



数据挖掘

Data Mining

主讲: 张仲楠 教授



廈門大學
XIAMEN UNIVERSITY



回归分析

目 录

01

基本概念

02

线性回归

03

非线性回归

04

逻辑回归

1. 基本概念

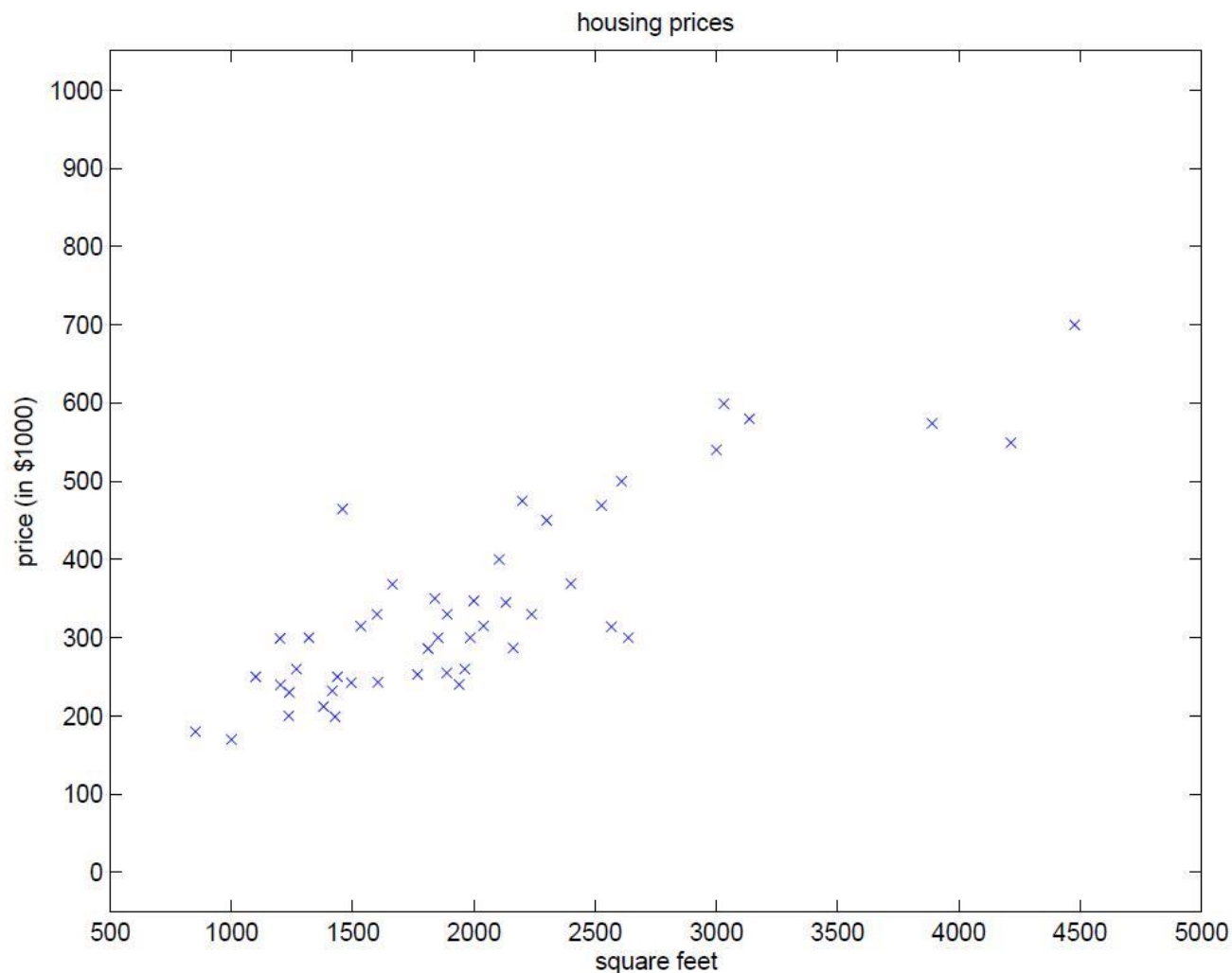
1. 基本定义

- 回归分析(regression analysis)是确定两种或两种以上变量间相互依赖的定量关系的一种统计分析方法。
- 作为一种预测模型，它基于观测数据建立变量间适当的依赖关系，以便分析数据间的内在规律，并应用于预测、控制等问题。

