

# 线性回归

---

最小二乘法

赵海涛

haitaozhao@ecust.edu.cn

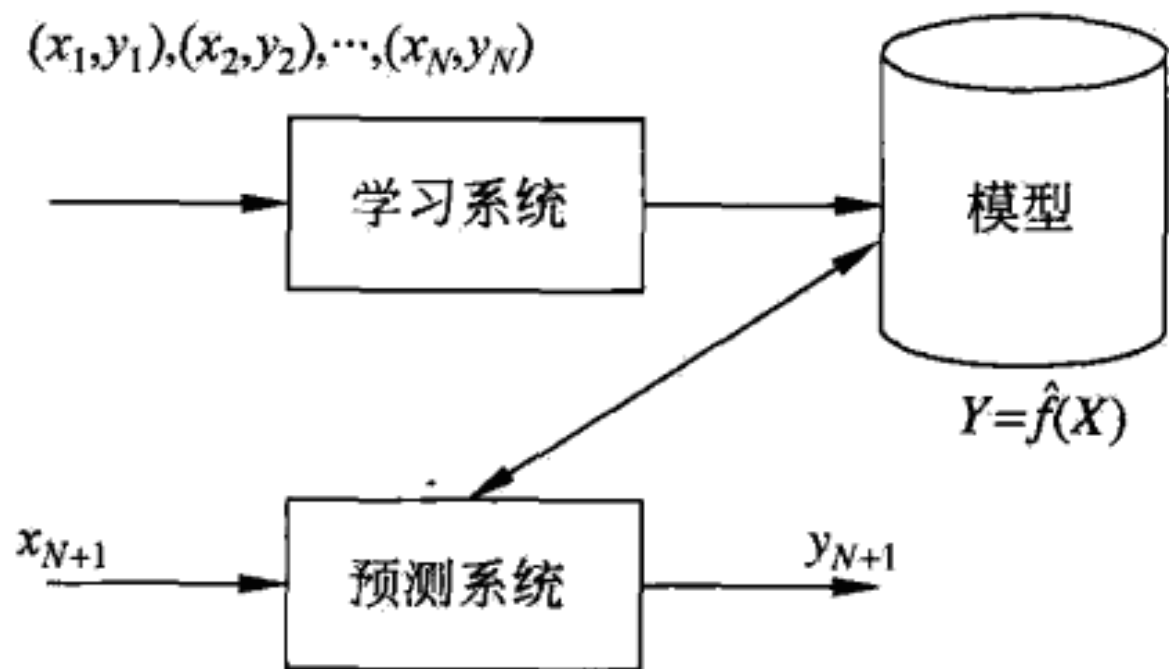
# 大纲

- 回归模型
- 回归的目标函数（最小二乘准则）
- 线性回归的求解
- 线性回归的扩展

# 回归问题

- 回归模型是表示从输入变量到输出变量之间映射的函数。回归问题的学习等价于函数拟合。
- 学习和预测两个阶段
- 训练集：

$$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$$

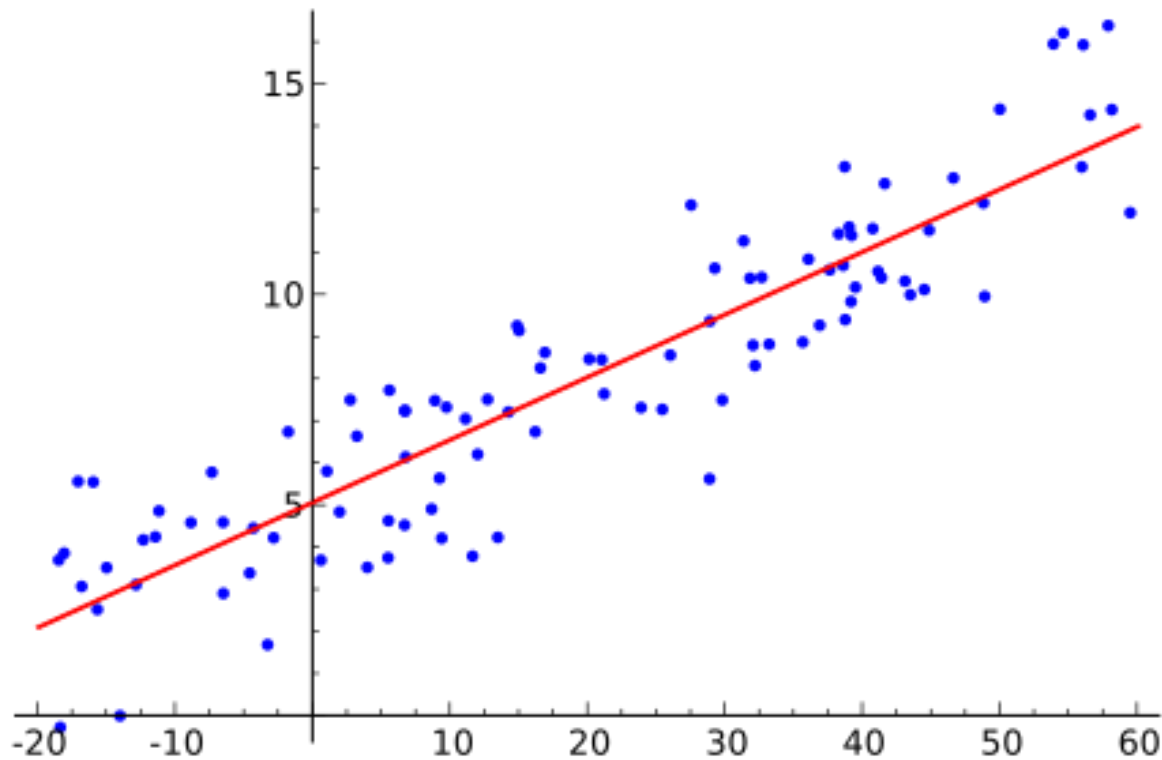
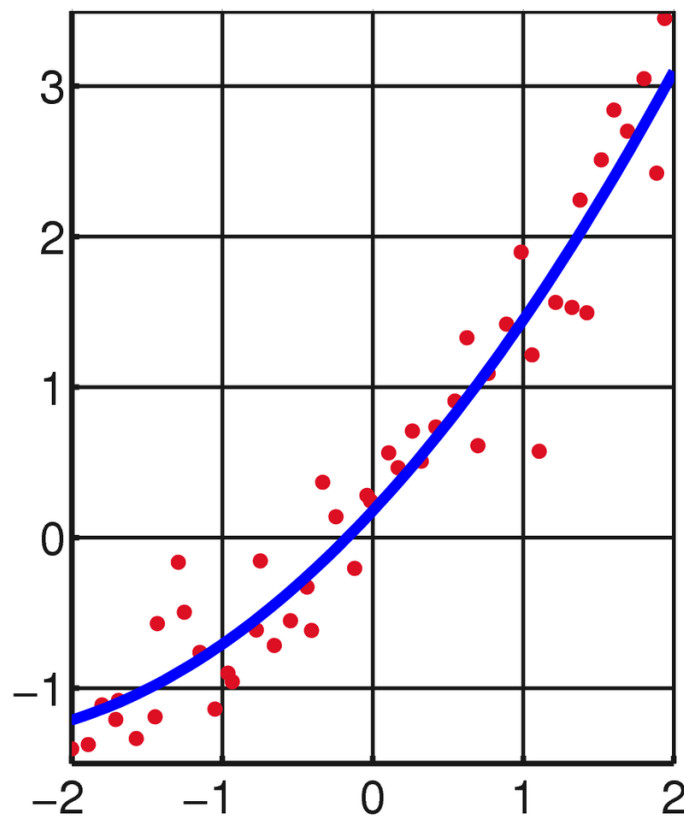


# 回归问题

- 回归学习最常用的损失函数是平方损失函数，在此情况下，回归问题可以由著名的最小二乘法(least squares)求解。

- 股价预测

- 房价预测



# 回归问题

- 假设我们有一个数据集包含 $n$ 个数据对 $(\mathbf{x}_i, y_i)$ ,  $i = 1, 2, \dots, n$ , 其中 $\mathbf{x}_i$ 是自变量,  $\mathbf{x}_i \in R^{1 \times d}$ ,  $y_i$ 是因变量,  $y_i \in R$ 。 $\mathbf{x}_i$ 可以认为是我们在生产过程中获得的一些传感器和仪表的数据构成的向量, 比如 $\mathbf{x}_i$ 的各个维度分别对应压力, 温度, 浓度等等,  $i$ 表示采样的时刻,  $y_i$ 是对应 $\mathbf{x}_i$ 的某个指标, 比如产品的质量指标等, 它是连续的一个数值。
- 我们希望有一个模型 $f(\mathbf{x}_i, \boldsymbol{\beta})$ 可以预测未来的 $y_i$ , 即 $y_i = f(\mathbf{x}_i, \boldsymbol{\beta})$  (其中 $\mathbf{x}_i \in R^{1 \times d}$ ,  $\boldsymbol{\beta} \in R^{d \times 1}$ )。简单地说, 这个模型应该“非常好”地适合于已有的数据, 我们令误差 $\varepsilon_i = y_i - f(\mathbf{x}_i, \boldsymbol{\beta})$ 。

