# 法律声明

☐ 课程详情请咨询
  ■ 微信公众号：小象
  ■ 新浪微博：ChinaHadoop

# 回归

# 主要内容

□ 线性回归
  ■ 高斯分布
  ■ 最大似然估计MLE
  ■ 最小二乘法的本质
□ Logistic回归
  ■ 分类问题的首选算法
□ 工具
  ■ 梯度下降算法
  ■ 极大似然估计

# 股价预测

- ☐ 方法：自回归
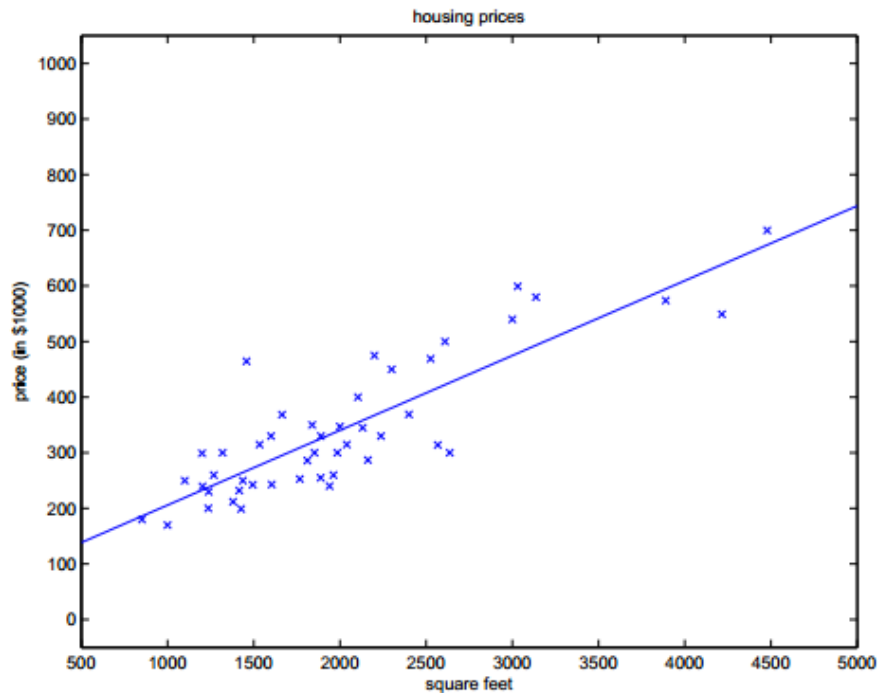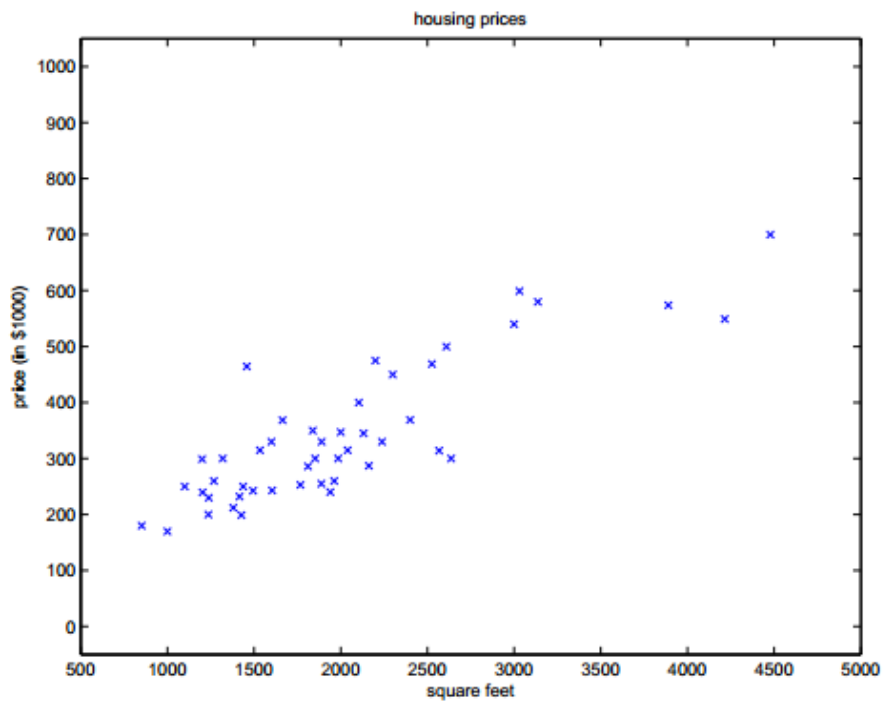- ☐ 参数：100阶



```python
if __name__ == "__main__":
    #  日期      开盘     最高     最低     收盘     成交量     成交额
    price = np.loadtxt('sh_600000.txt', delimiter='\t', skiprows=2, usecols=(4,))
    print "原始价格: \n", price
    n = 100            # 阶数
    y = price[n:]
    m = len(y)         # 样本个数
    print "预测价格: \n", y
    x = np.zeros((m, n+1))
    for i in range(m):
        x[i] = np.hstack((price[i:i+n], 1))
    print "自变量: \n", x
    theta = np.linalg.lstsq(x, y)[0]      # theta 为回归系数
    print theta
    show(theta, x, y)

    # 预测
    pn = 20        # 预测未来多少天
    x = price[-n:]
    y = np.hstack((price[-n:], np.zeros(pn)))
    for i in range(pn):
        y[n+i] = np.dot(theta, np.hstack((y[i:n+i], 1)))
    show_predict(y, pn)
```
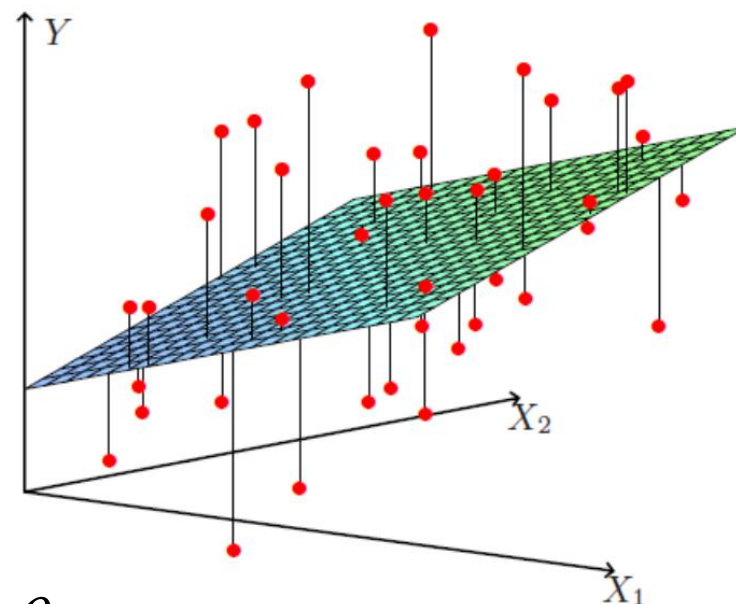
小象学院
ChinaHadoop.cn

# 线性回归

☐ y=ax+b

小象学院
ChinaHadoop.cn

# 多个变量的情形

□ **考虑两个变量**

| Living area (feet$^2$) | #bedrooms | Price (1000\$s) |
|:---:|:---:|:---:|
| 2104 | 3 | 400 |
| 1600 | 3 | 330 |
| 2400 | 3 | 369 |
| 1416 | 2 | 232 |
| 3000 | 4 | 540 |
| $\vdots$ | $\vdots$ | $\vdots$ |

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

$$h_\theta(x) = \sum_{i=0}^{n} \theta_i x_i = \theta^T x$$

小象学院
ChinaHadoop.cn

# 使用极大似然估计解释最小二乘

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

the $\epsilon^{(i)}$ are distributed IID (independently and identically distributed) according to a Gaussian distribution (also called a Normal distribution) with mean zero and some variance $\sigma^2$

theta

- 误差$\varepsilon^{(i)}(1 \leq i \leq m)$是独立同分布的，服从均值为0，方差为某定值$\sigma^2$的 高斯分布。
  - 原因：中心极限定理

小象学院
ChinaHadoop.cn

# 中心极限定理的意义

☐ 实际问题中，很多随机现象可以看做<span style="color:red">众多因素</span>的独立影响的综合反应，往往近似服从正态分布。

■ 城市耗电量：大量用户的耗电量总和

■ 测量误差：许多观察不到的、微小误差的总和

■ 注：应用前提是多个<span style="color:red">随机变量的和</span>，有些问题是乘性误差，则需要鉴别或者取对数后再使用。

小象学院
ChinaHadoop.cn

# 似然函数  $y^{(i)} = \theta^T x^{(i)} + \varepsilon^{(i)}$

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

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

$$
\begin{aligned}
L(\theta) &= \prod_{i=1}^{m} p(y^{(i)} \mid x^{(i)};\theta) \\
&= \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)
\end{aligned}
$$

互联网新技术在线教育领航者

小象学院
ChinaHadoop.cn

# 高斯的对数似然与最小二乘

$$
\begin{aligned}
\ell(\theta) &= \log L(\theta) \\
&= \log \prod_{i=1}^{m} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\
&= \sum_{i=1}^{m} \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\
&= m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^{m} (y^{(i)} - \theta^T x^{(i)})^2
\end{aligned}
$$

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

小象学院
ChinaHadoop.cn

# θ的解析式的求解过程

- □ **将M个N维样本组成矩阵X：**
  - ■ X的每一行对应一个样本，共M个样本(measurements)
  - ■ X的每一列对应样本的一个维度，共N维(regressors)
    - □ 还有额外的一维常数项，全为1

- □ 目标函数 $J(\theta) = \dfrac{1}{2}\sum_{i=1}^{m}\left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2 = \dfrac{1}{2}\left(X\theta - y\right)^T\left(X\theta - y\right)$

- □ 梯度：$\nabla_\theta J(\theta) = \nabla_\theta\left(\dfrac{1}{2}\left(X\theta - y\right)^T\left(X\theta - y\right)\right) = \nabla_\theta\left(\dfrac{1}{2}\left(\theta^T X^T - y^T\right)\left(X\theta - y\right)\right)$

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

$$= \dfrac{1}{2}\left(2X^T X\theta - X^T y - \left(y^T X\right)^T\right) = X^T X\theta - X^T y \xrightarrow{\ \text{求驻点}\ } 0$$