1. 基本概念

2. 分析案例

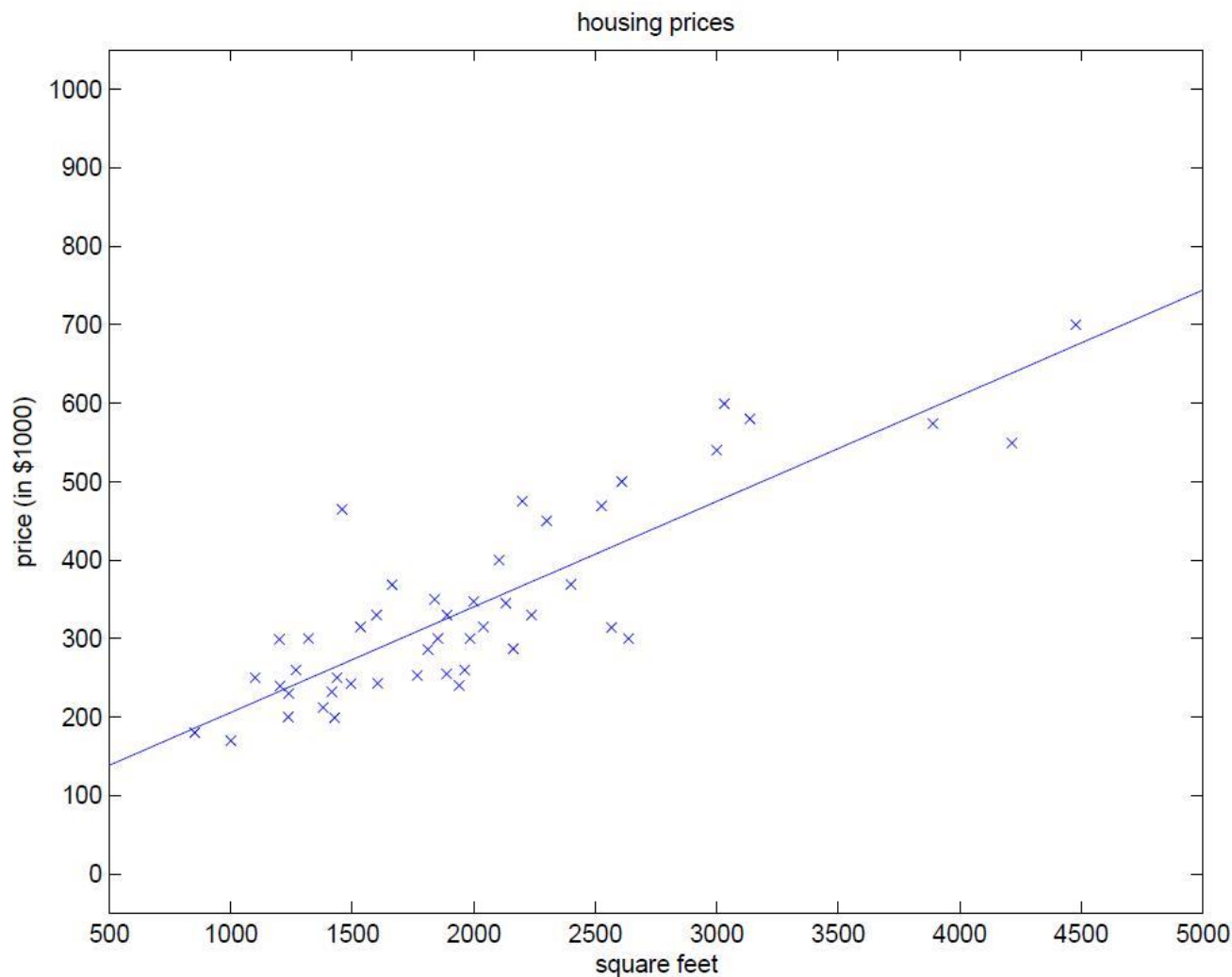
某地区的房屋面积(square feet)与价格(\$1000)的一个数据集，在该数据集中，只有一个自变量面积，和一个因变量价格



1. 基本概念

2. 分析案例

利用该数据集，我们可以训练一个**线性方程**，无限逼近所有**数据点**，然后利用该方程与给定的某一**自变量**（本例中为面积），可以**预测因变量**（本例中为房价）



1. 基本概念

3. 回归类别

- 回归与分类均为有监督学习(supervised learning)问题，其中输入 x 和输出 y 的数值是给定的，任务是学习从输入到输出的映射。
- 按照问题所涉及变量的多少，可将回归分析分为一元回归分析和多元回归分析。
- 按照自变量与因变量之间是否存在线性关系，分为线性回归分析和非线性回归分析。