# 回归问题的目标函数

- 回归的目标就是要求取 $\boldsymbol{\beta}$ 来极小化 $\varepsilon_i$ 的平方和

$$J(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \varepsilon_i^2 = \frac{1}{n} \sum_{i=1}^n [y_i - f(\mathbf{x}_i, \boldsymbol{\beta})]^2$$

上式即为回归问题中的目标函数。

- 通常，将 $f(\mathbf{x}_i, \boldsymbol{\beta})$ 进行简化，最简单的， $f(\mathbf{x}_i, \boldsymbol{\beta})$ 可以看作是一个超平面，即 $f(\mathbf{x}_i, \boldsymbol{\beta}) = \mathbf{x}_i \boldsymbol{\beta} + \beta_0$ ，其中 $\boldsymbol{\beta} = (\beta_1, \dots, \beta_d)^T$ 。则有方程

$$\begin{cases} y_1 = \beta_0 + \mathbf{x}_1 \boldsymbol{\beta} + \varepsilon_1 \\ y_2 = \beta_0 + \mathbf{x}_2 \boldsymbol{\beta} + \varepsilon_2 \\ \vdots \\ y_n = \beta_0 + \mathbf{x}_n \boldsymbol{\beta} + \varepsilon_n \end{cases}$$

# 线性回归

• 令  $Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$ ,  $X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & \ddots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix}$ ,  $E = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$ ,  $B = \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta} \end{bmatrix}$ , 则有

$$Y = XB + E$$

$$J(B) = \frac{1}{n} \sum_{i=1}^n \left[ y_i - \left( \sum x_{ij} \beta_j + \beta_0 \right) \right]^2 = \frac{1}{n} (Y - XB)^T (Y - XB)$$

- $\beta_0$  的作用主要是产生一个偏置, 如  $\sum \mathbf{x}_i = \mathbf{0}$ , 而  $\sum y_i \neq 0$  时, 那么在  $y$  轴上就会产生一个截距, 这是线性回归问题优化时要考虑的。

# 线性回归问题的求解

- 求解  $J(B)$  的最小值，可以直接求解，也可以用求导数的方式来求解，即用梯度下降法来求解。
- 梯度下降法和牛顿法求解
- 通过正交投影直接求解



# 梯度下降法

一元情况下的Taylor展开式：

$$f(x + \delta) = f(x) + f'(x)\delta + \frac{1}{2}f''(x)\delta^2 + \dots$$

推广到多元情况下的Taylor展开式，即 $\mathbf{x} \in R^{d \times 1}, f(\mathbf{x}) \in R$ 时：

$$f(\mathbf{x} + \Delta \mathbf{x}) = f(\mathbf{x}) + \nabla^T f(\mathbf{x}) \Delta \mathbf{x} + \frac{1}{2} \Delta \mathbf{x}^T \nabla^2 f(\mathbf{x}) \Delta \mathbf{x} + \dots$$

$$\text{其中, } \nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}, \quad \nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \frac{\partial^2 f}{\partial x_2 \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_d \partial x_1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_1 \partial x_d} & \dots & \dots & \frac{\partial^2 f}{\partial x_d \partial x_d} \end{bmatrix}.$$

# 梯度下降法

- 梯度下降法的目的就是求取 $f(\mathbf{x})$ 的极值。我们知道当 $\nabla f(\mathbf{x}) = \mathbf{0}$ ，且 $\nabla^2 f(\mathbf{x}) \succ 0$ 时(这里表示 $\nabla^2 f(\mathbf{x})$ 应为正定矩阵，正定矩阵的充要条件为它的所有特征值都大于零或者各阶主子式大于零),存在唯一极值点。
- 对于一阶Taylor展开式，当 $\Delta \mathbf{x}$ 足够小的时候有以下近似式：

$$f(\mathbf{x} + \Delta \mathbf{x}) \approx f(\mathbf{x}) + \nabla^T f(\mathbf{x}) \Delta \mathbf{x}$$

