



大数据导论

Introduction to Big Data



第7讲 回归预测：基础概念与算法

叶允明

计算机科学与技术学院

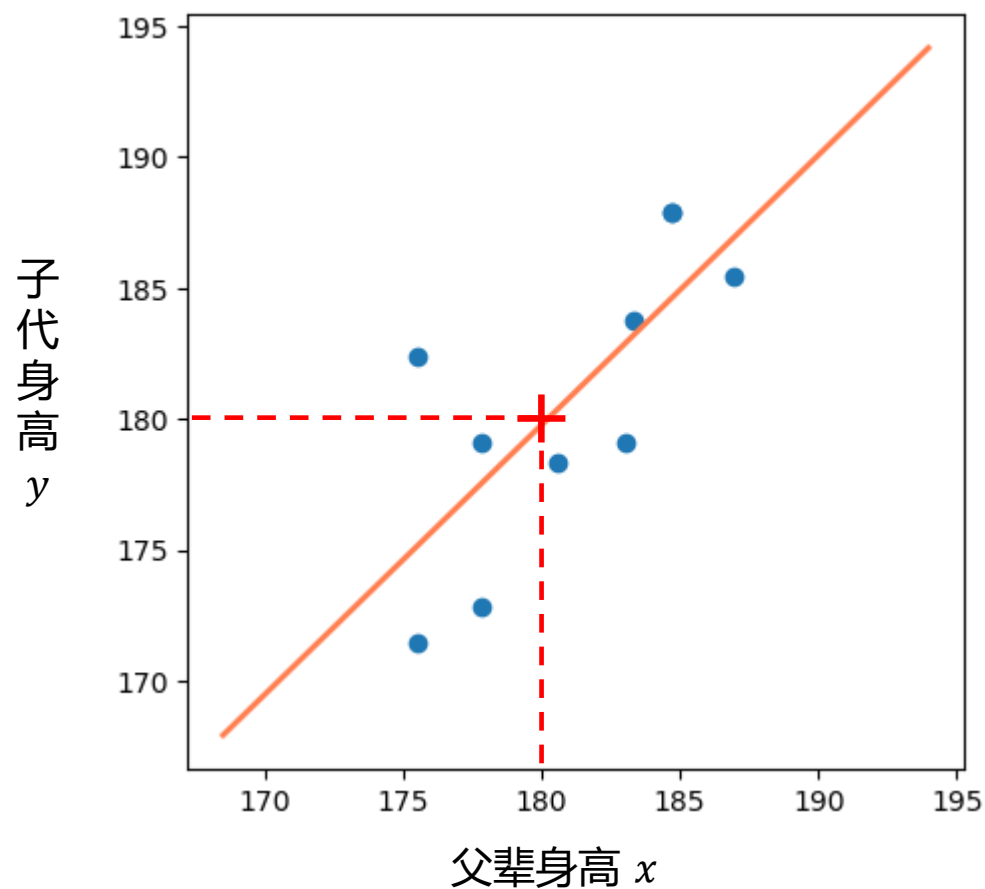
哈尔滨工业大学（深圳）

目录

- 回归预测的基本概念
- 线性回归
- 基于梯度下降法的线性回归
- 基于最小二乘法的线性回归 (optional)

回归预测任务

- 回归 (regression) : 预测给定数据对象的**目标值** (与分类的**类别**对应)