01

基本概念

02

线性回归

03

非线性回归

04

逻辑回归

2. 线性回归

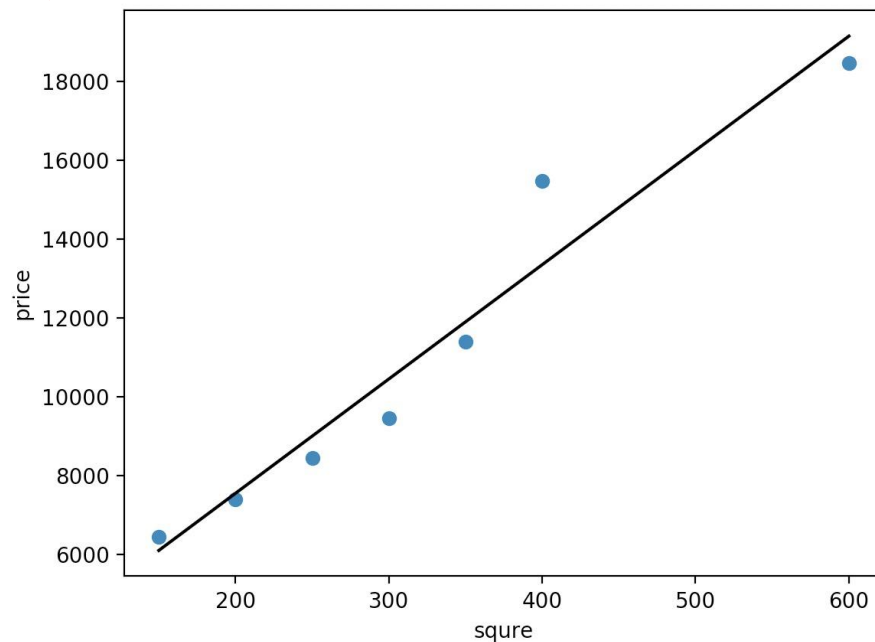
1. 基本概念

- 线性回归(Linear Regression): 因变量 (y) 和一个或多个自变量 (x) 之间建立一种线性方程关系
- 线性回归是常用的建模技术之一, 也通常是人们在学习预测模型时首选的技术之一
- 在这种技术中, 因变量是连续的, 自变量可以是连续的也可以是离散的

2. 线性回归

1. 基本概念

- 如果在某个回归分析问题中，**只有两个变量**，一个自变量和一个因变量，且自变量与因变量之间的函数关系能够**用一条直线来近似表示**，那么称其为**一元线性回归分析**。



2. 线性回归

2. 一元线性回归

- 一元线性回归包含一个自变量(x)和一个因变量(y)
- 一元线性回归方程：

$$E(y) = \beta_0 + \beta_1 x$$

- 这个方程对应的图像是一条直线，称作回归线，其中， β_0 是回归线的截距， β_1 是回归线的斜率， $E(y)$ 是在一个给定 x 值下 y 的期望值（均值）

2. 线性回归

2. 一元线性回归

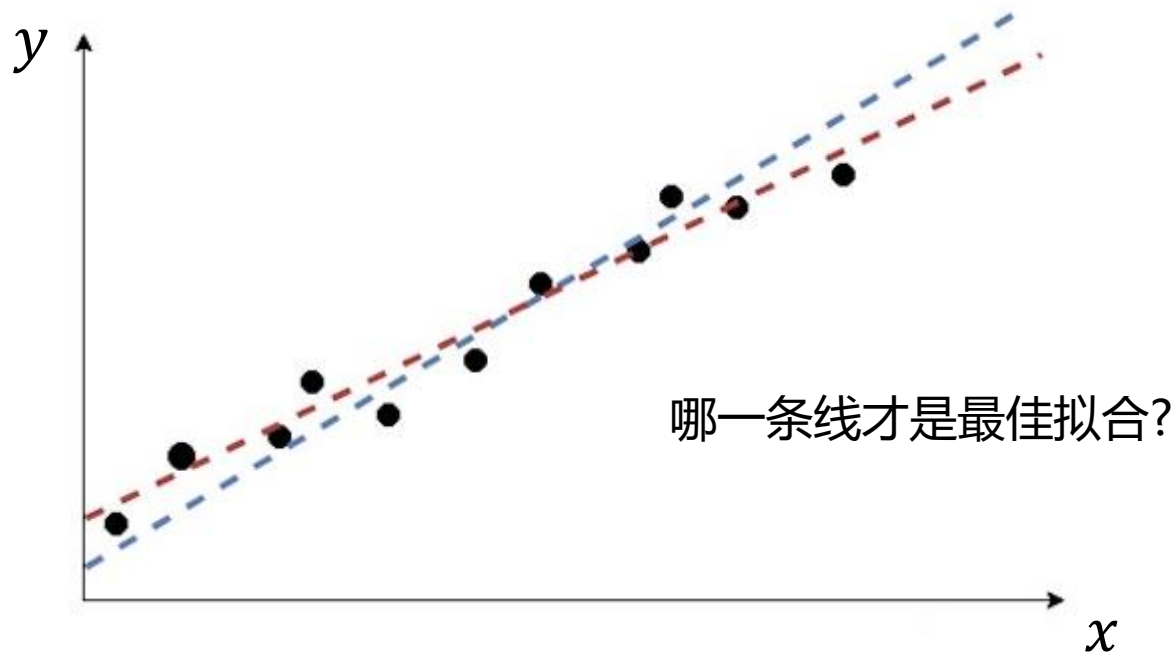
- 如何找出适合一元线性回归模型的最佳回归线?
- 一种拟合方法: **最小二乘法** (least square method)
- 最小二乘法的出发点是使实际测量数据 y_i 与拟合直线上对应的估计值 \hat{y}_i 的差(**残差**)的平方和为最小, 即:

$$\min \sum (y_i - \hat{y}_i)^2$$

2. 线性回归

2. 一元线性回归

- 法国数学家阿德里安-马里·勒让德 (1752 - 1833) 提出让总的误差的平方最小的 y 就是真值，这是基于“如果误差是随机的，应该围绕真值上下波动”



2. 线性回归

2. 一元线性回归

■ 假设我们的线性方程为： $f(x) = \beta_0 + \beta_1 x$

■ 样本的误差为：

$$e_i = y_i - f(x_i) = y_i - \beta_0 - \beta_1 x_i$$

■ 根据最小二乘法思想，总误差为：

$$Q = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

■ 通过最小化 Q 来确定直线方程，此时该问题变成了求函数极值的问题。

2. 线性回归

2. 一元线性回归

- 求关于未知参数 β_0 和 β_1 的偏导数:

$$\begin{cases} \frac{\partial Q}{\partial \beta_0} = 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)(-1) \\ \frac{\partial Q}{\partial \beta_1} = 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)(-x_i) \end{cases}$$

- 通过令偏导数为0, 可求解极值点, 即:

$$\beta_0 = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2} \quad \beta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

- 将样本 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ 代入, 即可得到参数 β_0 和 β_1 的值

2. 线性回归

2. 一元线性回归

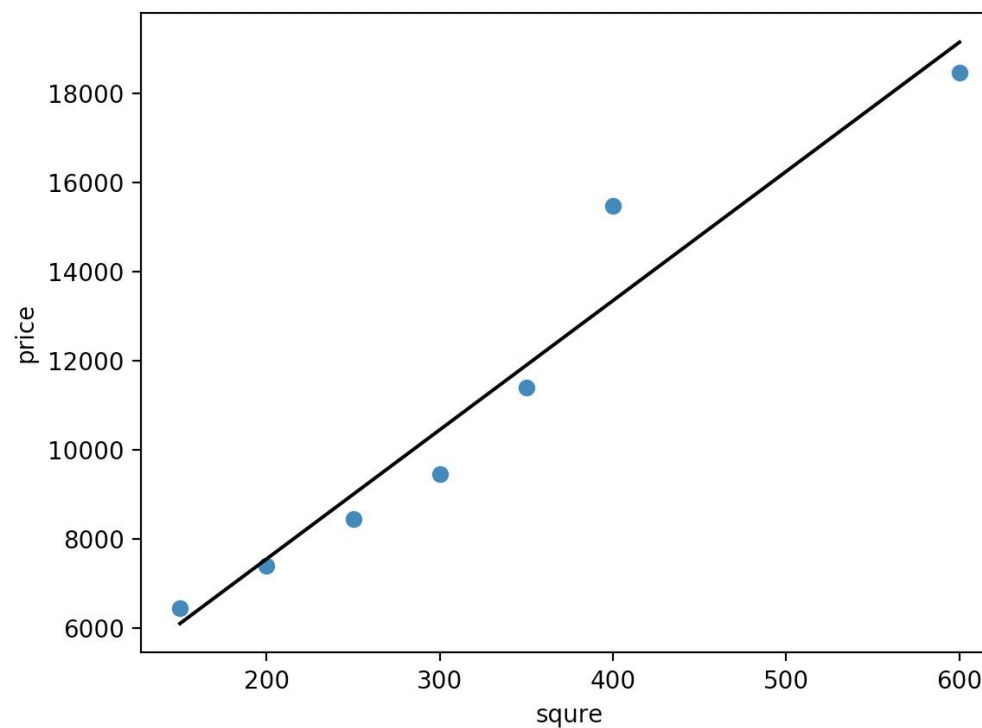
x	y
150	6450
200	7450
250	8450
300	9450
350	11450
400	15450
600	18450

$$\beta_0 = \frac{\sum x_i^2 \sum y_i - \sum x_i \sum x_i y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\beta_0 = 1771.80851064$$

$$\beta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\beta_1 = 28.77659574$$



$$f(x) = 1771.81 + 28.78x$$

2. 线性回归

2. 多元线性回归

- 给定数据集 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$, 其中 \mathbf{x}_i 由 d 个属性描述, 表示为 $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{id})$, x_{ij} 是 \mathbf{x}_i 在第 j 个属性 X_j 上的取值, $y_i \in \mathbb{R}$
- 线性模型(linear model)试图学得一个通过属性的线性组合来进行预测的函数, 即