从而

$$f(\mathbf{x} + \Delta \mathbf{x}) - f(\mathbf{x}) \approx \nabla^T f(\mathbf{x}) \Delta \mathbf{x}$$

# 梯度下降法

如果已经有一个 $\mathbf{x}_t$ ，希望 $\mathbf{x}_t + \Delta\mathbf{x}_t$ 使 $f(\mathbf{x}_t + \Delta\mathbf{x}_t)$ 小于 $f(\mathbf{x}_t)$ ，那么可以令

$$\Delta\mathbf{x}_t = -\nabla f(\mathbf{x}_t)。$$

此时就有

$$f(\mathbf{x}_t + \Delta\mathbf{x}_t) - f(\mathbf{x}_t) = -\nabla^T f(\mathbf{x}_t) \nabla f(\mathbf{x}_t) = -\sum_{i=1}^d \left( \frac{\partial f(\mathbf{x}_t)}{\partial x_i} \right)^2 < 0$$

在Taylor展开中，通常需要对 $\Delta\mathbf{x}_t$ 有一定的限制，在梯度下降法中一般取

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \alpha \Delta\mathbf{x}_t，$$

其中 $\alpha$ 为学习率或步长。

# 梯度下降法的步骤

下面给出进行梯度下降法求解的一般步骤：

- ① 给定初值  $\mathbf{x}_1$ ,  $i = 1$ , 学习率  $\alpha$ , 给定阈值  $\varepsilon$ ;
- ② 求  $\nabla f(\mathbf{x}_i)$ , 令  $\Delta \mathbf{x}_i = -\nabla f(\mathbf{x}_i)$ ;
- ③ 求  $\mathbf{x}_{i+1} = \mathbf{x}_i + \alpha \Delta \mathbf{x}_i$ ;
- ④ 若  $\|f(\mathbf{x}_{i+1}) - f(\mathbf{x}_i)\|_2^2 \leq \varepsilon$ , 转向⑤, 否则,  $i = i + 1$ , 转向②
- ⑤ 输出  $\mathbf{x}_{i+1}$ .

# 牛顿法

- 如果考虑二阶Taylor展开, 则:

$$f(\mathbf{x} + \Delta\mathbf{x}) \approx f(\mathbf{x}) + \nabla^T f(\mathbf{x})\Delta\mathbf{x} + \frac{1}{2}\Delta\mathbf{x}^T \nabla^2 f(\mathbf{x})\Delta\mathbf{x}$$

注: 有很多书和论文用  $\mathbf{g} = \nabla f(\mathbf{x})$ ,  $H = \nabla^2 f(\mathbf{x})$  来代替书写。

- 要求解  $f(\mathbf{x} + \Delta\mathbf{x})$  的最小值, 可以看作求  $\Delta\mathbf{x}$  使得

$$f(\mathbf{x}) + \nabla^T f(\mathbf{x})\Delta\mathbf{x} + \frac{1}{2}\Delta\mathbf{x}^T \nabla^2 f(\mathbf{x})\Delta\mathbf{x}$$

最小化。对  $\Delta\mathbf{x}$  求导可得

$$\Delta\mathbf{x} = -(\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x}) = -H^{-1} \mathbf{g}$$