$x = 180$ 时, $y = ?$

回归预测的应用领域

- 几乎每个人工智能应用领域都涉及到预测问题
 - 股票预测
 - 贷款额度估计
 - 视频预测
 - 销售业绩预测
 - 医学诊断
 - 欺诈检测
 -

回归任务的定义

- 回归任务可以用一个形式化函数表示：

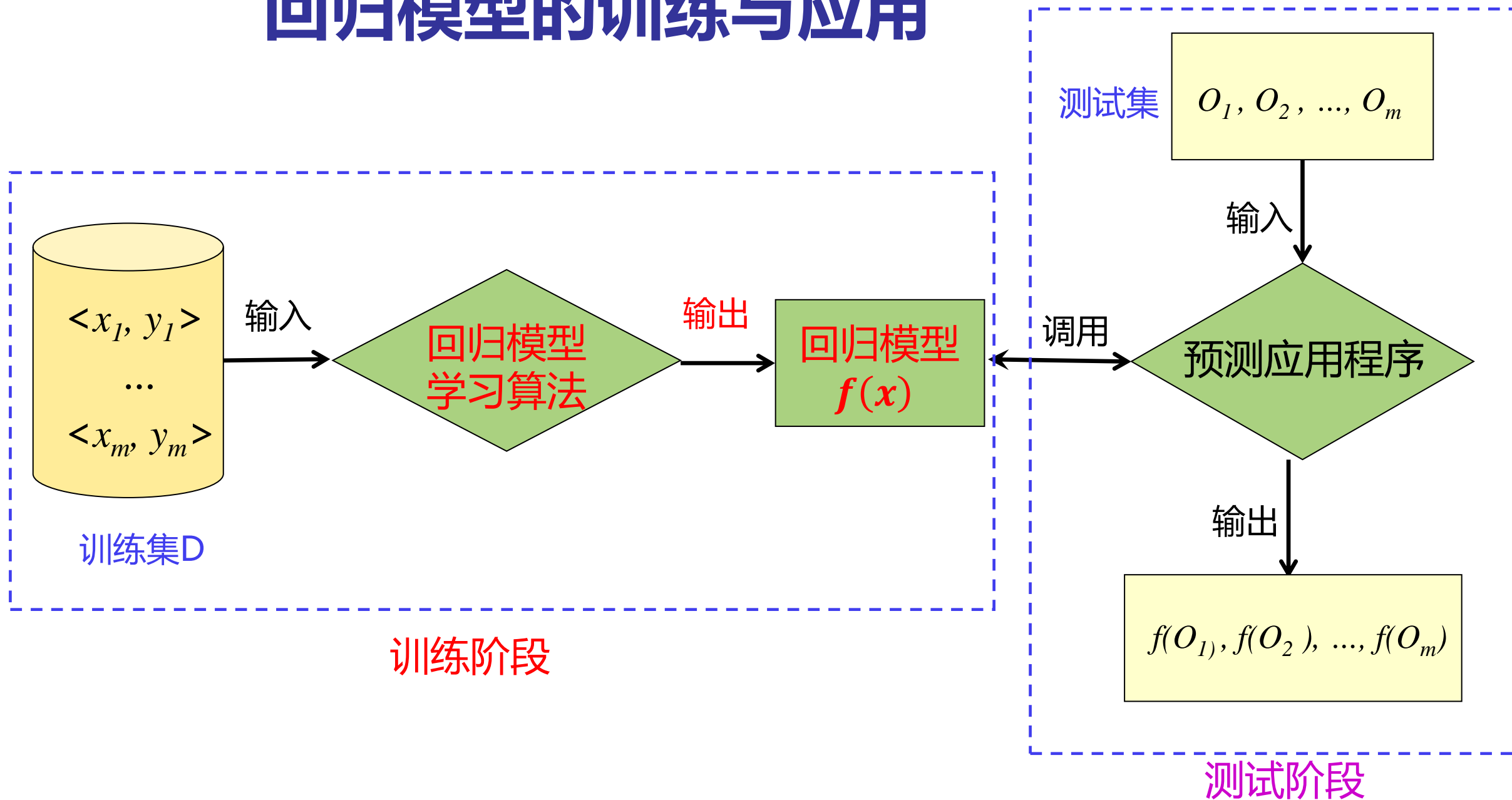
$$y = f(x),$$

其中 $x \in \mathbb{D}, y \in \mathbb{R}$

- 回归函数 $f(x)$ 经过运算可以输出一个连续的实数值 y ，即“回归模型”

如何构造回归函数 $f(x)$ 呢？

回归模型的训练与应用



线性回归

应用案例

- 预测下一天的日均气温 $y = f(\mathbf{x})$, 其中 $\mathbf{x} \in \mathbb{D}, y \in \mathbb{R}$
 - Kaggle数据集 (daily climate time series data)

日期	日均气温 (mean temp)	相对湿度 (humidity)	风速 (wind speed)	气压 (pressure)
...
2017-01-02	7.40	92.00	2.980	1017.80
2017-01-03	7.17	87.00	4.63	1018.67

\mathbf{x}

f

y

线性回归模型

- 给定数据集 $\mathbb{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$, 其中, 每个样本 \mathbf{x}_i 有 d 个特征:

$$\mathbf{x}_i = (x_{i,1}; x_{i,2}; \dots; x_{i,d})^T, \quad y_i \in \mathbb{R}.$$

- 线性模型的目的是学习一个关于 \mathbf{x} 的线性函数 $f(\mathbf{x})$, 来尽可能准确地预测 y

$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b, \quad \mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_d \end{pmatrix}$$
$$y \approx f(\mathbf{x})$$

线性回归模型的优劣评价方法

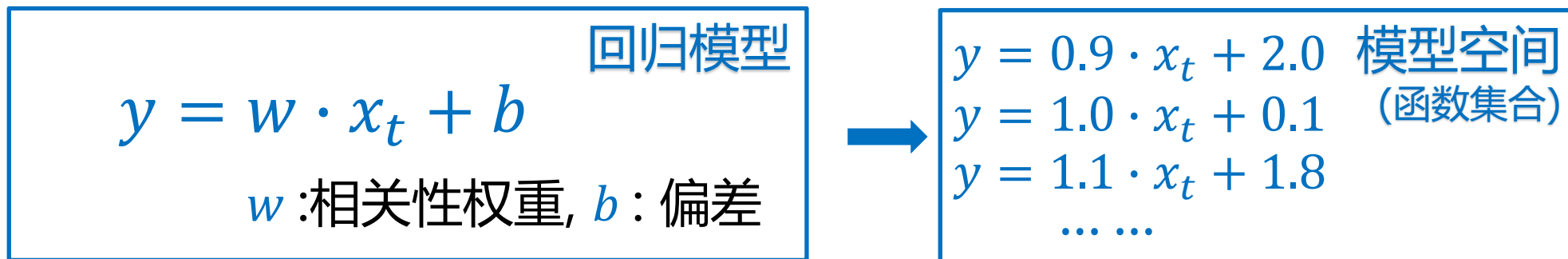
- 在求解 \mathbf{w} 和 b 之前，需要给出衡量 $f(\mathbf{x})$ 与 y 之间误差的损失函数
- 在回归方法中，均方误差是比较常用的损失函数，其损失函数定义如下：

$$L(\mathbf{w}, b) = \frac{1}{2n} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2$$

单变量线性回归

- 第一步：确定模型空间

➤ 直观想法：下一天的日均气温很可能与前一天的日均气温相关



线性回归模型: $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$

数据: (\mathbf{x}_i, y_i)

	日均气温 (mean temp)	相对湿度 (humidity)	风速 (wind speed)	气压 (pressure)
	7.40	92.00	2.980	1017.80

$\mathbf{x}_i = (x_{i,t}, x_{i,h}, x_{i,w}, x_{i,p})^T$

损失函数

- 第二步：模型优劣的评价标准

$$L(w, b) = \sum_{i=1}^n \left[y_i - \underbrace{(w \cdot x_{i,t} + b)}_{\text{预测值}} \right]^2$$

预测误差

模型空间
(函数集合)

$$\begin{aligned} y &= 0.9 \cdot x_t + 2.0 \\ y &= 1.0 \cdot x_t + 0.1 \\ y &= 1.1 \cdot x_t - 1.8 \\ &\dots \end{aligned}$$

- 为了便于计算，随机选取了三十天的数据作为训练集(单位, °C)