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_d X_d$$

- $\beta_1, \beta_2, \dots, \beta_d$ 称为偏回归系数, β_i 的意义是当其他自变量 $X_j (j \neq i)$ 都固定时, 自变量 X_i 每变化一个单位而使因变量平均改变的数值

2. 线性回归

2. 多元线性回归

- 样本的误差为:

$$e_i = y_i - f(x_i) = y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \cdots - \beta_d x_{id}$$

- 用矩阵表示:

$$\begin{array}{c} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ 1 & x_{21} & \cdots & x_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nd} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_d \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \\ \mathbf{Y} \qquad \qquad \mathbf{X} \qquad \qquad \mathbf{\beta} \qquad \qquad \mathbf{e} \\ n \times 1 \qquad n \times (d+1) \qquad (d+1) \times 1 \qquad n \times 1 \end{array}$$

2. 线性回归

2. 多元线性回归

- 样本的误差为：

$$e_i = y_i - f(x_i) = y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \cdots - \beta_d x_{id}$$

- 根据最小二乘法思想，总误差为：

$$Q = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \cdots - \beta_d x_{id})^2$$

- 通过最小化 Q 来确定直线方程，求关于未知参数 $\beta_0, \beta_1, \dots, \beta_d$ 的偏导数，并令它们为0

$$\frac{\partial Q}{\partial \beta_j} = 0 (0 \leq j \leq d)$$

2. 线性回归

2. 多元线性回归

$$\begin{cases} \frac{\partial Q}{\partial \beta_0} = 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \cdots - \beta_d x_{id})(-1) = 0 & \longrightarrow \sum e_i = 0 \\ \frac{\partial Q}{\partial \beta_j} = 2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \cdots - \beta_d x_{id})(-x_{ij}) = 0 \quad (1 \leq j \leq d) & \longrightarrow \sum x_{ij} e_i = 0 \quad (1 \leq j \leq d) \end{cases}$$

■ 对上述方程整理后，可以得到：

$$\begin{cases} n\beta_0 + \beta_1 \sum x_{i1} + \beta_2 \sum x_{i2} + \cdots + \beta_d \sum x_{id} = \sum y_i \\ \beta_0 \sum x_{ij} + \beta_1 \sum x_{i1} x_{ij} + \beta_2 \sum x_{i2} x_{ij} + \cdots + \beta_d \sum x_{id} x_{ij} = \sum y_i x_{ij} \quad (1 \leq j \leq d) \end{cases}$$