# 最小二乘意义下的参数最优解

☐ 参数的解析式

X m*n

n*n    n*m    m*1    n * 1    n*1

$$\theta = \left(X^T X\right)^{-1} X^T y$$

☐ 若 $\text{X}^\text{T}\text{X}$ 不可逆或防止过拟合，增加λ扰动

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

☐ "简便"方法记忆结论

$$X\theta = y \Rightarrow X^T X \theta = X^T y$$

$$\Rightarrow \theta = \left(X^T X\right)^{-1} X^T y$$

小象学院
ChinaHadoop.cn

# 加入λ扰动后

☐ XᵀX半正定：对于任意的非零向量u

$$uX^T Xu = (Xu)^T Xu \xrightarrow{\;\;令v=Xu\;\;} v^T v \geq 0$$

☐ 所以，对于任意的实数λ>0，$X^T X + \lambda I$正定，从而可逆。保证回归公式一定有意义。

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

Ridge
lambda

小象学院
ChinaHadoop.cn

# 线性回归的复杂度惩罚因子

☐ 线性回归的目标函数为：

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

☐ 将目标函数增加平方和损失：

$$J\left(\vec{\theta}\right) = \frac{1}{2}\sum_{i=1}^{m}\left(h_{\vec{\theta}}\left(x^{(i)}\right) - y^{(i)}\right)^2 + \lambda\sum_{j=1}^{n}\theta_j^2$$

☐ 本质即为假定参数θ服从高斯分布。

小象学院
ChinaHadoop.cn

# $|\mathrm{w}|$与$|\mathrm{w}|^2$

# L1-norm如何处理梯度？

- 目标函数： $J(\vec{\theta}) = \dfrac{1}{2}\sum_{i=1}^{m}\left(h_{\vec{\theta}}\left(x^{(i)}\right) - y^{(i)}\right)^2 + \lambda\sum_{j=1}^{n}\left|\theta_j\right|$

- 给定： $f(x;\alpha) = x + \dfrac{1}{\alpha}\log(1 + \exp(-\alpha x)),\ x \geq 0$

- 近似： $|x| \approx f(x;\alpha) + f(-x;\alpha) = \dfrac{1}{\alpha}\log(1 + \exp(-\alpha x) + 1 + \exp(\alpha x))$

- 梯度： $\nabla|x| \approx \dfrac{1}{1 + \exp(-\alpha x)} - \dfrac{1}{1 + \exp(\alpha x)}$

- 二阶导： $\nabla^2|x| \approx \dfrac{2\alpha\exp(\alpha x)}{(1 + \exp(\alpha x))^2}$

- 实践中，对于一般问题，如取： $\alpha = 10^6$

# 机器学习与数据使用

| 训练数据 → θ |
|:---:|

| 训练数据 → θ | 测试数据 |
|:---:|:---:|

| 训练数据 → θ | 验证数据 → λ | 测试数据 |
|:---:|:---:|:---:|

☐ 交叉验证

■ 如：十折交叉验证

# 广义逆矩阵(伪逆)

- 若A为非奇异矩阵，则线性方程组Ax=b的解为 $x = (A^T A)^{-1} A^T \cdot b$ ，从方程解的直观意义上，可以定义： $\quad A^+ = (A^T A)^{-1} A^T$

- 若A为可逆方阵， $A^+ = (A^T A)^{-1} A^T$ 即为 $A^{-1}$
$$\left(A^T A\right)^{-1} A^T = A^{-1} \left(A^T\right)^{-1} A^T = A^{-1}$$

- 当A为矩阵(非方阵)时，称$A^+$称为A的广义逆(伪逆)。

  - 奇异值分解SVD

小象学院
ChinaHadoop.cn

# 梯度下降算法 $J(\theta) = \dfrac{1}{2}\sum_{i=1}^{m}\left(h_\theta\left(x^{(i)}\right) - y^{(i)}\right)^2$

- □ 初始化θ(随机初始化)
- □ 沿着负梯度方向迭代，更新后的θ使J(θ)更小

$$\theta = \theta - \alpha \cdot \frac{\partial J(\theta)}{\partial \theta}$$

- ■ α：学习率、步长



Gradient Descent

小象学院
ChinaHadoop.cn

# 梯度方向

$$\begin{aligned}
\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} \left(h_\theta(x) - y\right)^2 \\
&= 2 \cdot \frac{1}{2} \left(h_\theta(x) - y\right) \cdot \frac{\partial}{\partial \theta_j} \left(h_\theta(x) - y\right) \\
&= \left(h_\theta(x) - y\right) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^{n} \theta_i x_i - y \right) \\
&= \left(h_\theta(x) - y\right) x_j
\end{aligned}$$

小象学院
ChinaHadoop.cn

# 批量梯度下降算法

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^{m} \left(y^{(i)} - h_\theta(x^{(i)})\right) x_j^{(i)}$$

}

theta

**gradient descent**. Note that, while gradient descent can be susceptible to local minima in general, the optimization problem we have posed here for linear regression has only one global, and no other local, optima; thus gradient descent always converges (assuming the learning rate $\alpha$ is not too large) to the global minimum. Indeed, $J$ is a convex quadratic function.

# 批量梯度下降图示

# 随机梯度下降算法

Loop {

    for i=1 to m, {

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

    }

}

This algorithm is called **stochastic gradient descent** (also **incremental gradient descent**). Whereas batch gradient descent has to scan through the entire training set before taking a single step—a costly operation if $m$ is large—stochastic gradient descent can start making progress right away, and continues to make progress with each example it looks at. Often, stochastic gradient descent gets $\theta$ "close" to the minimum much faster than batch gradient descent. (Note however that it may never "converge" to the minimum, and the parameters $\theta$ will keep oscillating around the minimum of $J(\theta)$; but in practice most of the values near the minimum will be reasonably good approximations to the true minimum.[2]) For these reasons, particularly when the training set is large, stochastic gradient descent is often preferred over batch gradient descent.

# 折中：mini-batch

☐ 如果不是每拿到一个样本即更改梯度，而是若干个样本的平均梯度作为更新方向，则是mini-batch梯度下降算法。

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^m \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

}

Loop {

  for i=1 to m, {

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

  }

}

# 回归Code

```python
def calcCoefficient(data, listA, listW, listLostFunction):
    N = len(data[0])  # 维度
    w = [0 for i in range(N)]
    wNew = [0 for i in range(N)]
    g = [0 for i in range(N)]

    times = 0
    alpha = 100.0   # 学习率随意初始化
    while times < 10000:
        j = 0
        while j < N:
            g[j] = gradient(data, w, j)
            j += 1
        normalize(g)   # 正则化梯度
        alpha = calcAlpha(w, g, alpha, data)
        numberProduct(alpha, g, wNew)

        print "times,alpha,fw,w,g:\t", times, alpha, fw(w, data), w, g
        if isSame(w, wNew):
            break
        assign2(w, wNew)   # 更新权值
        times += 1

    listA.append(alpha)
    listW.append(assign(w))
    listLostFunction.append(fw(w, data))
    return w
```
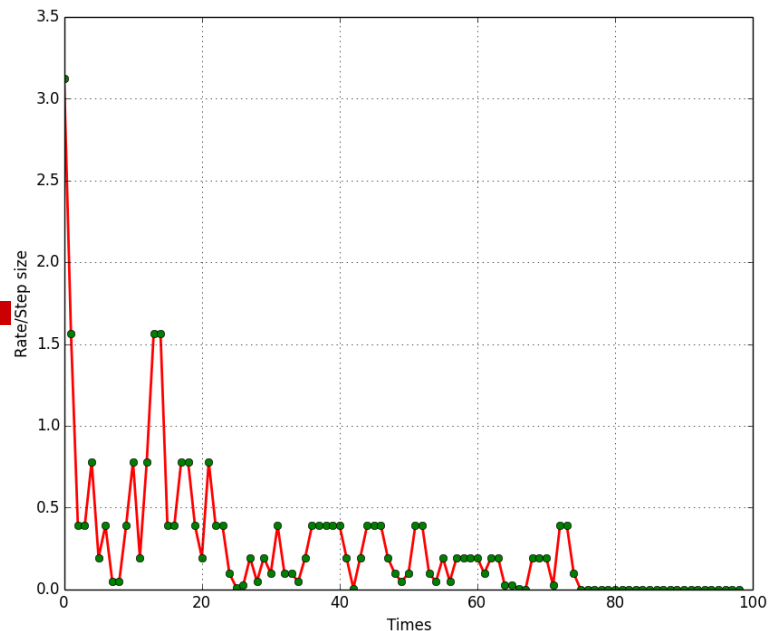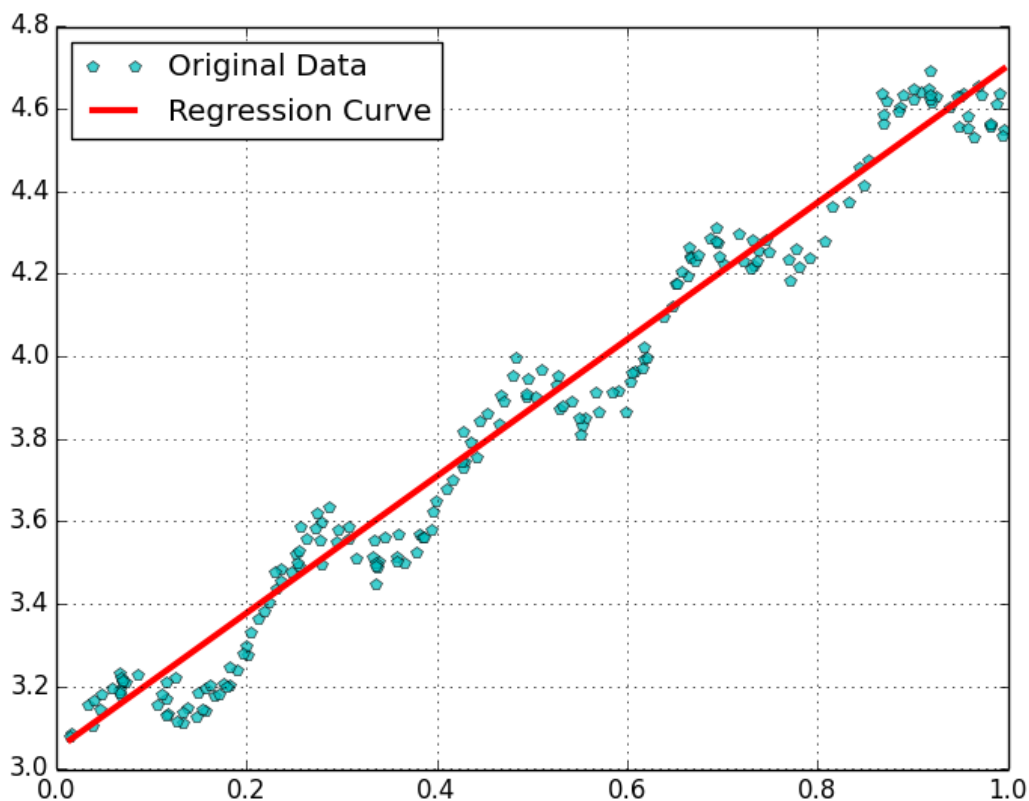
小象学院
ChinaHadoop.cn