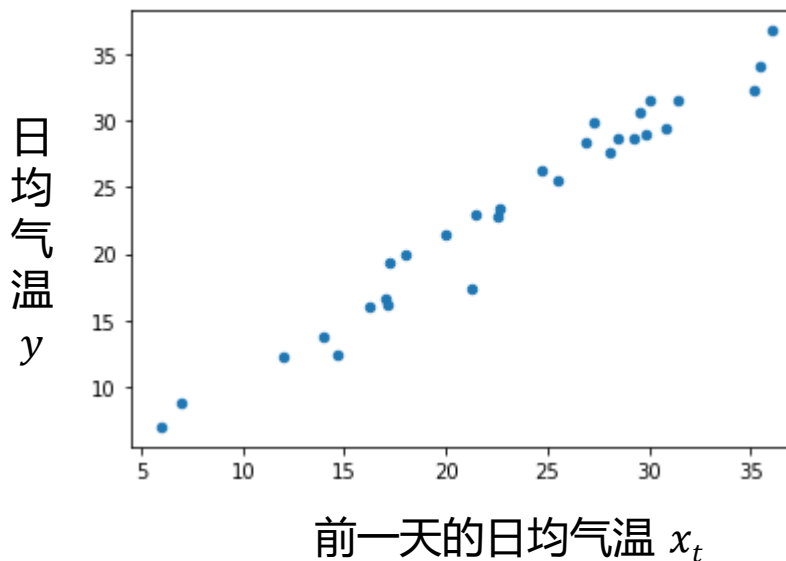
$(x_{1,t}, y_1)$

$(x_{2,t}, y_2)$

\vdots

$(x_{30,t}, y_{30})$

散点图

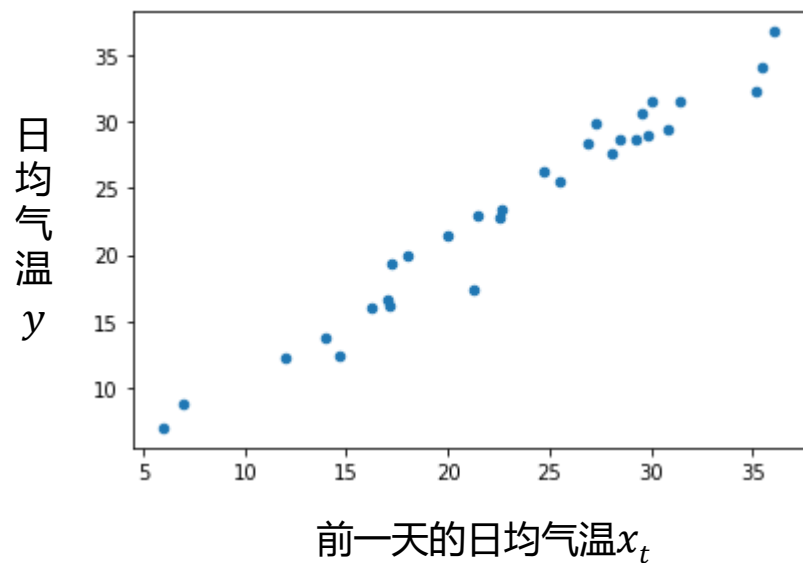


损失函数

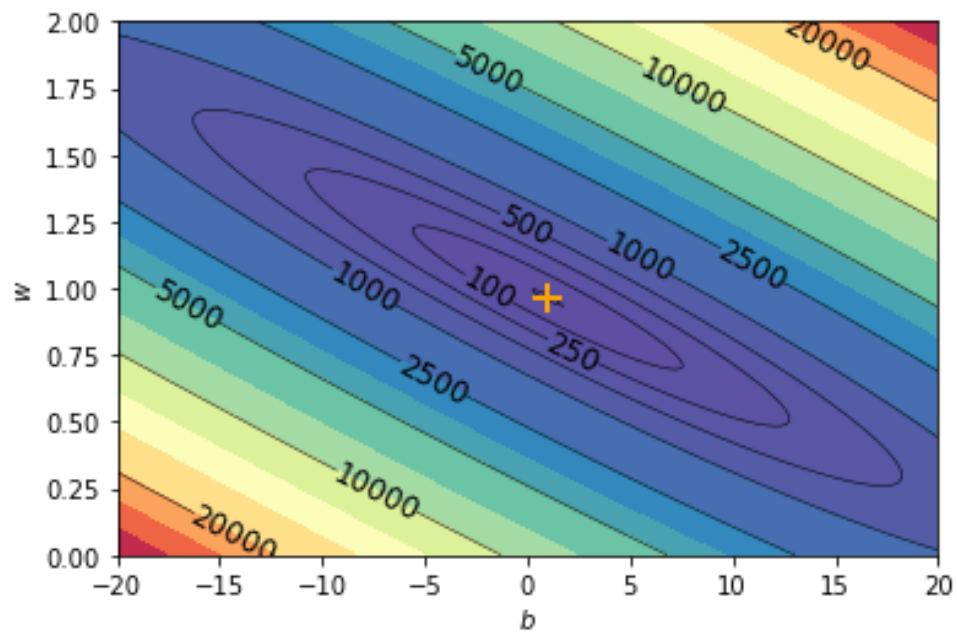
- 损失函数等值线图

$$L(w, b) = \sum_{i=1}^n (y_i - w \cdot x_{i,t} - b)^2$$

如何寻找最小值点？



模型空间



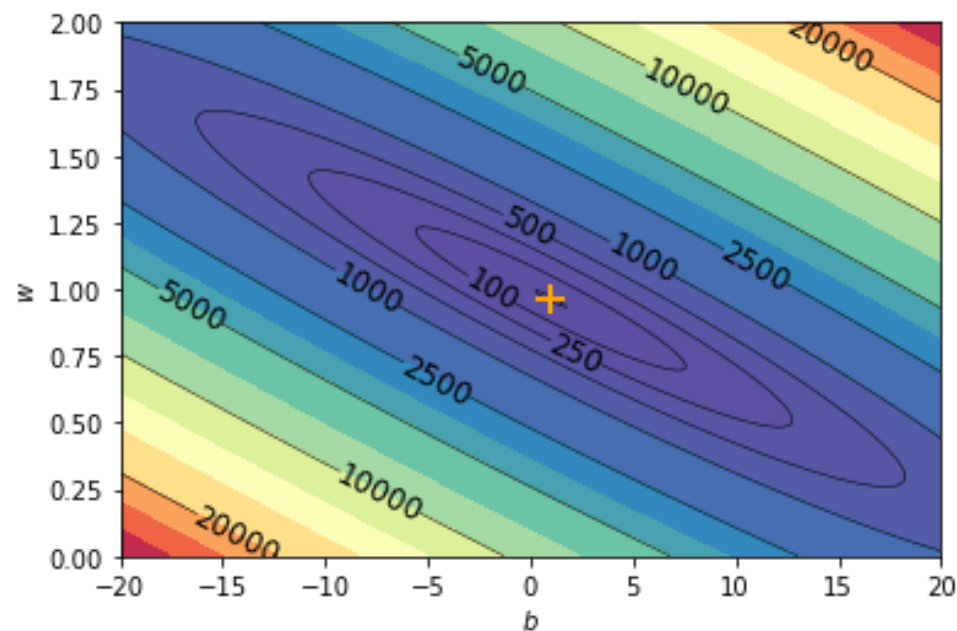
最优化算法

- 第三步：找到“最优”模型
- 寻找一组最优的参数 (w^*, b^*) ，能够最小化损失函数，即

$$\begin{aligned} w^*, b^* &= \arg \min_{b, w} L(w, b) \\ &= \arg \min_{b, w} \sum_{i=1}^n [y_i - (w \cdot x_{i,t} + b)]^2 \end{aligned}$$