■ 参数的最小二乘估计量为： $\beta = (X'X)^{-1}X'Y$



01

基本概念

02

线性回归

03

非线性回归

04

逻辑回归

3. 非线性回归

1. 基本概念

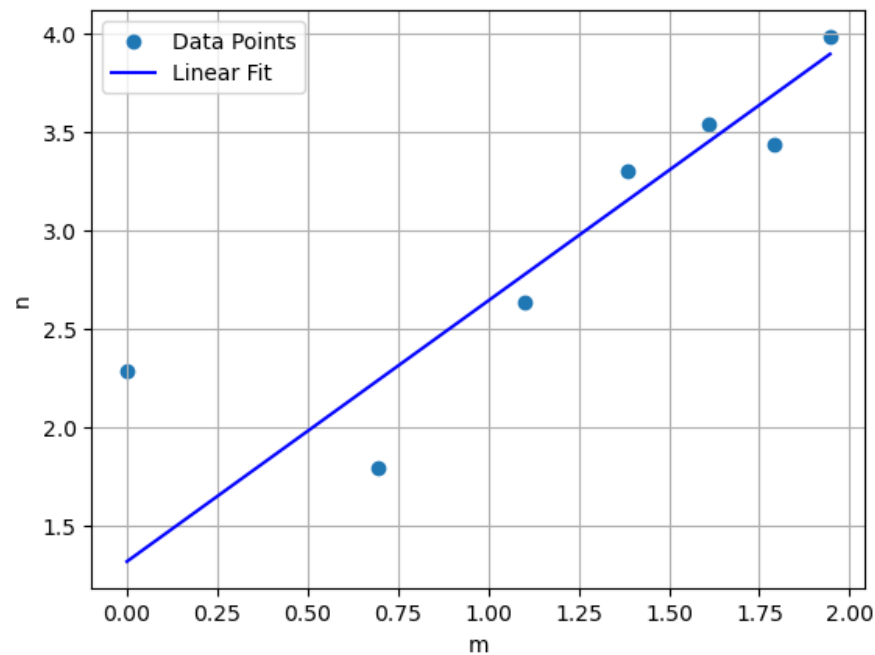
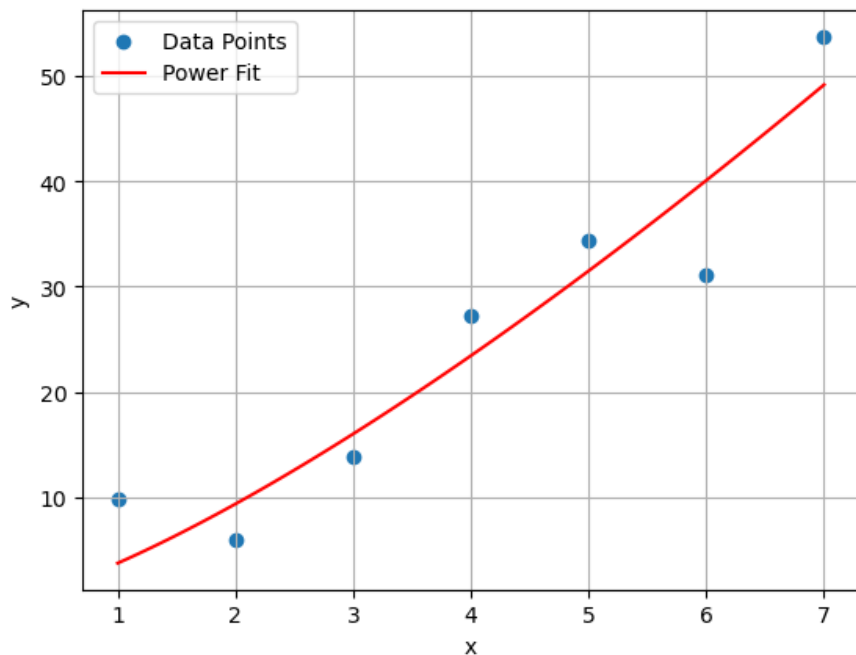
- 两个现象变量之间的相关关系并非线性关系，而呈现某种非线性的曲线关系，如：双曲线、二次曲线、三次曲线、幂函数曲线、指数函数曲线、S型曲线、对数曲线等。
- 对于非线性回归问题，常采用适当的变量代换，把问题转化为线性回归问题，求出线性回归模型后代回，得到非线性回归方程。

3. 非线性回归

x	1	2	3	4	5	6	7
y	9.82	6.0	13.894	27.204	34.338	31.114	53.750
$\ln x$	0	0.693	1.099	1.386	1.609	1.792	1.946
$\ln y$	2.284	1.792	2.631	3.303	3.536	3.438	3.984

3. 二次函数型

- 样本点分布在某幂函数曲线 $y = c_3 x^{c_4}$ 的周围, 其中 c_3, c_4 是待定参数。
- 变量代换: 令 $m = \ln y$, $n = \ln x$, 变换后样本点应该分布在直线 $m = bn + a$ 的周围, 其中 $a = \ln c_3$, $b = c_4$