# 附：学习率Code

```python
# w当前值；g当前梯度方向；a当前学习率；data数据
def calcAlpha(w, g, a, data):
    c1 = 0.3
    now = fw(w, data)
    wNext = assign(w)
    numberProduct(a, g, wNext)
    next = fw(wNext, data)
    # 寻找足够大的a，使得h(a)>0
    count = 30
    while next < now:
        a *= 2
        wNext = assign(w)
        numberProduct(a, g, wNext)
        next = fw(wNext, data)
        count -= 1
        if count == 0:
            break

    # 寻找合适的学习率a
    count = 50
    while next > now - c1*a*dotProduct(g, g):
        a /= 2
        wNext = assign(w)
        numberProduct(a, g, wNext)
        next = fw(wNext, data)

        count -= 1
        if count == 0:
            break
    return a
```

小象学院
ChinaHadoop.cn

# 线性回归、rate、Loss

# SGD与学习率

```python
# w当前值；g当前梯度方向；a当前学习率；data数据
def calcAlphaStochastic(w, g, a, data):
    c1 = 0.01      # 因为是每个样本都下降，所以参数运行度大些，即：激进
    now = fwStochastic(w, data)
    wNext = assign(w)
    numberProduct(a, g, wNext)
    next = fwStochastic(wNext, data)
    # 寻找足够大的a，使得h(a)>0
    count = 30
    while next < now:
        if a < 1e-10:
            a = 0.01
        else:
            a *= 2
        wNext = assign(w)
        numberProduct(a, g, wNext)
        next = fwStochastic(wNext, data)
        count -= 1
        if count == 0:
            break

    # 寻找合适的学习率a
    count = 50
    while next > now - c1*a*dotProduct(g, g):
        a /= 2
        wNext = assign(w)
        numberProduct(a, g, wNext)
        next = fwStochastic(wNext, data)

        count -= 1
        if count == 0:
            break
    return a
```

```python
def calcCoefficient(data, listA, listW, listLostFunction):
    M = len(data)          # 样本数目
    N = len(data[0])       # 维度
    w = [0 for i in range(N)]
    wNew = [0 for i in range(N)]
    g = [0 for i in range(N)]

    times = 0
    alpha = 100.0  # 学习率随意初始化
    same = False
    while times < 10000:
        i = 0
        while i < M:
            j = 0
            while j < N:
                g[j] = gradientStochastic(data[i], w, j)
                j += 1
            normalize(g)  # 正则化梯度
            alpha = calcAlphaStochastic(w, g, alpha, data[i])
            #alpha = 0.01
            numberProduct(alpha, g, wNew)

            print "times,i, alpha,fw,w,g:\t", times, i, alpha, fw(w, data), w, g
            if isSame(w, wNew):
                if times > 5:   #防止训练次数过少
                    same = True
                    break
            assign2(w, wNew)   # 更新权值

            listA.append(alpha)
            listW.append(assign(w))
            listLostFunction.append(fw(w, data))

            i += 1
        if same:
            break
        times += 1
    return w
```