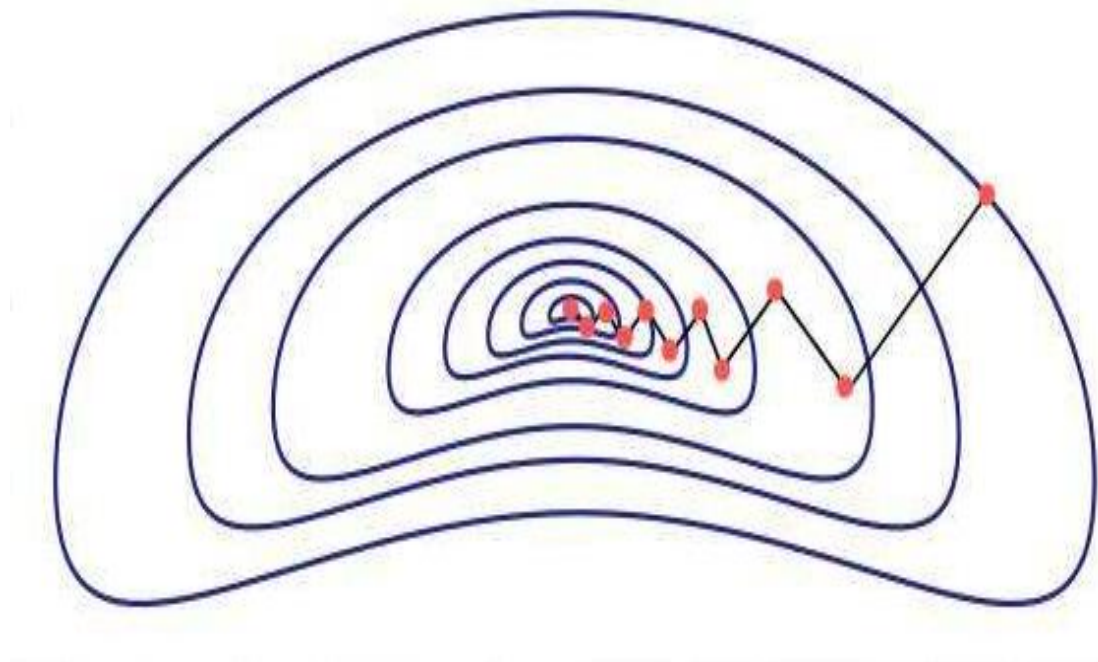
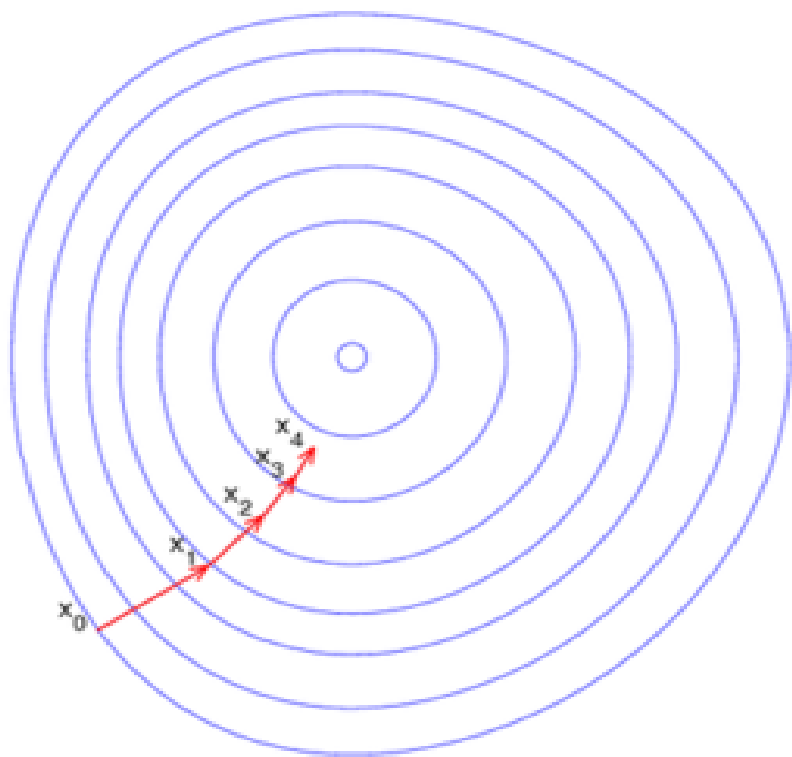
# 梯度下降法的收敛性

- 梯度下降法收敛速度慢与关于 $\mathbf{x}_k$ 的Hessian矩阵的“条件”有关。梯度下降法第 $k$ 次的迭代过程与之前迭代过程之间的关系可以通过下式来表示：

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(\frac{1-r}{1+r}\right)^2 (f(\mathbf{x}_k) - f(\mathbf{x}^*))$$

- 其中， $r$ 是关于 $\mathbf{x}_k$ 的Hessian矩阵（用 $H_k$ 表示）最小特征值与最大特征值的比值。从上式我们可以看出当 $r$ 较小或者 $H_k$ 的条件数较大的时候，梯度下降法的收敛速度就可能会较慢。经典的梯度下降法是线性收敛的，而牛顿法是二次收敛的。
- 另外，梯度下降法容易出现锯齿效应。