$$m = 1.324n + 1.319$$

$$e^{1.319} = 3.74$$

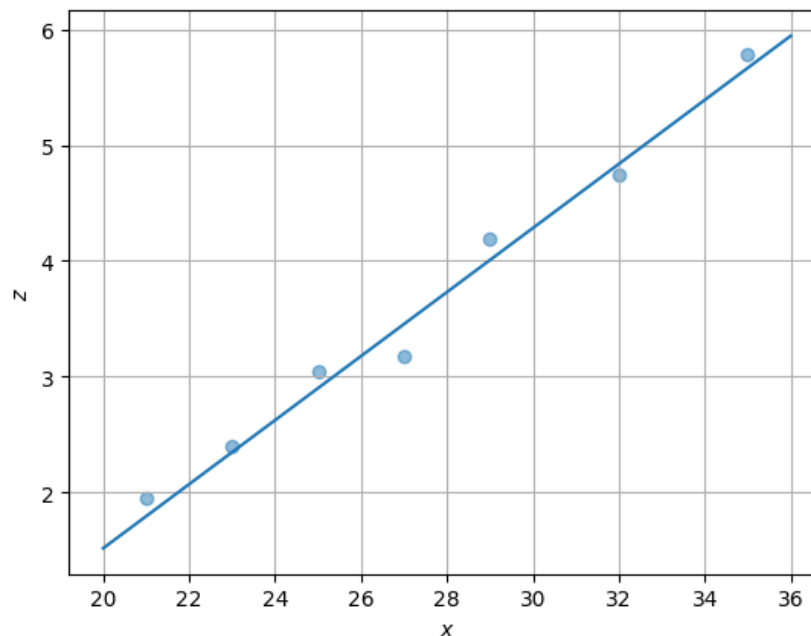
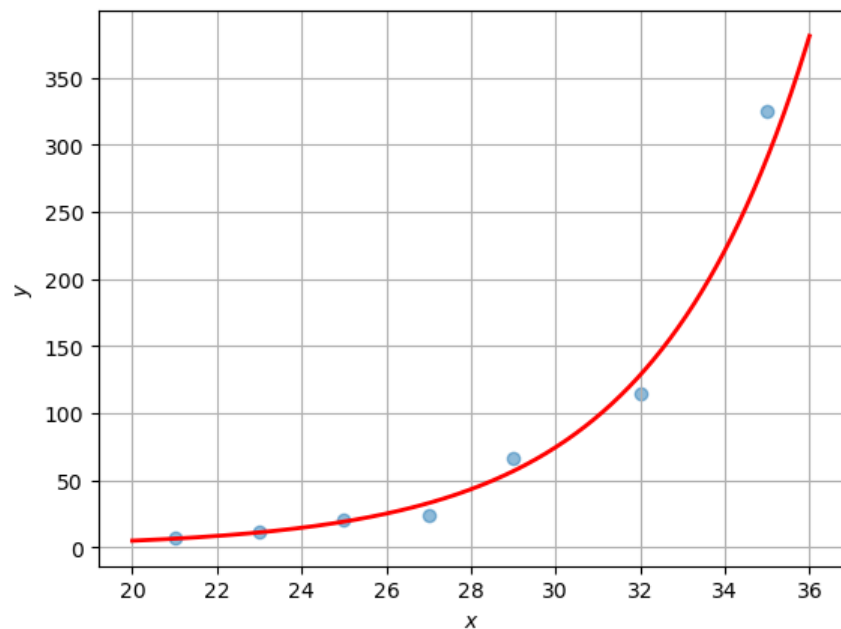
$$y = 3.74x^{1.324}$$

3. 非线性回归

x	21	23	25	27	29	32	35
y	7	11	21	24	66	115	325
z	1.946	2.398	3.045	3.178	4.190	4.745	5.784

2. 指数函数型

- 样本点分布在某一条指数函数曲线 $y = c_1 e^{c_2 x}$ 的周围, 其中 c_1, c_2 是待定参数。
- 变量代换: 令 $z = \ln y$, 变换后样本点应该分布在直线 $z = bx + a$ 的周围, 其中 $a = \ln c_1, b = c_2$



$$z = 0.272x - 3.849$$

$$y = e^{0.272x - 3.849}$$



01

基本概念

02

线性回归

03

非线性回归

04

逻辑回归

4. 逻辑回归

1. 基本概念

- logistic 回归(logistic regression): 一种概率判别模型, 它直接利用其属性值来估计数据实例 x 的概率
- logistic 回归的基本思想是使用线性预测器 $z = \mathbf{w}^T \mathbf{x} + b$ 表示 x 的概率:

$$\frac{P(y=1|\mathbf{x})}{P(y=0|\mathbf{x})} = e^z = e^{\mathbf{w}^T \mathbf{x} + b}$$