小象学院
ChinaHadoop.cn

# 随机梯度下降SGD

# 批量与随机梯度下降
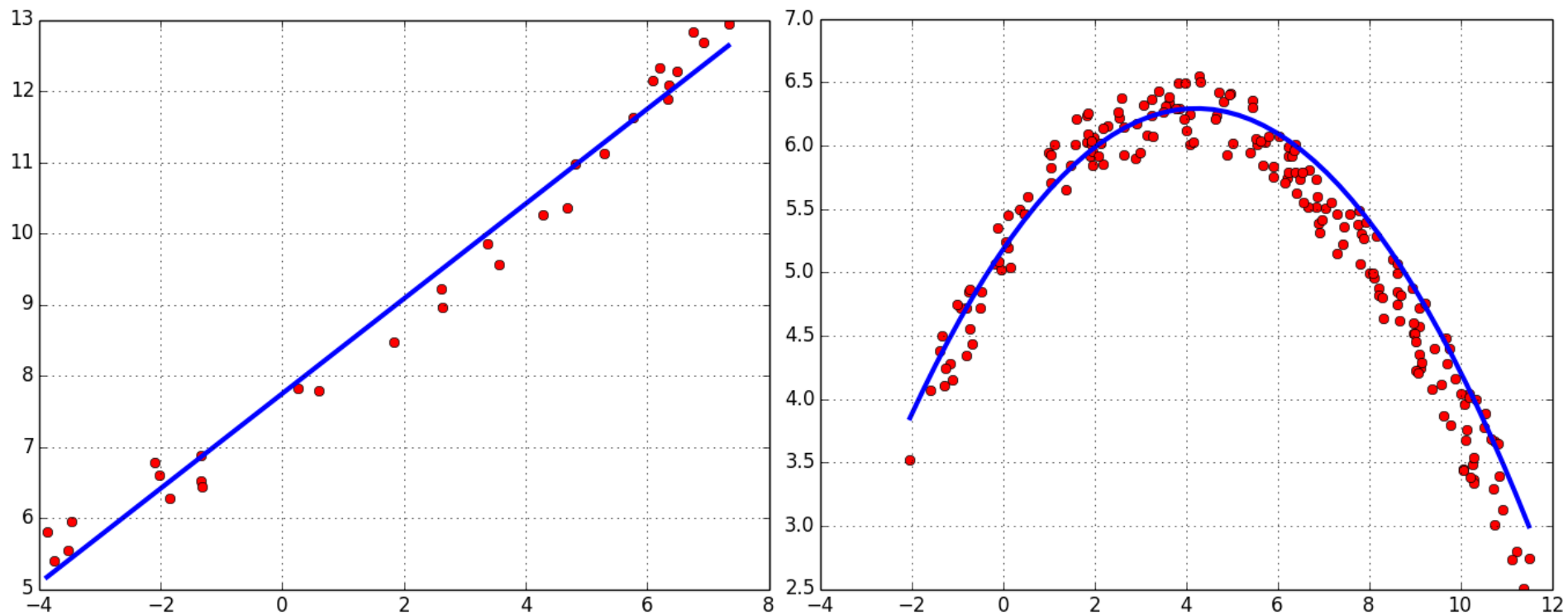
# 线性回归的进一步分析

☐ 可以对样本是非线性的，只要对参数θ线性



$$y = \theta_0 + \theta_1 x + \theta_2 x^2$$

# Code

```python
def regression(data, alpha, lamda):
    n = len(data[0]) - 1
    theta = np.zeros(n)
    for times in range(100):
        for d in data:
            x = d[:-1]
            y = d[-1]
            g = np.dot(theta, x) - y
            theta = theta - alpha * g * x + lamda * theta
        print times, theta
    return theta
```
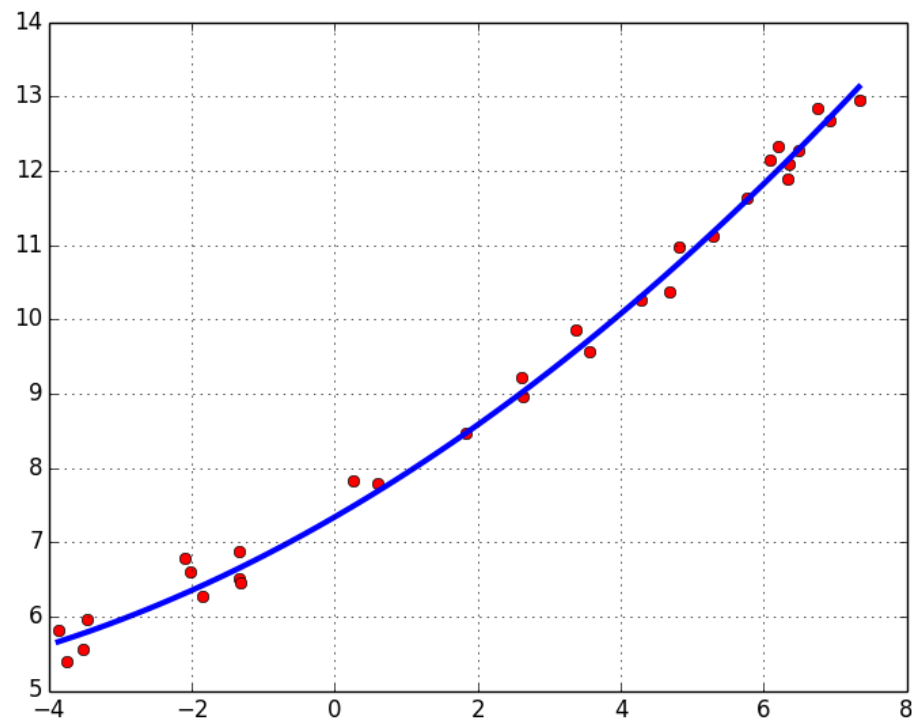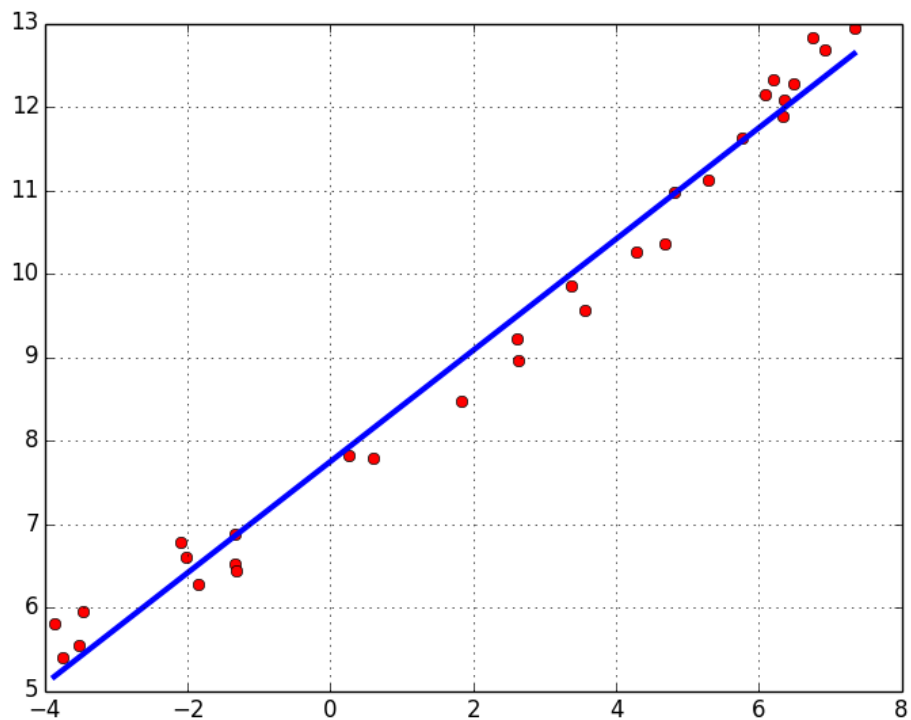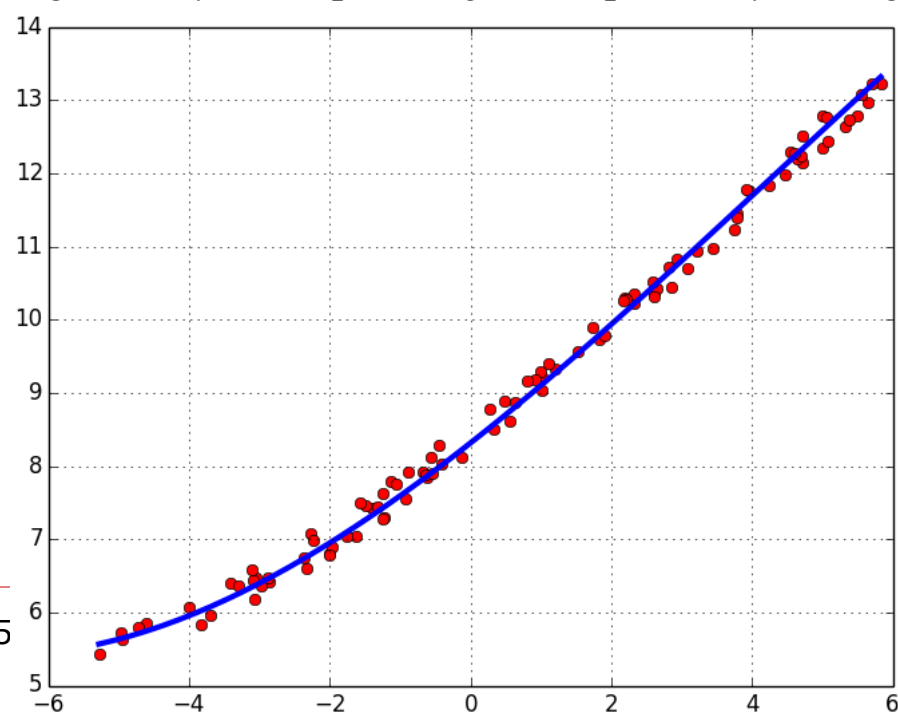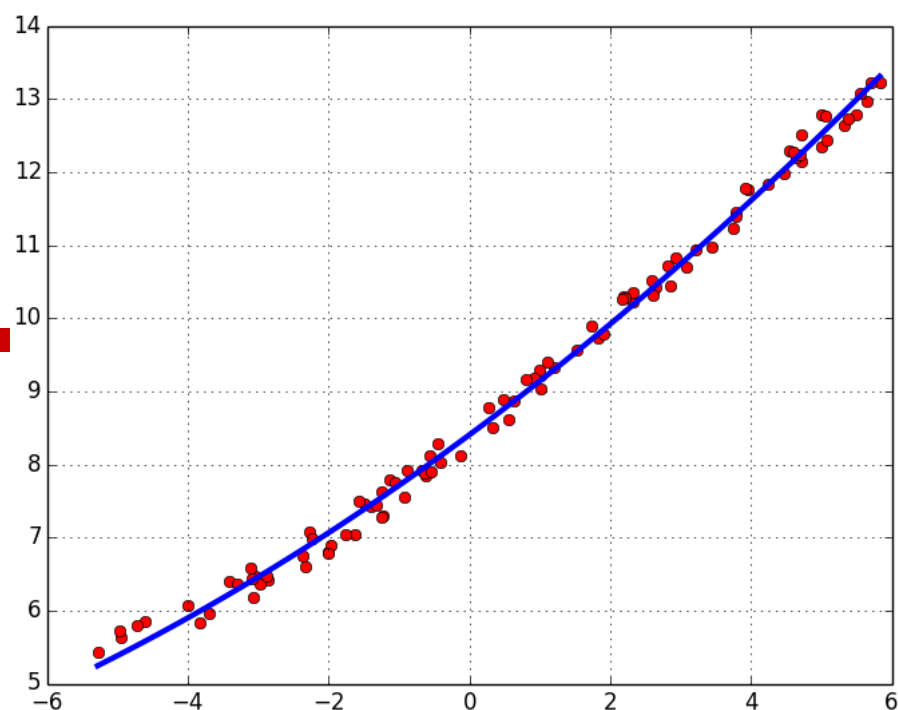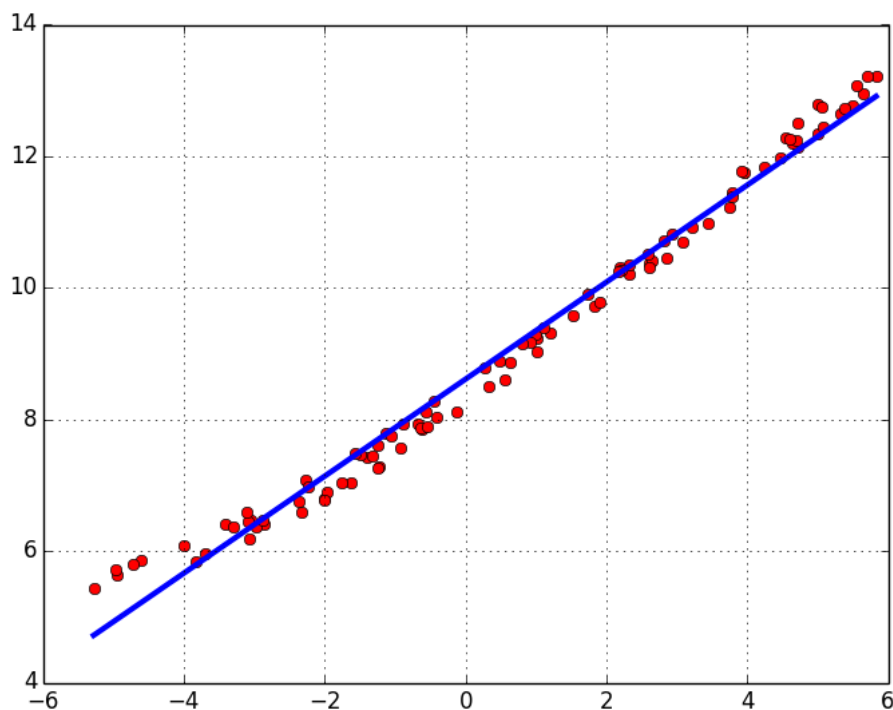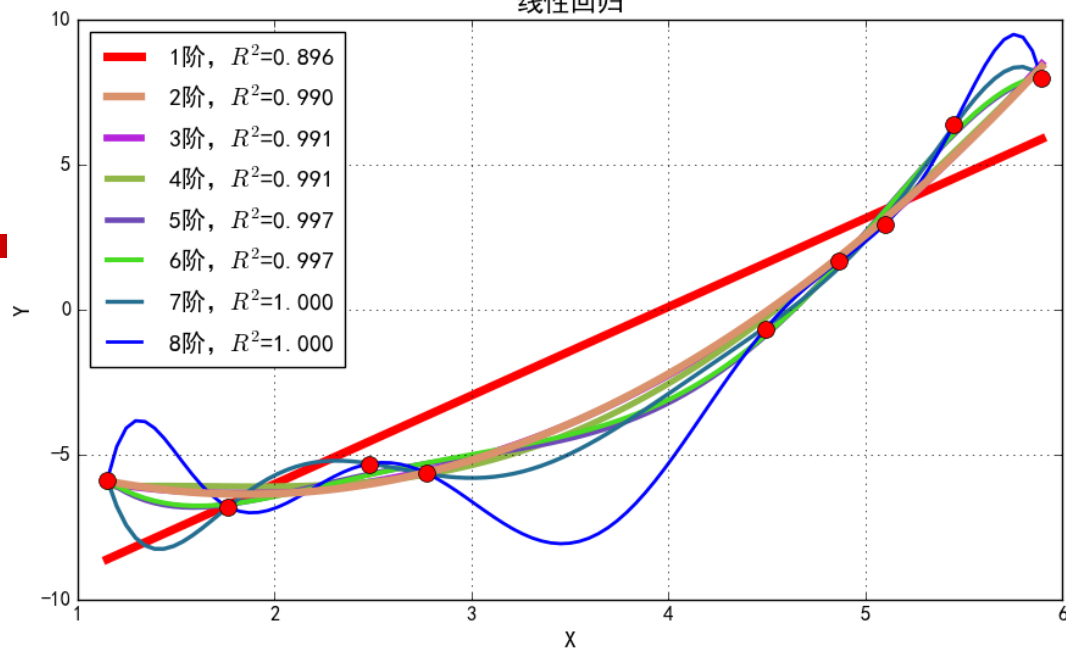
# 线性回归

# 线性回归

# 特征选择

# 超参与过拟合



多项式曲线拟合

线性回归

| | |
|---|---|
| 1阶，$R^2=0.896$ | |
| 2阶，$R^2=0.990$ | |
| 3阶，$R^2=0.991$ | |
| 4阶，$R^2=0.991$ | |
| 5阶，$R^2=0.997$ | |
| 6阶，$R^2=0.997$ | |
| 7阶，$R^2=1.000$ | |
| 8阶，$R^2=1.000$ | |

Ridge回归

| | |
|---|---|
| 1阶，$R^2=0.891$ | |
| 2阶，$R^2=0.990$ | |
| 3阶，$R^2=0.990$ | |
| 4阶，$R^2=0.987$ | |
| 5阶，$R^2=0.974$ | |
| 6阶，$R^2=0.997$ | |
| 7阶，$R^2=0.997$ | |
| 8阶，$R^2=0.996$ | |