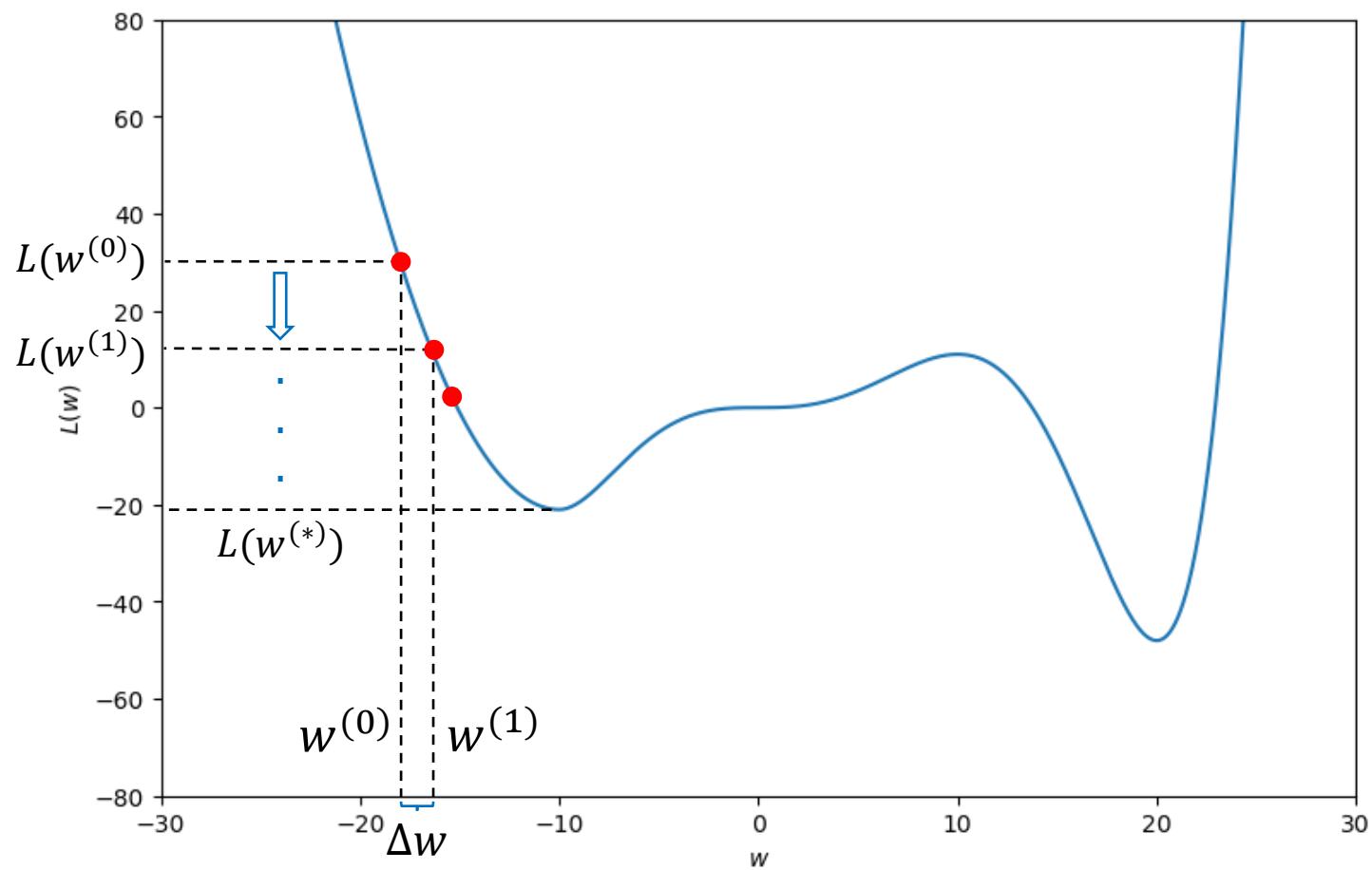
- 常用算法：

- 梯度下降法
- 最小二乘法（小规模数据）



基于梯度下降法的线性回归模型

梯度下降法的基本思想



梯度下降法的数学原理？

梯度下降法的算法流程（单个参数）

- 以一个仅有单个参数 w 的光滑的损失函数 $L(w)$ 为例，梯度下降法：

随机选取初始值 $w^{(0)}$

$$w^{(1)} \leftarrow w^{(0)} - \eta \frac{dL}{dw} \Big|_{w=w^{(0)}}$$

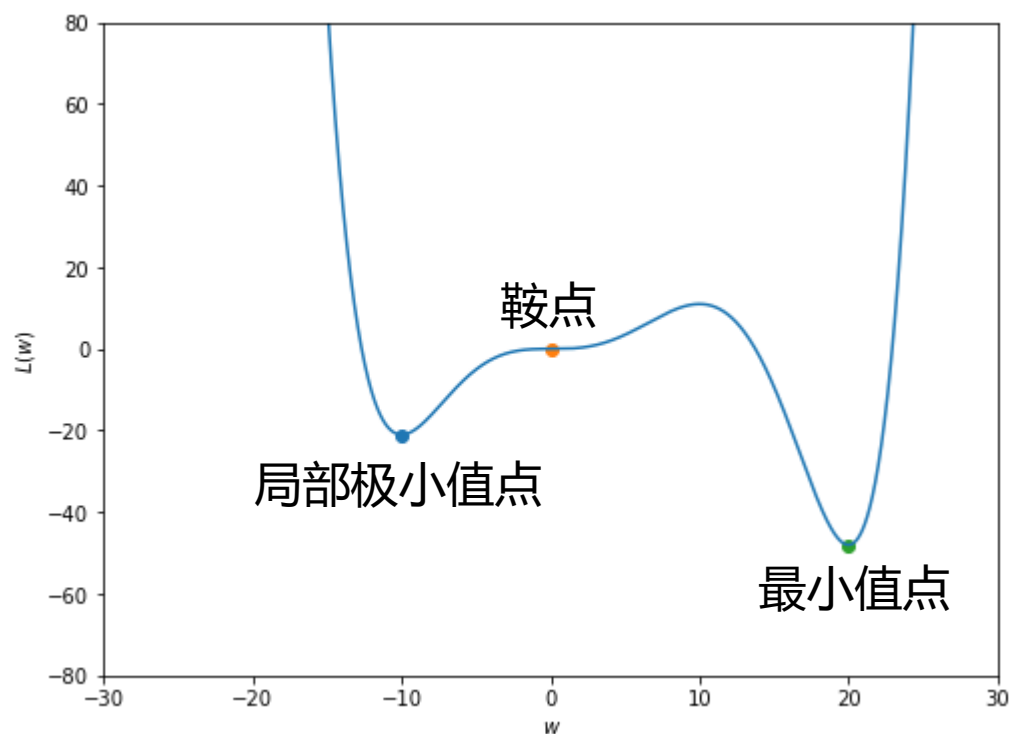
$$w^{(2)} \leftarrow w^{(1)} - \eta \frac{dL}{dw} \Big|_{w=w^{(1)}}$$

...

$$w^{(j+1)} \leftarrow w^{(j)} - \eta \frac{dL}{dw} \Big|_{w=w^{(j)}}$$

...

直到 $|w^{(n+1)} - w^{(n)}| < \varepsilon$, ε 为终止条件



梯度下降法的算法流程（两个参数）

- 案例的损失函数 $L(w, b)$ 含有两个参数 $L(w, b) = \sum_{i=1}^N (y_i - w \cdot x_{i,t} - b)^2$

➤ （随机）选取两个初始值 $b^{(0)}, w^{(0)}$

➤ 计算 $\frac{\partial L}{\partial w} \big|_{w=w^{(0)}, b=b^{(0)}}$, $\frac{\partial L}{\partial b} \big|_{w=w^{(0)}, b=b^{(0)}}$, 更新 b, w

$$\nabla L(b, w) = \begin{pmatrix} \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial w} \end{pmatrix}$$

$$w^{(1)} \leftarrow w^{(0)} - \eta \frac{\partial L}{\partial w} \big|_{w=w^{(0)}, b=b^{(0)}}, \quad b^{(1)} \leftarrow b^{(0)} - \eta \frac{\partial L}{\partial b} \big|_{w=w^{(0)}, b=b^{(0)}}$$

➤ 计算 $\frac{\partial L}{\partial w} \big|_{w=w^{(1)}, b=b^{(1)}}$, $\frac{\partial L}{\partial b} \big|_{w=w^{(1)}, b=b^{(1)}}$, 更新 b, w

$$w^{(2)} \leftarrow w^{(1)} - \eta \frac{\partial L}{\partial w} \big|_{w=w^{(1)}, b=b^{(1)}}, \quad b^{(2)} \leftarrow b^{(1)} - \eta \frac{\partial L}{\partial b} \big|_{w=w^{(1)}, b=b^{(1)}}$$

... ..