# 梯度下降法的收斂性



# 梯度下降法的收敛性

- 梯度下降法收敛速度慢与关于 $\mathbf{x}_k$ 的Hessian矩阵的“条件”有关。梯度下降法第 $k$ 次的迭代过程与之前迭代过程之间的关系可以通过下式来表示：

$$f(\mathbf{x}_{k+1}) - f(\mathbf{x}^*) \leq \left(\frac{1-r}{1+r}\right)^2 (f(\mathbf{x}_k) - f(\mathbf{x}^*))$$

- 其中， $r$ 是关于 $\mathbf{x}_k$ 的Hessian矩阵（用 $H_k$ 表示）最小特征值与最大特征值的比值。从上式我们可以看出当 $r$ 较小或者 $H_k$ 的条件数较大的时候，梯度下降法的收敛速度就可能会较慢。经典的梯度下降法是线性收敛的，而牛顿法是二次收敛的。
- 另外，梯度下降法容易出现锯齿效应。



# 梯度下降法求解线性回归

- $$\begin{aligned} J(B) &= \frac{1}{n} (Y - XB)^T (Y - XB) \\ &= \frac{1}{n} (Y^T Y - Y^T X B - B^T X^T Y + B^T X^T X B) \\ &= \frac{1}{n} (Y^T Y - 2B^T X^T Y + B^T X^T X B) \end{aligned}$$
- $$\nabla J(B) = \frac{2}{n} (-X^T Y + X^T X B)$$

线性回归问题的一般步骤：

- ① 给定初值  $B_1$ ,  $i = 1$ , 学习率  $\alpha$ , 给定阈值  $\varepsilon$ ;
- ② 求  $B_{i+1} = B_i - \alpha \nabla J(B_i)$  ;
- ③ 若  $\|J(B_{i+1}) - J(B_i)\|_2^2 \leq \varepsilon$ , 转向④, 否则,  $i = i + 1$ , 转向②
- ④ 输出  $B_{i+1}$ .

# 梯度下降法求解线性回归

- 一种常见的方法是针对不同维度进行自适应改进学习率，对不同维度采用不同的  $\alpha$ ，比较常用的是RMSProp(root mean square propagation)。首先对每个维度进行移动平均  $\mathbf{v}_t = \mu \mathbf{v}_{t-1} + (1 - \mu)[\nabla J(B_t)]^2$ （平方是每个元素的平方）。 $\mu$ 一般选为0.9,0.99或0.999。参数的更新变为

$$B_{t+1} = B_t - \frac{\alpha \nabla J(B_t)}{\sqrt{\mathbf{v}_t} + \varepsilon}$$

（开根号是对每个元素进行，分子除以分母也是对每个元素进行）

这里可发现和经典的梯度下降法相比，多了一个类似权重的项用于不同的维度上。

# 梯度下降法求解线性回归

- 另一类较常见的是Adam(adaptive moment estimation)

$$\mathbf{m}_t = \mu_0 \mathbf{m}_{t-1} + (1 - \mu_0) \nabla J(B_t)$$

$$\mathbf{v}_t = \mu_1 \mathbf{v}_{t-1} + (1 - \mu_1) [\nabla J(B_t)]^2$$

$$B_{t+1} = B_t - \frac{\alpha \mathbf{m}_t}{\sqrt{\mathbf{v}_t} + \varepsilon}$$