# Coefficient of Determination

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^{m}(\hat{y}_i - y_i)^2}{\sum_{i=1}^{m}(y_i - \bar{y})^2}$$
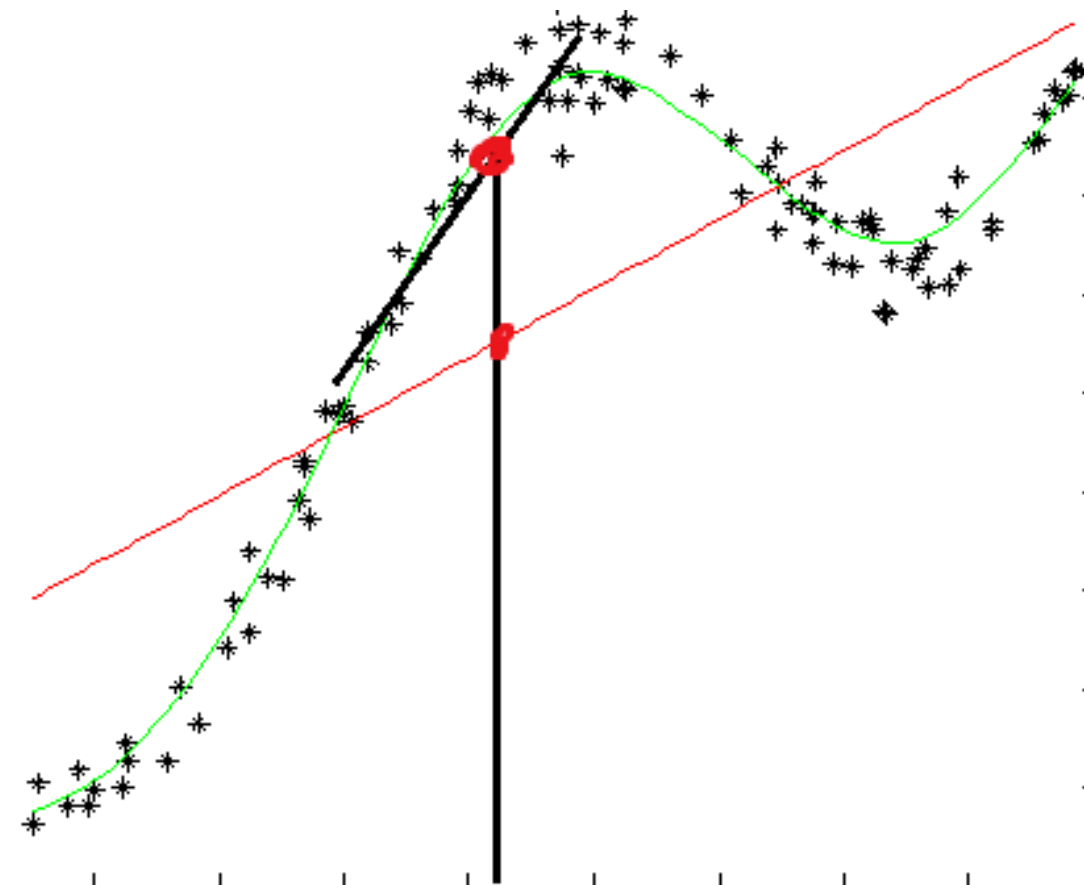
☐ 对于m个样本 $(\vec{x}_1, y_1), (\vec{x}_2, y_2), \cdots, (\vec{x}_m, y_m)$

☐ 某模型的估计值为 $(\vec{x}_1, \hat{y}_1), (\vec{x}_2, \hat{y}_2), \cdots, (\vec{x}_m, \hat{y}_m)$

☐ 计算样本的总平方和TSS(Total Sum of Squares)：$TSS = \sum_{i=1}^{m}(y_i - \bar{y})^2$
  - ■ 即样本伪方差的m倍 $Var(Y) = TSS/m$

☐ 计算残差平方和RSS(Residual Sum of Squares)：$RSS = \sum_{i=1}^{m}(\hat{y}_i - y_i)^2$
  - ■ 注：RSS即误差平方和SSE(Sum of Squares for Error)

☐ 定义 $R^2 = 1 - RSS/TSS$
  - ■ $R^2$越大，拟合效果越好
  - ■ $R^2$的最优值为1；若模型预测为随机值，$R^2$有可能为负
  - ■ 若预测值恒为样本期望，$R^2$为0

☐ 亦可定义ESS(Explained Sum of Squares)：$ESS = \frac{1}{m}\sum_{i=1}^{m}(\hat{y}_i - \bar{y})^2$
  - ■ TSS = ESS + RSS
  - ■ ESS又称回归平方和SSR(Sum of Squares for Regression)

# 局部加权回归

- 黑色是样本点

- 红色是线性回归曲线

- 绿色是局部加权回归曲线

# 局部加权线性回归

- ☐ LWR：Locally Weighted linear Regression

1. Fit $\theta$ to minimize $\sum_i (y^{(i)} - \theta^T x^{(i)})^2$

2. Output $\theta^T x$.

1. Fit $\theta$ to minimize $\sum_i w^{(i)}(y^{(i)} - \theta^T x^{(i)})^2$

2. Output $\theta^T x$.

# 权值的设置

☐ ω的一种可能的选择方式(高斯核函数)：

$$w^{(i)} = \exp\left(-\frac{(x^{(i)} - x)^2}{2\tau^2}\right)$$
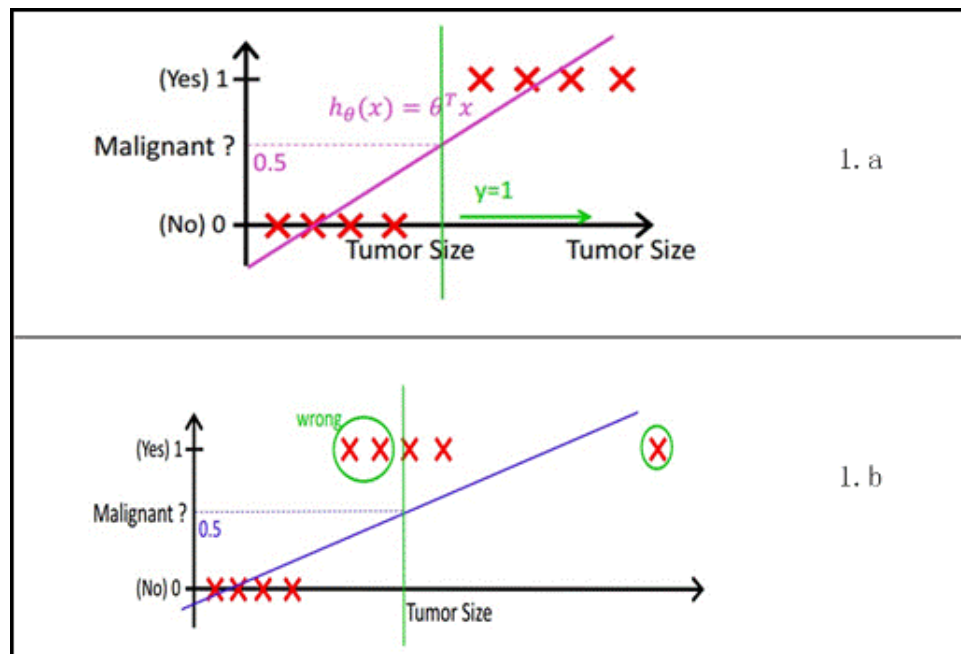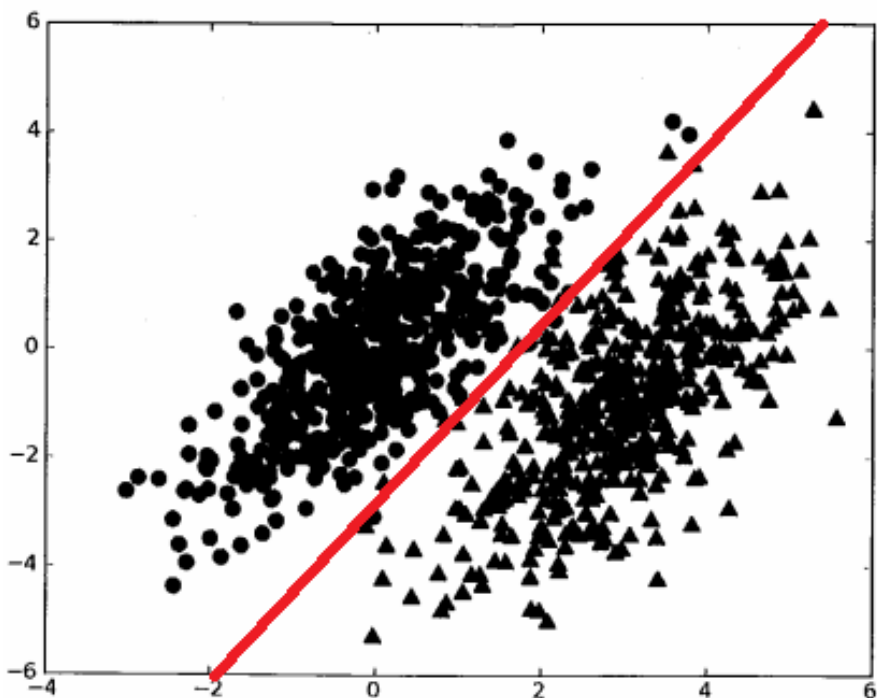
☐ τ称为带宽，它控制着训练样本随着与x$^{(i)}$距离的衰减速率。

☐ 多项式核函数

$$\kappa(\boldsymbol{x}_1, \boldsymbol{x}_2) = (\langle \boldsymbol{x}_1, \boldsymbol{x}_2 \rangle + R)^d$$