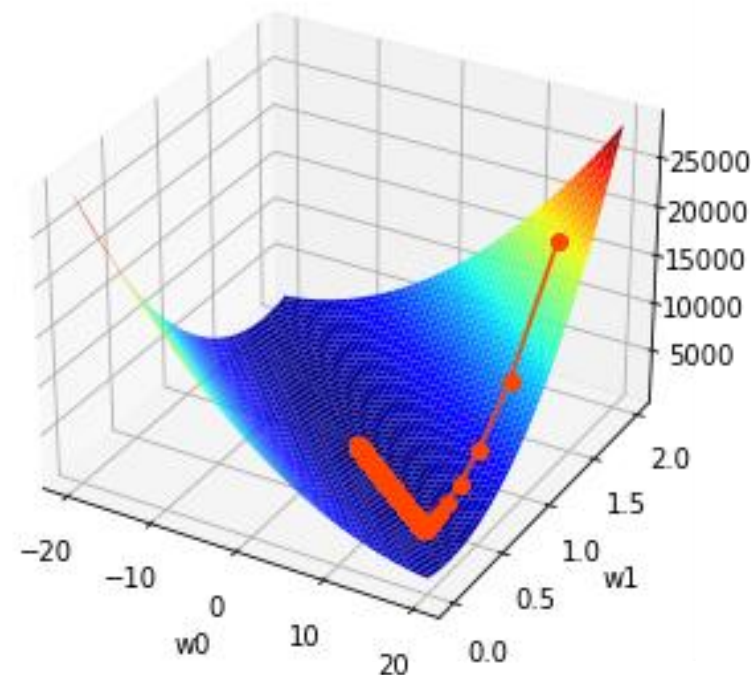
梯度下降法的应用示例

- 计算梯度

$$L(w, b) = \sum_{i=1}^n (y_i - \underline{w \cdot x_{i,t}} - \underline{b})^2$$

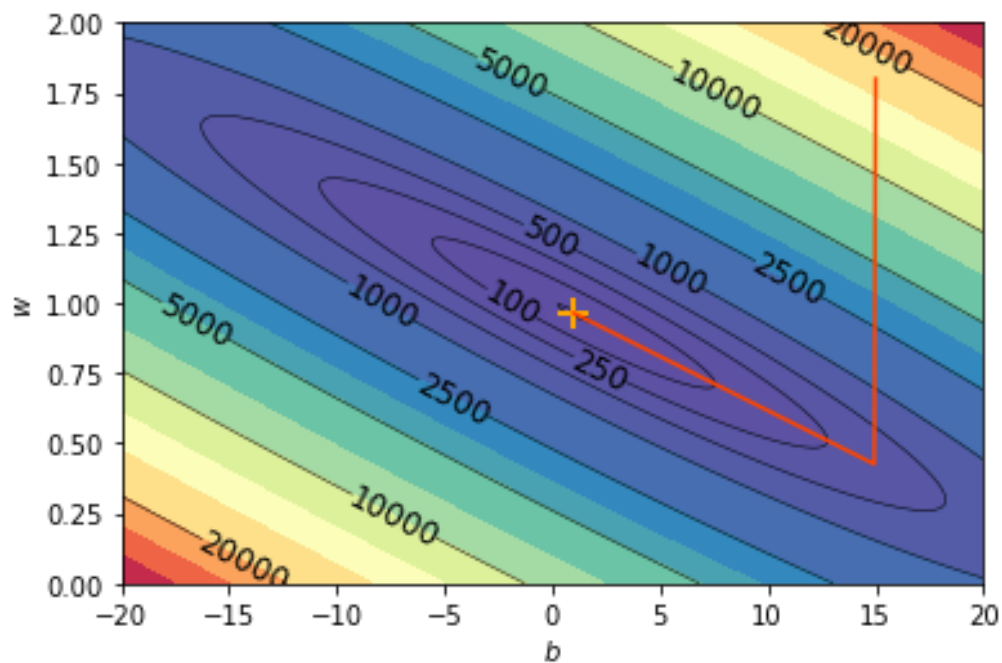
$$\frac{\partial L}{\partial w} = \sum_{i=1}^n 2(y_i - w \cdot x_{i,t} - b)(-x_{i,t})$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^n 2(y_i - w \cdot x_{i,t} - b)(-1)$$

