aligned}\ell(\theta) &= \ln L(\theta) \\&= \ln \prod_{i=1}^m \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\&= \sum_{i=1}^m \ln \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\&= m \ln \frac{1}{\sigma \sqrt{2\pi}} - \frac{1}{\sigma^2} \bullet \frac{1}{2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 \\loss(y_j, \hat{y}_j) &= J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2\end{aligned}$$

θ 的求解过程

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2} (X\theta - Y)^T (X\theta - Y) \rightarrow \min_{\theta} J(\theta)$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \left(\frac{1}{2} (X\theta - Y)^T (X\theta - Y) \right) = \nabla_{\theta} \left(\frac{1}{2} (\theta^T X^T - Y^T) (X\theta - Y) \right)$$

$$= \nabla_{\theta} \left(\frac{1}{2} (\theta^T X^T X\theta - \theta^T X^T Y - Y^T X\theta + Y^T Y) \right)$$

$$= \frac{1}{2} (2X^T X\theta - X^T Y - (Y^T X)^T)$$

$$= X^T X\theta - X^T Y$$

$$\theta = (X^T X)^{-1} X^T Y$$

最小二乘法的参数最优解

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

- 参数解析式

$$\theta = (X^T X)^{-1} X^T Y$$

- 最小二乘法的使用要求矩阵 $X^T X$ 是**可逆**的；为了防止不可逆或者过拟合的问题存在，可以增加额外数据影响，导致最终的矩阵是可逆的：

$$\theta = (X^T X + \lambda I)^{-1} X^T Y$$

- 最小二乘法直接求解的难点：矩阵逆的求解是一个难处

普通最小二乘法线性回归案例

- 现有一批描述家庭用电情况的数据，对数据进行算法模型预测，并最终得到预测模型（每天各个时间段和功率之间的关系、功率与电流之间的关系等）
 - 数据来源: [Individual household electric power consumption Data Set](#)
 - 建议: 使用python的sklearn库的linear_model中LinearRegression来获取算法

Individual household electric power consumption Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Measurements of electric power consumption in one household with a one-minute sampling rate over a period of almost 4 years. Different electrical quantities and some sub-metering values are available.

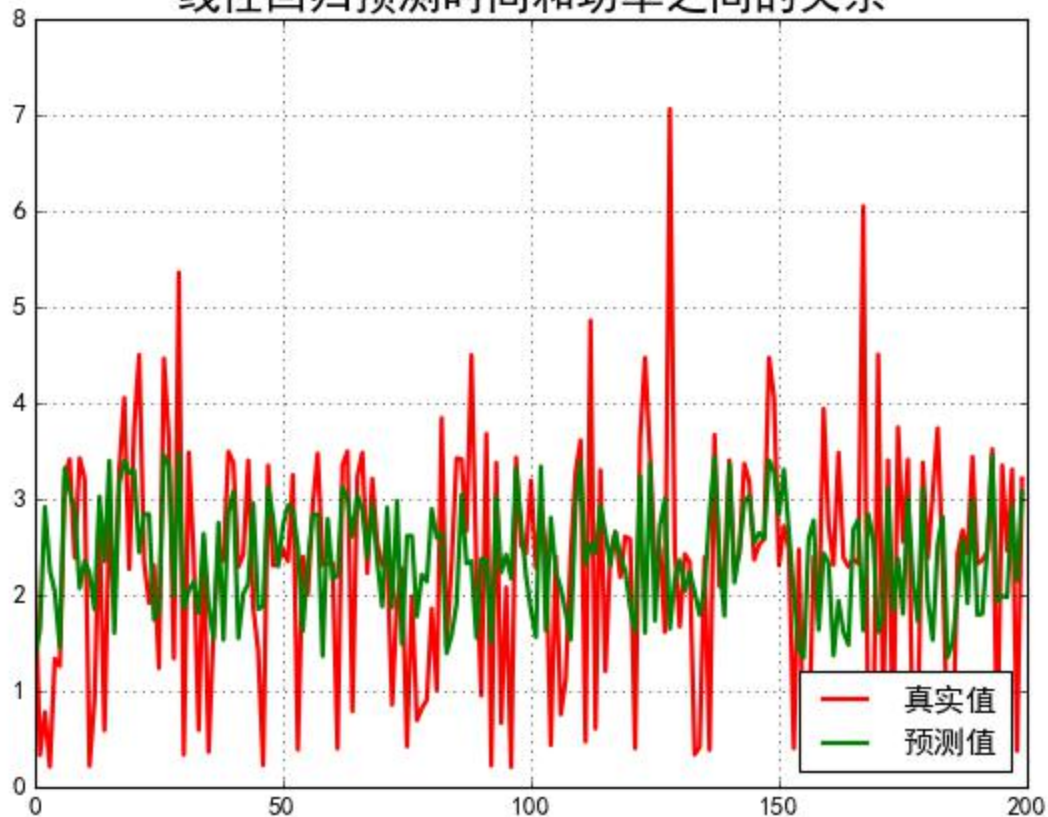
Data Set Characteristics:	Multivariate, Time-Series	Number of Instances:	2075259	Area:	Physical
Attribute Characteristics:	Real	Number of Attributes:	9	Date Donated	2012-08-30
Associated Tasks:	Regression, Clustering	Missing Values?	Yes	Number of Web Hits:	135342