■ 在SVM章节继续核函数的讨论。

小象学院
ChinaHadoop.cn

# 用回归解决分类问题，如何？

# 线性回归-Logistic回归

- □ **紫色：**
  - ■ 线性回归
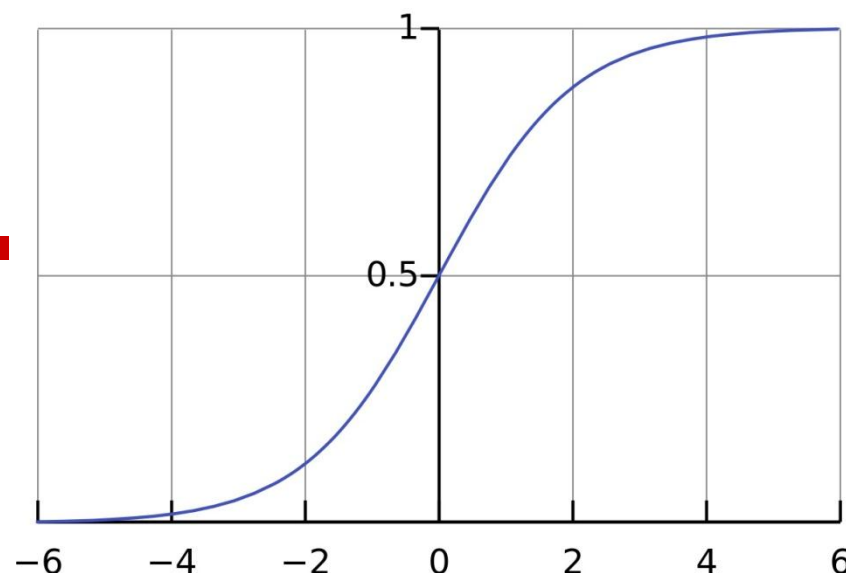- □ **绿色：**
  - ■ Logistic回归



- □ **左侧：**
  - ■ 线性回归
- □ **右侧：**
  - ■ Softmax回归

# Logistic回归

□ Logistic/sigmoid函数

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

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(x) = \left(\frac{1}{1 + e^{-x}}\right)' = \frac{e^{-x}}{\left(1 + e^{-x}\right)^2}$$

$$= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} = \frac{1}{1 + e^{-x}} \cdot \left(1 - \frac{1}{1 + e^{-x}}\right)$$

$$= g(x) \cdot (1 - g(x))$$

小象学院
ChinaHadoop.cn

# Logistic回归参数估计

□ 假定：

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

$$p(y \mid x; \theta) = (h_\theta(x))^y (1 - h_\theta(x))^{1-y}$$

$$
\begin{aligned}
L(\theta) &= p(\vec{y} \mid X; \theta) \\
&= \prod_{i=1}^{m} p(y^{(i)} \mid x^{(i)}; \theta) \\
&= \prod_{i=1}^{m} \left(h_\theta(x^{(i)})\right)^{y^{(i)}} \left(1 - h_\theta(x^{(i)})\right)^{1-y^{(i)}}
\end{aligned}
$$

小象学院
ChinaHadoop.cn

# 对数似然函数

$$\ell(\theta) = \log L(\theta)$$

$$= \sum_{i=1}^{m} y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))$$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x)$$

$$= \left( y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x)(1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x$$

$$= \left( y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x) \right) x_j$$

$$= (y - h_\theta(x)) x_j$$

# 参数的迭代

□ Logistic回归参数的学习规则：

$$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

□ 比较上面的结果和线性回归的结论的差别：

■ 它们具有相同的形式！

Repeat until convergence {

$$\theta_j := \theta_j + \alpha \sum_{i=1}^{m} \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$

}

Loop {
  for i=1 to m, {
    $$\theta_j := \theta_j + \alpha \left( y^{(i)} - h_\theta(x^{(i)}) \right) x_j^{(i)}$$
  }
}

# 对数线性模型

☐ 一个事件的几率odds，是指该事件发生的概率与该事件不发生的概率的比值。

☐ 对数几率：logit函数

$$P(y = 1 \mid x; \theta) = h_\theta(x)$$
$$P(y = 0 \mid x; \theta) = 1 - h_\theta(x)$$

$$\log it(p) = \log \frac{p}{1-p} = \log \frac{h_\theta(x)}{1-h_\theta(x)} = \log\left(\frac{\frac{1}{1+e^{-\theta^T x}}}{\frac{e^{-\theta^T x}}{1+e^{-\theta^T x}}}\right) = \theta^T x$$

# Logistic回归的损失函数 $y_i \in \{0,1\}$

$$L(\theta) = \prod_{i=1}^{m} p_i^{y_i}(1-p_i)^{1-y_i}$$

$$\hat{y}_i = \begin{cases} p_i & y_i = 1 \\ 1-p_i & y_i = 0 \end{cases}$$

$$\Rightarrow l(\theta) = \sum_{i=1}^{m} \ln\left[ p_i^{y_i}(1-p_i)^{1-y_i} \right]$$

$$\xrightarrow{p_i = \frac{1}{1+e^{-f_i}}} l(\theta) = \sum_{i=1}^{m} \ln\left[ \left(\frac{1}{1+e^{-f_i}}\right)^{y_i}\left(\frac{1}{1+e^{f_i}}\right)^{1-y_i} \right]$$

$$\therefore\ loss(y_i, \hat{y}_i) = -l(\theta)$$

$$= \sum_{i=1}^{m}\left[ y_i \ln\left(1+e^{-f_i}\right) + (1-y_i)\ln\left(1+e^{f_i}\right) \right]$$

小象学院
ChinaHadoop.cn

# Logistic回归的损失：$y_i \in \{-1,1\}$

$$L(\theta) = \prod_{i=1}^{m} p_i^{(y_i+1)/2} (1-p_i)^{-(y_i-1)/2} \Rightarrow l(\theta) = \sum_{i=1}^{m} \ln\left[ p_i^{(y_i+1)/2} (1-p_i)^{-(y_i-1)/2} \right]$$

$$\xrightarrow{p_i = \frac{1}{1+e^{-f_i}}} l(\theta) = \sum_{i=1}^{m} \ln\left[ \left(\frac{1}{1+e^{-f_i}}\right)^{(y_i+1)/2} \left(\frac{1}{1+e^{f_i}}\right)^{-(y_i-1)/2} \right]$$

$$\therefore \ loss(y_i, \hat{y}_i) = -l(\theta)$$

$$= \sum_{i=1}^{m} \left[ \frac{1}{2}(y_i+1)\ln(1+e^{-f_i}) - \frac{1}{2}(y_i-1)\ln(1+e^{f_i}) \right]$$

$$= \begin{cases} \sum_{i=1}^{m} \left[\ln(1+e^{-f_i})\right] & y_i = 1 \\ \sum_{i=1}^{m} \left[\ln(1+e^{f_i})\right] & y_i = -1 \end{cases} \Rightarrow loss(y_i, \hat{y}_i) = \sum_{i=1}^{m} \left[\ln(1+e^{-y_i \cdot f_i})\right]$$

$$y_i = \{-1,1\}$$

$$\hat{y}_i = \begin{cases} p_i & y_i = 1 \\ 1-p_i & y_i = -1 \end{cases}$$

小象学院
ChinaHadoop.cn

# 分类：Logistic回归

☐ 沿似然函数正梯度上升

☐ 维度提升

```python
def logistic_regression(data, alpha, lamda):
    n = len(data[0]) - 1
    w = [0 for x in range(n)]
    w2 = [0 for x in range(n)]
    for times in range(10000):
        for d in data:
            for i in range(n):
                w2[i] = w[i] + alpha * (d[n] - h(w,d))*d[i] + lamda * w[i]
            for i in range(n):
                w[i] = w2[i]
        print w
    return w
```

```python
def feature_whole(x1, x2, e):
    d = [1]
    for n in range(1,e+1):
        for i in range(n+1):
            d.append(pow(x1, n-i) * pow(x2, i))
    return d
```

小象学院
ChinaHadoop.cn

# 异或



x1        x2        y
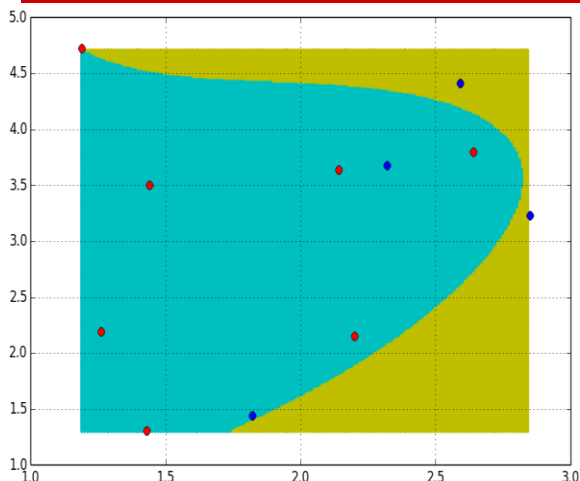0          0         0
0          1         1
1          0         1
1          1         0

# 数据升维： "选取特征"

# 复习：指数族

- 指数族概念的目的，是为了说明广义线性模型Generalized Linear Models

- 凡是符合指数族分布的随机变量，都可以用GLM回归分析

# 广义线性模型GLM

- □ y不再只是正态分布，而是扩大为指数族中的任一分布；
- □ 变量x → g(x) → y
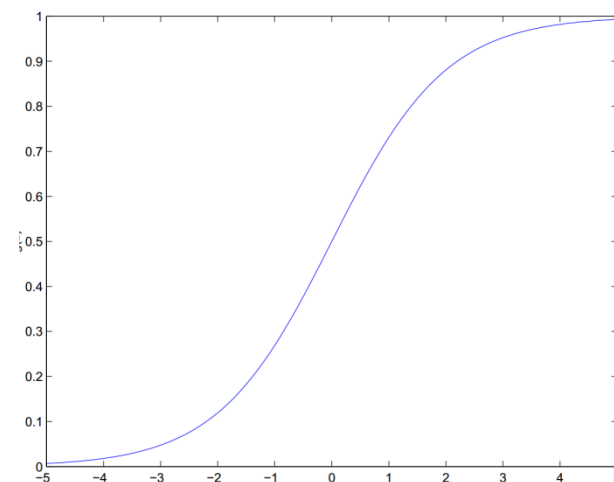  - ■ 连接函数g
- □ 连接函数g单调可导
  - ■ 如Logistic回归中的 $g(z) = \dfrac{1}{1+e^{-z}}$
  - ■ 拉伸变换：
  $$g(z) = \dfrac{1}{1+e^{-\lambda z}}$$

# Softmax回归

- K分类，第k类的参数为$\vec{\theta}_k$，组成二维矩阵$\theta_{k \times n}$

- 概率：$p(c=k \mid x;\theta) = \dfrac{\exp(\theta_k^T x)}{\sum\limits_{l=1}^{K} \exp(\theta_l^T x)}, \quad k=1,2,\cdots,K$

- 似然函数：$L(\theta) = \prod\limits_{i=1}^{m} \prod\limits_{k=1}^{K} p(c=k \mid x^{(i)};\theta)^{y_k^{(i)}} = \prod\limits_{i=1}^{m} \prod\limits_{k=1}^{K} \dfrac{\exp(\theta_k^T x^{(i)})^{y_k^{(i)}}}{\sum\limits_{l=1}^{K} \exp(\theta_l^T x^{(i)})}$