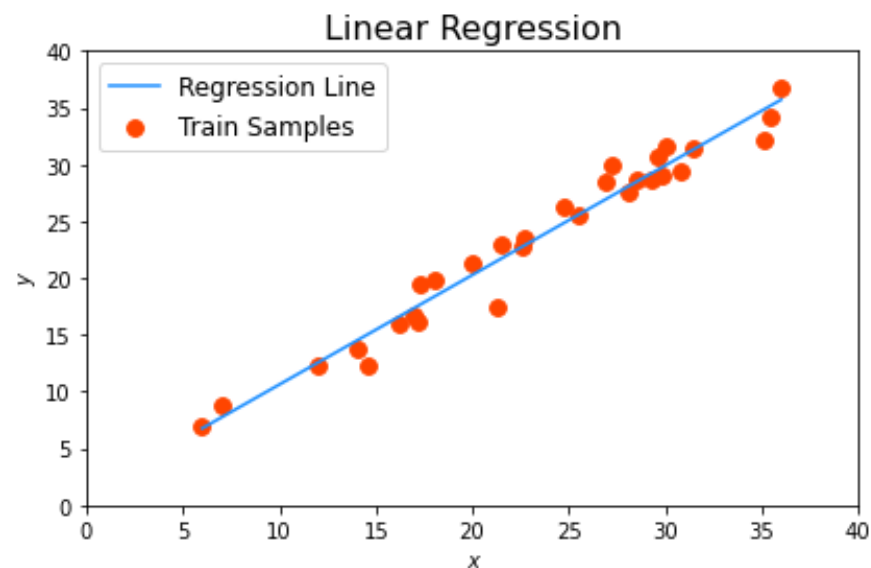


梯度下降法示例

- 沿梯度下降的方向寻找极小值



$$w^* = 0.965, b^* = 0.980$$



$$y = 0.965x_t + 0.980$$

测试模型效果

- 训练集均方误差

$$\frac{1}{30} \sum_{i=1}^{30} (y_i - w^* \cdot x_{i,t} - b^*)^2 = 2.134$$

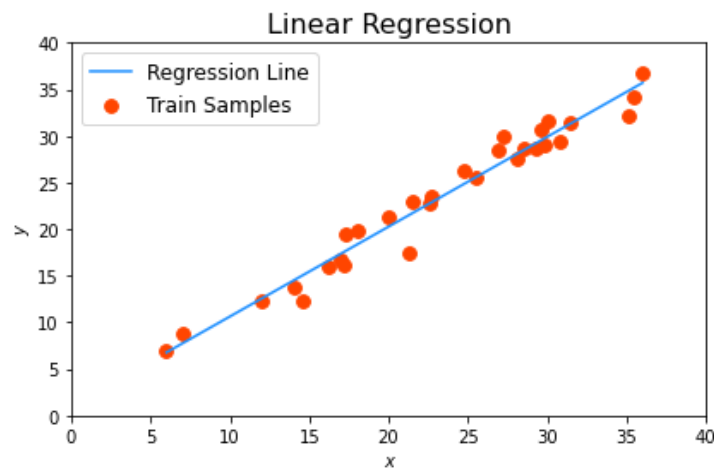
- 在测试数据中随机选取十天的数据
作为测试集，测试模型的泛化性能

- 测试集均方误差

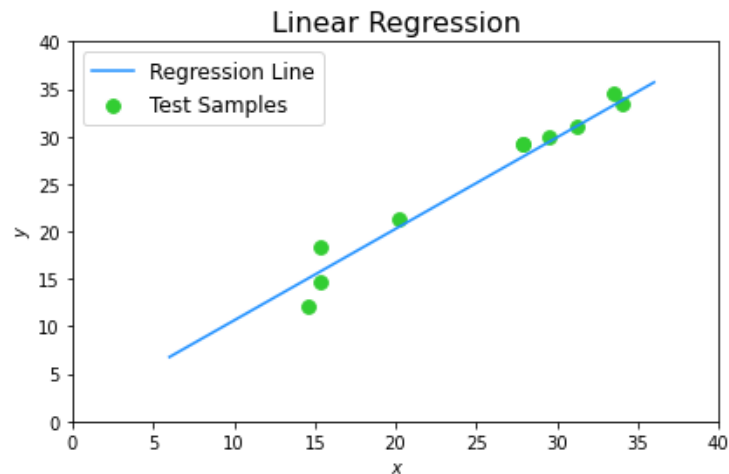
$$\frac{1}{10} \sum_{i=1}^{10} (y_i - w^* \cdot x_{i,t} - b^*)^2 = 2.294$$

- 测试集平均误差

$$\frac{1}{10} \sum_{i=1}^{10} |y_i - w^* \cdot x_{i,t} - b^*| = 1.229$$



$$w^* = 0.965, b^* = 0.980$$



改进模型：增加二次项特征

- 回归模型：

$$y = w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

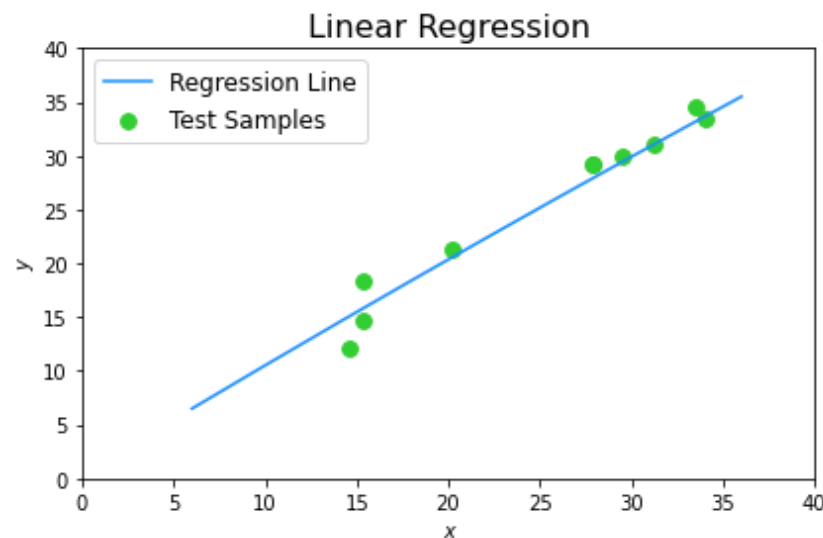
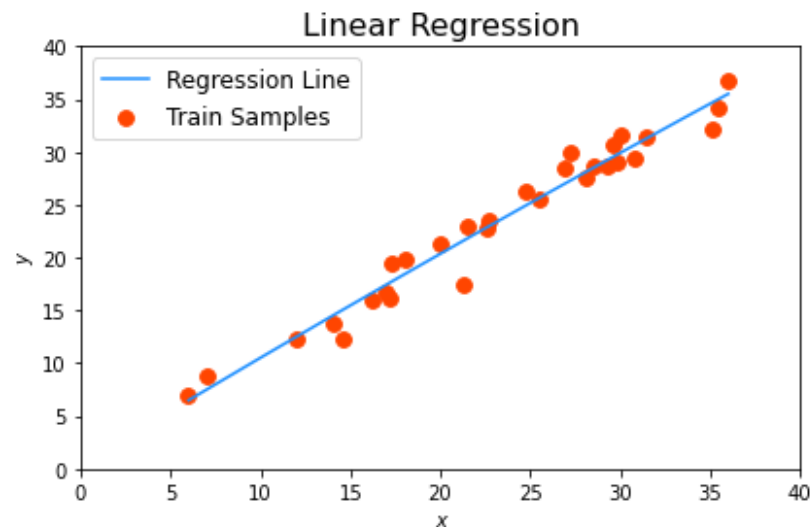
- 计算模型参数，可得