- 如果 $\mathbf{w}^T \mathbf{x} + b > 0$, 那么 x 属于第 1 类, 因为它的概率大于 1; 否则, x 属于第 0 类

4. 逻辑回归

1. 基本概念

■ 由于 $P(y = 0|\mathbf{x}) + P(y = 1|\mathbf{x}) = 1$,

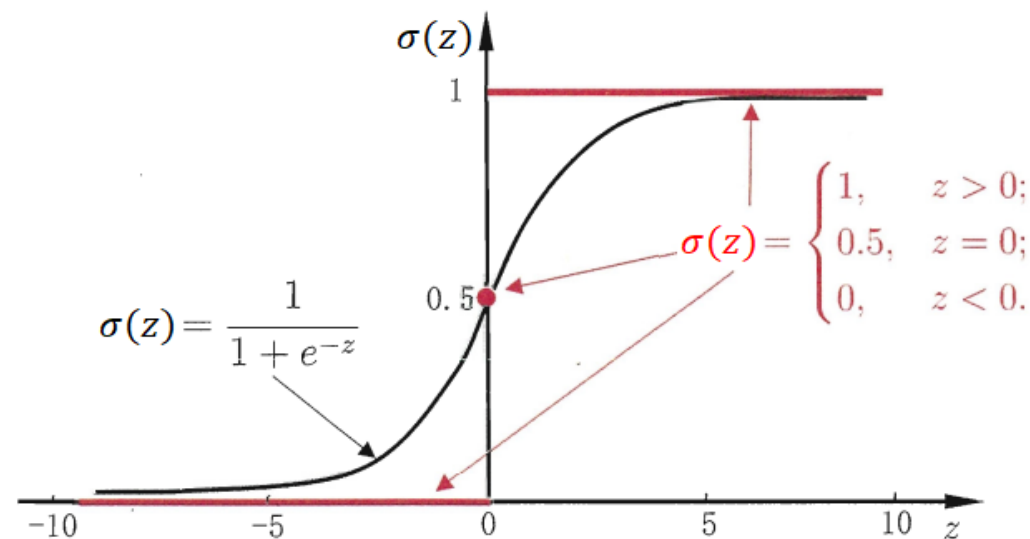
可以重写: $\frac{P(y=1|\mathbf{x})}{1-P(y=1|\mathbf{x})} = e^z$

■ 进一步简化, 将 $P(y = 1|\mathbf{x})$ 表示为 z 的函数: $P(y = 1|\mathbf{x}) = \frac{1}{1+e^{-z}} = \sigma(z)$

■ 函数 $\sigma(\cdot)$ 称为 logistic 或者 S 形函数(Sigmoid function)

■ $P(y = 0|\mathbf{x})$ 可以表示为: $P(y = 0|\mathbf{x}) = 1 - \sigma(z) = \frac{1}{1+e^z}$

■ 如果知道了参数 \mathbf{w} 和 b 的合适值, 可以用上式来估计任何数据实例 \mathbf{x} 的后验概率, 并确定其类别标签



4. 逻辑回归

2. 学习模型参数

- logistic 回归的参数 (\mathbf{w}, b) 是在训练过程中使用最大似然估计法来估计的。
这种方法需要计算观察给定 (\mathbf{w}, b) 的训练数据的可能性, 然后确定最大似然下的模型参数 (\mathbf{w}^*, b^*) 。
- 让 $D.train = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$, 表示一组 n 个训练实例, 其中 y_i 是一个二分类变量(0或1)。

4. 逻辑回归

2. 学习模型参数

- 给定 \mathbf{x}_i , \mathbf{w} 和 b , 可以表示观察到 y_i 的可能性:

$$\begin{aligned} P(y_i|\mathbf{x}_i, \mathbf{w}, b) &= P(y = 1|\mathbf{x}_i)^{y_i} \times P(y = 0|\mathbf{x}_i)^{1-y_i} \\ &= (\sigma(\mathbf{w}^T \mathbf{x}_i + b))^{y_i} \times (1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b))^{1-y_i} \end{aligned}$$

- 所有训练实例 $\mathcal{L}(\mathbf{w}, b)$ 的可能性可以通过取单个似然积(假设训练实例中的独立性)来计算:
- $\mathcal{L}(\mathbf{w}, b) = \prod_{i=1}^n P(y_i|\mathbf{x}_i, \mathbf{w}, b) = \prod_{i=1}^n P(y = 1|\mathbf{x}_i)^{y_i} \times P(y = 0|\mathbf{x}_i)^{1-y_i}$

4. 逻辑回归

2. 学习模型参数

■ 考虑似然函数的负对数(以e为底), 也称为交叉熵:

$$\begin{aligned} -\log \mathcal{L}(\mathbf{w}, b) &= -\sum_{i=1}^n y_i \log P(y=1|\mathbf{x}_i) + (1-y_i) \log P(y=0|\mathbf{x}_i) \\ &= -\sum_{i=1}^n y_i \log(\sigma(\mathbf{w}^T \mathbf{x}_i + b)) + (1-y_i) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}_i + b)) \end{aligned}$$

■ 我们想找到模型参数 (\mathbf{w}^*, b^*) 使得交叉熵 $-\log \mathcal{L}(\mathbf{w}^*, b^*)$ 最小:

$$(\mathbf{w}^*, b^*) = \operatorname{argmin}_{(\mathbf{w}, b)} -\log \mathcal{L}(\mathbf{w}, b)$$

4. 逻辑回归

3. logistic回归模型的特点

- 1) logistic回归是一种用来直接计算概率的判别模型，它不做任何关于条件概率的假设。因此，它是相当通用的，可以应用于不同的应用程序。它也可以轻松地扩展到多分类，那时，它被称为多项式logistic回归(multinomial logistic regression)。然而，它的表达能力仅限于学习线性决策边界。
- 2) 因为每个属性都有不同的权重(参数)，因此可以分析logistic回归的学习参数来理解属性和类别标签之间的关系。