- 对数似然：$J_m(\theta) = \ln L(\theta) = \sum\limits_{i=1}^{m} \sum\limits_{k=1}^{K} \left( y_k^{(i)} \cdot \theta_k^T x^{(i)} - \ln \sum\limits_{l=1}^{K} \exp(\theta_l^T x^{(i)}) \right)$

- 随机梯度：$J(\theta) = \sum\limits_{k=1}^{K} y_k \cdot \left( \theta_k^T x - \ln \sum\limits_{l=1}^{K} \exp(\theta_l^T x) \right)$

$$\frac{\partial J(\theta)}{\partial \theta_k} = (y_k - p(y_k \mid x;\theta)) \cdot x$$
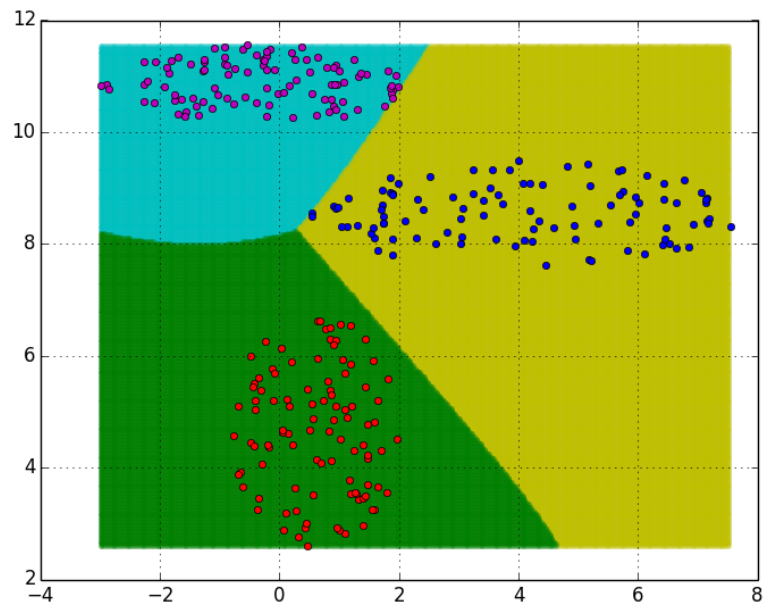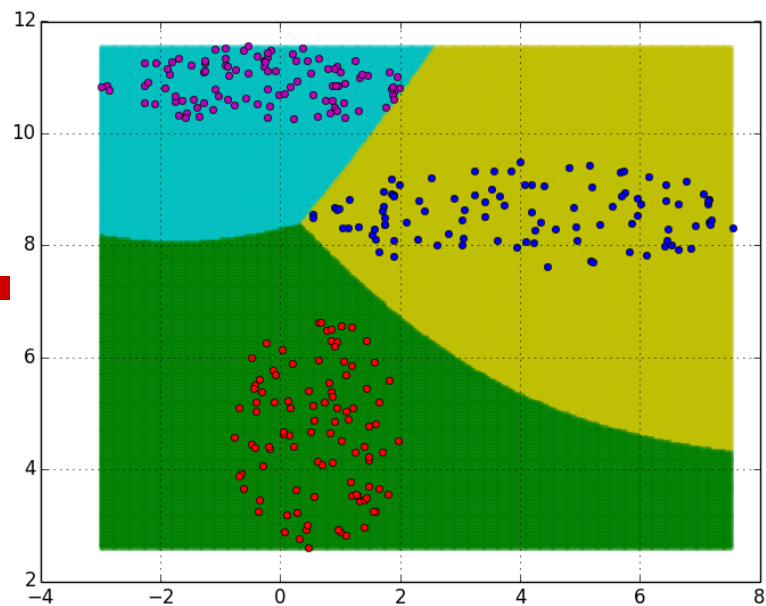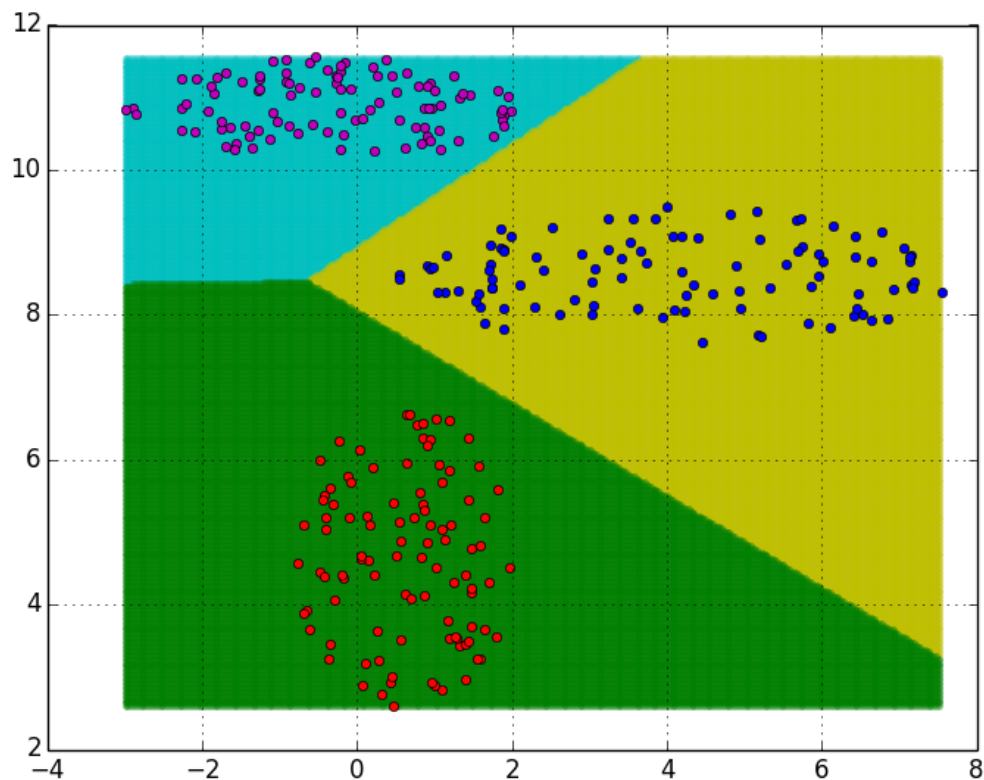
# Code

```python
def soft_max(data, K, alpha, lamda):
    n = len(data[0]) - 1       # 样本维度
    w = np.zeros((K, n))       # 当前权值：每个类别有各自的权值
    wNew = np.zeros((K, n))       # 临时权值：迭代过程中的权值
    for times in range(1000):
        for d in data:
            x = d[:-1]                    # 输入向量
            for k in range(K):           # K：类别数目
                y = 0                     # 期望输出(标量)
                if int(d[-1] + 0.5) == k:
                    y = 1
                p = predict(w, x, k)
                g = (y - p) * x          # 梯度(n维列向量)
                wNew[k] = w[k] + (alpha * g + lamda * w[k])
            w = wNew.copy()
    return w
```
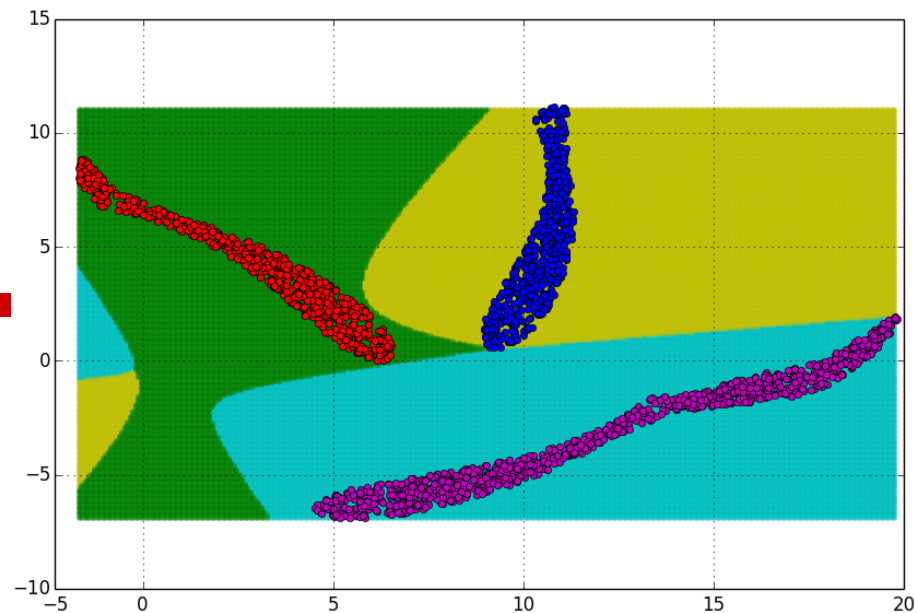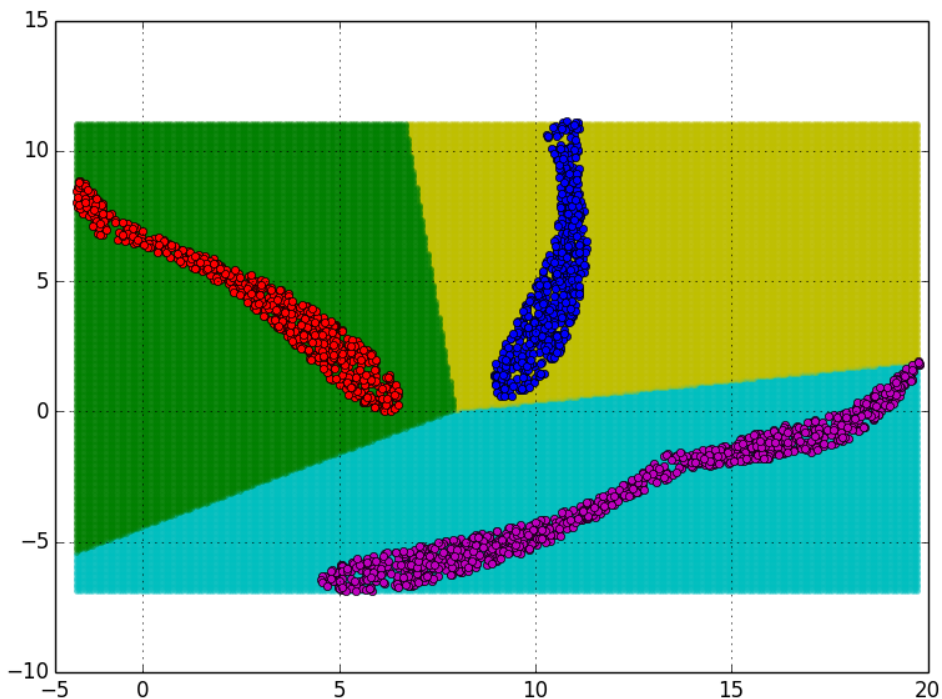
小象学院
ChinaHadoop.cn