Attribute Information:

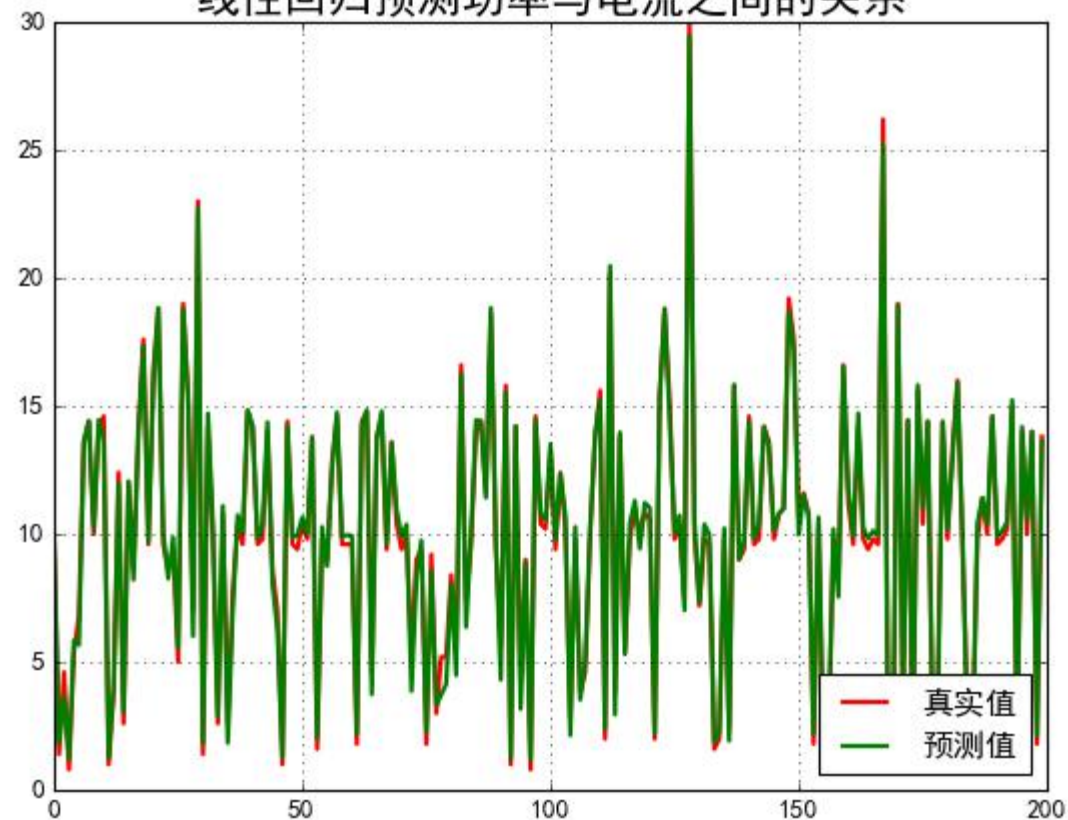
- 1.date: Date in format dd/mm/yyyy
- 2.time: time in format hh:mm:ss
- 3.global_active_power: household global minute-averaged active power (in kilowatt)
- 4.global_reactive_power: household global minute-averaged reactive power (in kilowatt)
- 5.voltage: minute-averaged voltage (in volt)
- 6.global_intensity: household global minute-averaged current intensity (in ampere)
- 7.sub_metering_1: energy sub-metering No. 1 (in watt-hour of active energy). It corresponds to the kitchen, containing mainly a dishwasher, an oven and a microwave (hot plates are not electric but gas powered).
- 8.sub_metering_2: energy sub-metering No. 2 (in watt-hour of active energy). It corresponds to the laundry room, containing a washing-machine, a tumble-drier, a refrigerator and a light.
- 9.sub_metering_3: energy sub-metering No. 3 (in watt-hour of active energy). It corresponds to an electric water-heater and an air-conditioner.

普通最小二乘法线性回归案例

线性回归预测时间和功率之间的关系



线性回归预测功率与电流之间的关系



目标函数(loss/cost function)

- 0-1损失函数 $J(\theta) = \begin{cases} 1, Y \neq f(X) \\ 0, Y = f(X) \end{cases}$
- 感知损失函数 $J(\theta) = \begin{cases} 1, |Y - f(X)| > t \\ 0, |Y - f(X)| \leq t \end{cases}$
- 平方和损失函数 $J(\theta) = \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- 绝对值损失函数 $J(\theta) = \sum_{i=1}^m |h_{\theta}(x^{(i)}) - y^{(i)}|$
- 对数损失函数 $J(\theta) = -\sum_{i=1}^m (y^{(i)} \log h_{\theta}(x^{(i)}))$