实际就是在梯度方向的选择上，也运用移动平均，既考虑当前的梯度方向，又考虑上一步的梯度方向，通过加权平均可减小锯齿效应。

注：其它改进的方法还有很多，主要也是针对方向和学习率。常用的这两个方法效果已经相当好。

# 梯度下降法求解线性回归

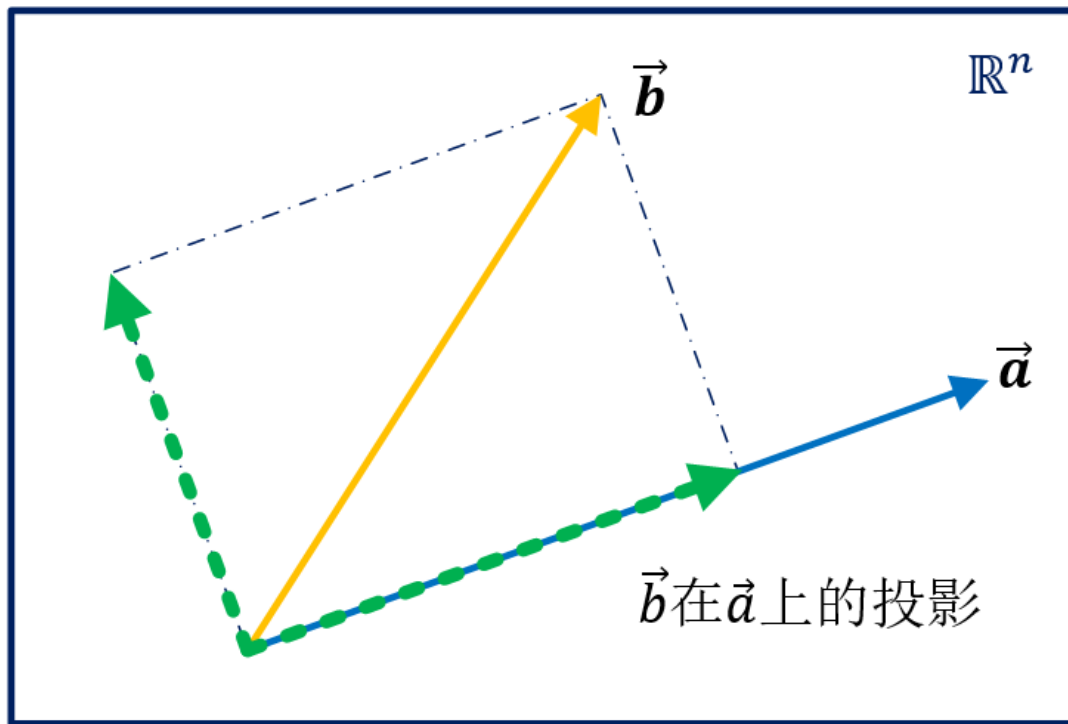
- 考虑牛顿法，即 $\Delta B = -(\nabla^2 J(B))^{-1} \nabla J(B) = -(X^T X)^{-1}(-X^T Y + X^T X B)$   
则，

$$B^* = B + \Delta B = B - B + (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T Y$$

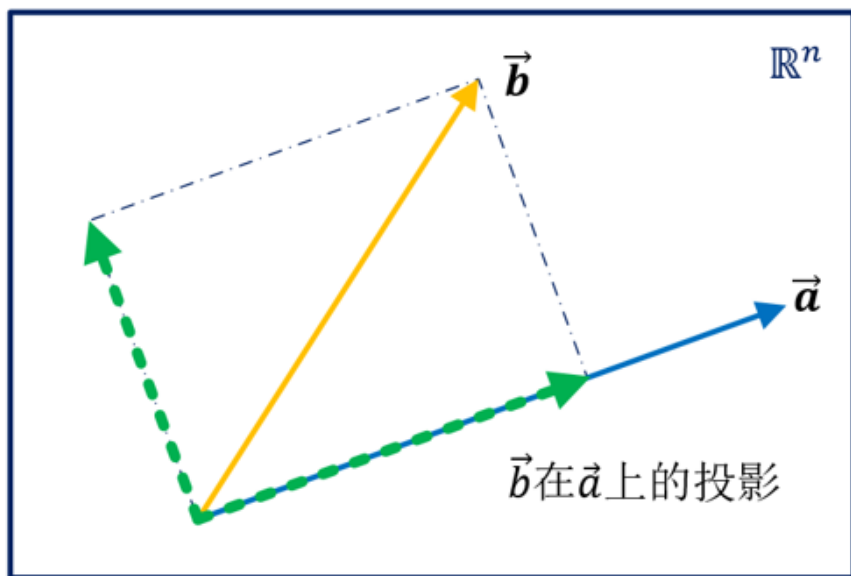
可以发现 $B^*$ 即为最优解，这是因为目标函数为 $B$ 的二次函数，Taylor展开到二阶即为精确展开。一次牛顿法迭代就可得到最优解。

# 线性回归最优解的几何解释

考虑在 $n$ 维空间 $\mathbb{R}^n$ 中的向量 $\vec{a} \in \mathbb{R}^n$ , 对于任意给定的向量 $\vec{b} \in \mathbb{R}^n$ , 则向量 $\vec{b}$ 在向量 $\vec{a}$ 方向上的长度矢量, 称为 $\vec{b}$ 在 $\vec{a}$ 方向上的投影。