# Softmax分类

# 特征选择

# 总结和思考

☐ Logistic/Softmax回归是实践中解决分类问题的最重要方法。

   ■ 方法简单、容易实现、效果良好、易于解释

   ■ 不止是分类：推荐系统

☐ 特征选择很重要，除了人工选择，还可以用其他机器学习方法，如随机森林、PCA、LDA等。

☐ 梯度下降算法是参数优化的重要手段，尤其SGD。

   ■ 适用于在线学习

   ■ 跳出局部极小值

   ■ 思考：计算可逆方阵的逆，可否使用梯度下降算法？

---
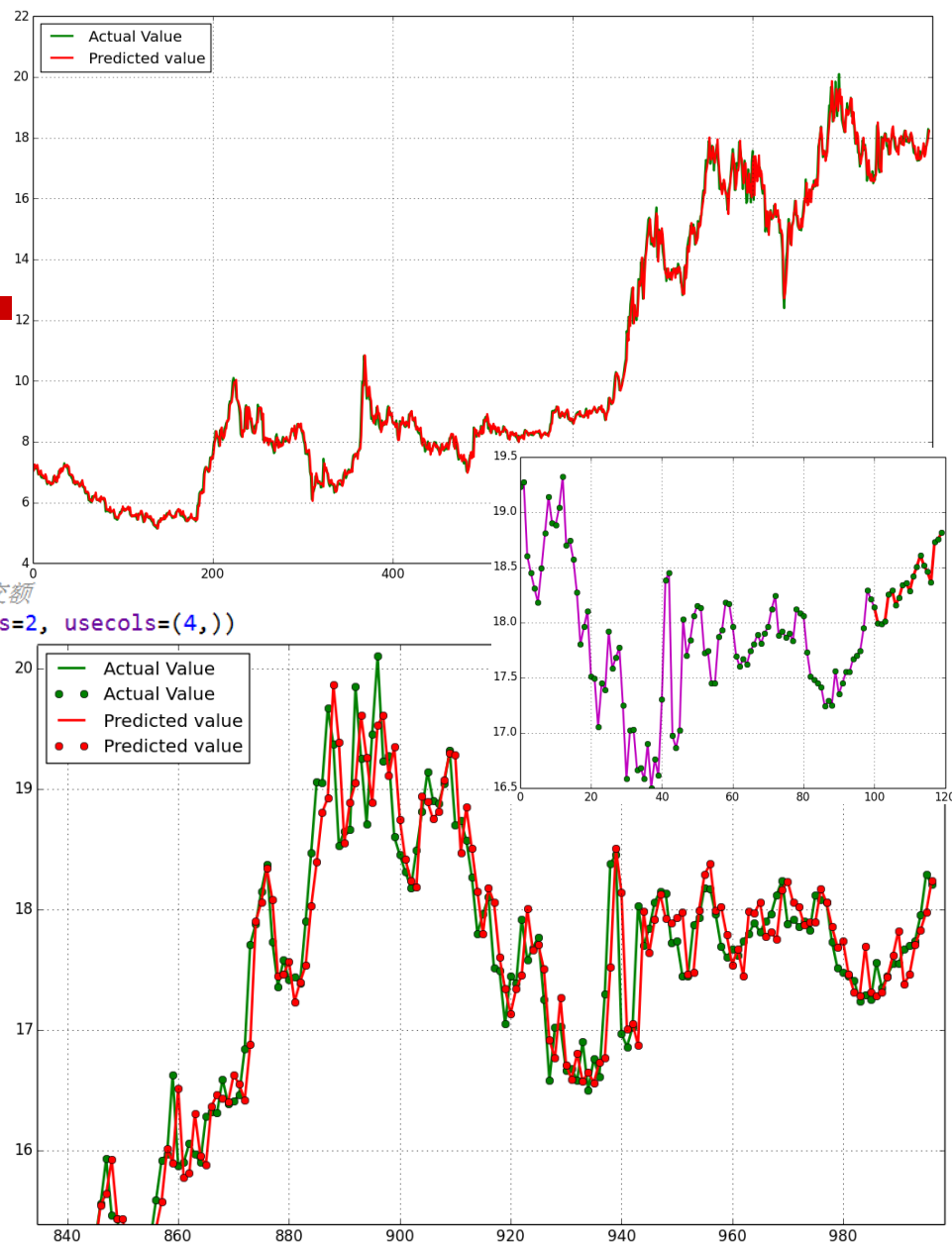
小象学院
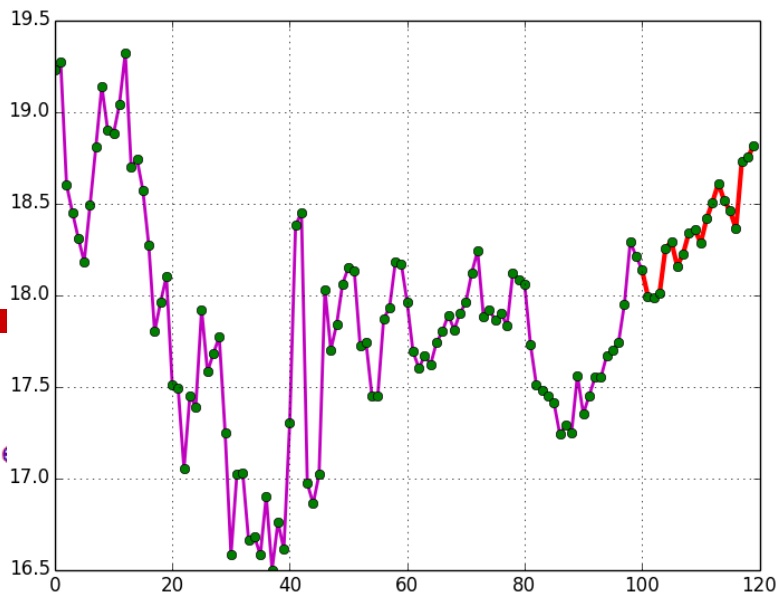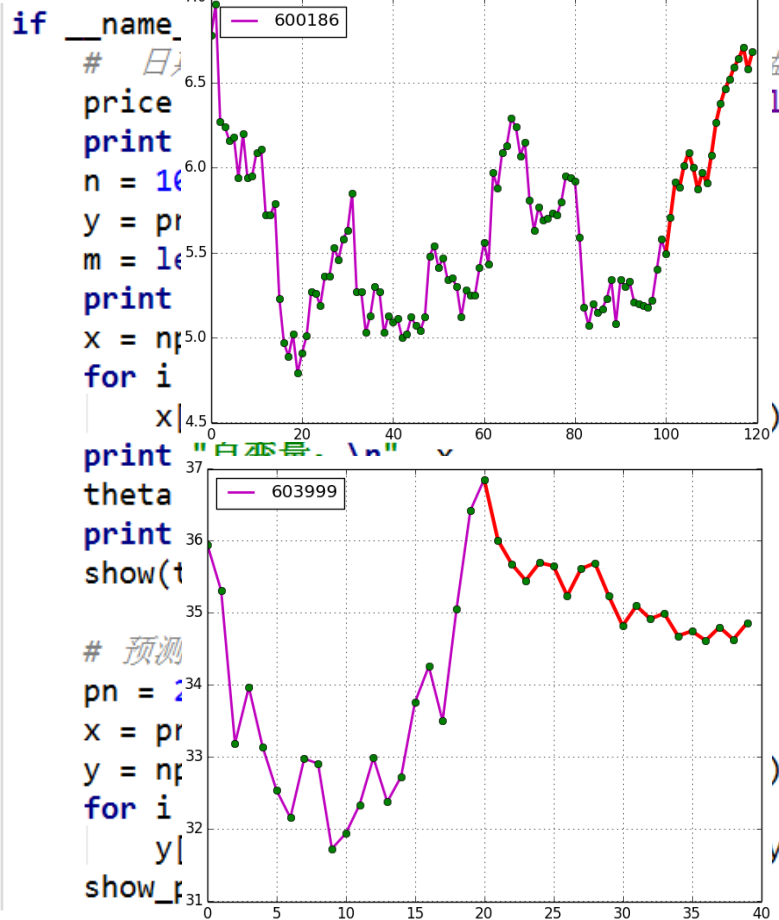ChinaHadoop.cn

# 股价拟合

□ 方法：自回归

□ 参数：100阶

```python
if __name__ == "__main__":
    # 日期     开盘     最高     最低     收盘     成交量     成交额
    price = np.loadtxt('sh_600000.txt', delimiter='\t', skiprows=2, usecols=(4,))
    print "原始价格: \n", price
    n = 100            # 阶数
    y = price[n:]
    m = len(y)         # 样本个数
    print "预测价格: \n", y
    x = np.zeros((m, n+1))
    for i in range(m):
        x[i] = np.hstack((price[i:i+n], 1))
    print "自变量: \n", x
    theta = np.linalg.lstsq(x, y)[0]     # theta为回归系数
    print theta
    show(theta, x, y)

    # 预测
    pn = 20        # 预测未来多少天
    x = price[-n:]
    y = np.hstack((price[-n:], np.zeros(pn)))
    for i in range(pn):
        y[n+i] = np.dot(theta, np.hstack((y[i:n+i], 1)))
    show_predict(y, pn)
```

# 股价预测

# 作业

- 解释线性回归中使用误差平方和作为目标函数的原因。
- 请描述BGD和SGD的区别，并指出SGD的优势有哪些？
- 推导Logistic回归的梯度。

小象学院
ChinaHadoop.cn

# 参考文献

- Prof. Andrew Ng. *Machine Learning*. Stanford University

- 李航，统计学习方法，清华大学出版社，2012

- Stephen Boyd, Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004
    - 中译本：王书宁，许鋆，黄晓霖，凸优化，清华大学出版社，2013

# 我们在这里



- http://wenda.ChinaHadoop.cn
  - 视频/课程/社区
- 微博
  - @ChinaHadoop
  - @邹博_机器学习
- 微信公众号
  - 小象
  - 大数据分析挖掘

# 感谢大家！

# 恳请大家批评指正！