# 线性回归最优解的几何解释



由与  $\vec{b} - \mu\vec{a}$  部分与向量  $\mu\vec{a}$  是正交的，则有

$$\mu\vec{a}^T(\vec{b} - \mu\vec{a}) = 0$$

进而有

$$\vec{a}^T(\vec{b} - \mu\vec{a}) = 0$$

$$\vec{a}^T\vec{b} - \vec{a}^T\mu\vec{a} = 0$$

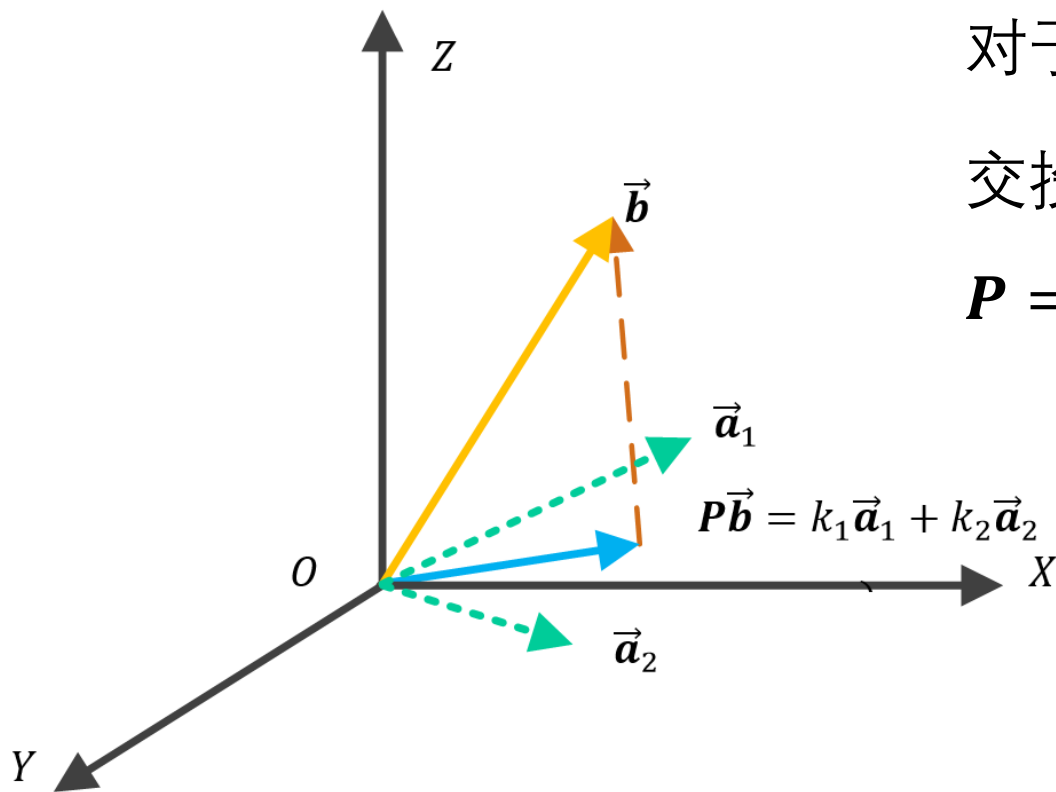
$$\vec{a}^T\vec{b} = \vec{a}^T\mu\vec{a} = \mu\vec{a}^T\vec{a}$$

$$\mu = (\vec{a}^T\vec{a})^{-1}\vec{a}^T\vec{b}$$

因此，从  $P\vec{b} = \mu\vec{a}$  可得

$$P = a(a^T a)^{-1}a^T$$

# 线性回归最优解的几何解释



对于任意的向量 $\vec{b}$ ，若要求出其向空间 $A$ 的正交投影结果，则就是对其做变换 $P\vec{b}$ ，其中 $P = A(A^T A)^{-1} A^T$ 。

# 线性回归最优解

- 最小二乘解就是在寻找本体位置和投影位置之差的二范数最小的解。最小二乘问题的目标函数为

$$\hat{\mathbf{x}} = \underset{\mathbf{x}}{\operatorname{argmin}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$$

- 通过正交投影的知识可以知道，对正交投影  $\mathbf{P} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ ， $\mathbf{y}$  在  $\mathbf{A}$  的列空间中的投影为

$$\mathbf{P}\mathbf{y} = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

- 可得

$$\hat{\mathbf{x}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$

- 如果采用伪逆的写法则记为

$$\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{y}$$

这就是从正交投影的理解方式，导出的最小二乘解。



# 线性回归问题的扩展

- 加权最小二乘法
- 偏最小二乘法
- 正则项（惩罚项）
- 非线性扩展

谢谢各位同学!