$$w_2^* = -1.50 \times 10^{-3}, w_1^* = 1.030, b^* = 0.361$$

- 模型误差：

- 训练集均方误差： $2.123 < 2.134$

- 测试集均方误差： $2.278 < 2.294$



改进模型：增加三次项特征

- 回归模型：

$$y = w_3 \cdot x_t^3 + w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

- 计算模型参数，可得

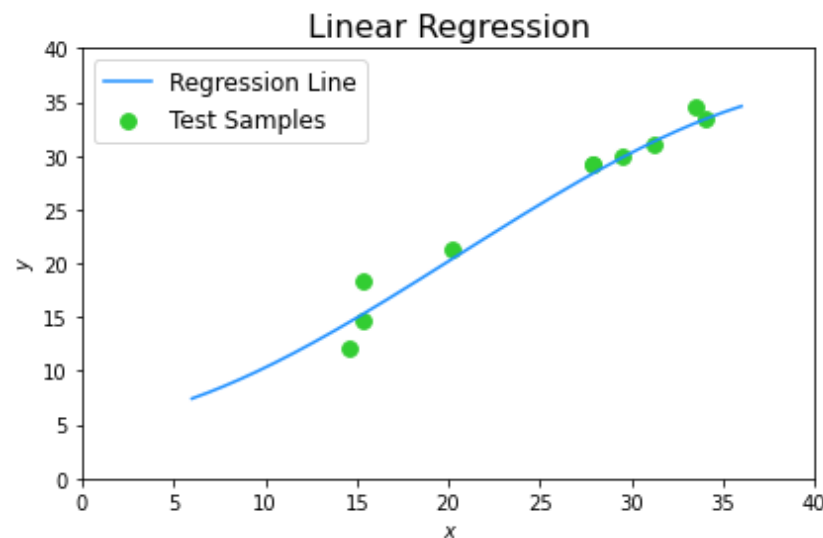
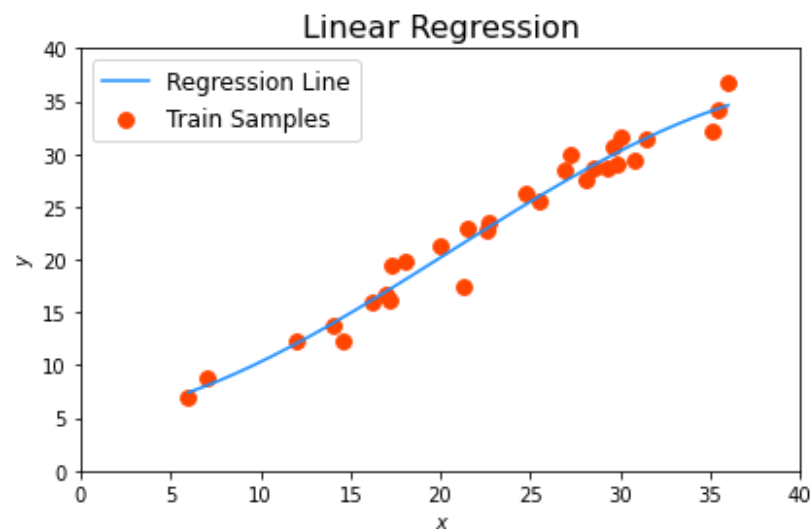
$$w_3^* = -7.43 \times 10^{-4}, w_2^* = 0.046,$$

$$w_1^* = 0.136, b^* = 5.123$$

- 模型误差：

- 训练集均方误差： $1.913 < 2.123$

- 测试集均方误差： $2.042 < 2.278$



改进模型：增加四次项特征

- 回归模型：

$$y = w_4 \cdot x_t^4 + w_3 \cdot x_t^3 + w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

- 计算模型参数，可得

$$w_4^* = 4.75 \times 10^{-5}, w_3^* = -4.83 \times 10^{-3},$$

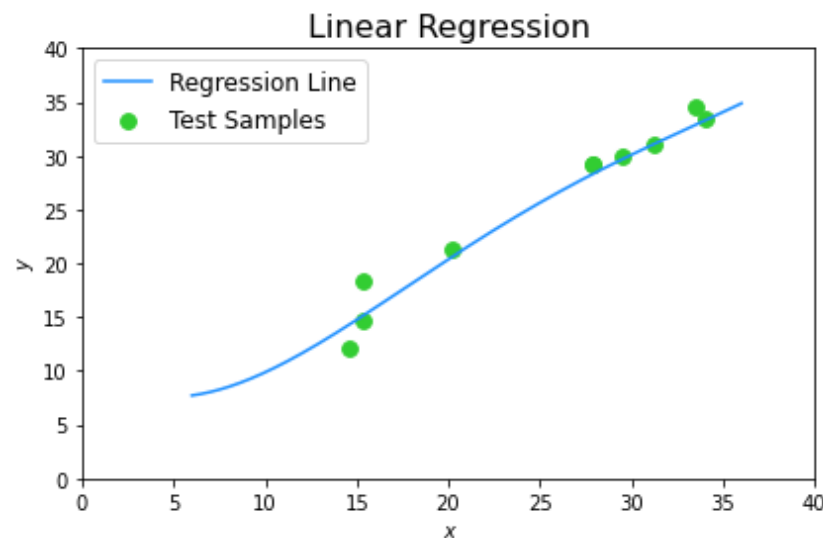
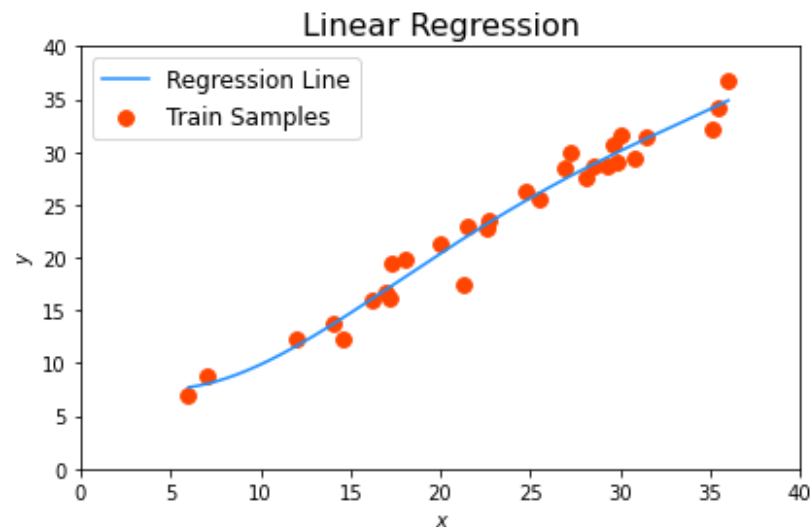
$$w_2^* = 0.167, w_1^* = -1.290, b^* = 10.43$$

- 模型误差：

➤ 训练集均方误差：1.878 < 1.913

➤ 测试集均方误差：2.053 > 2.042

过拟合



改进模型：增加五次项特征

- 回归模型：

$$y = w_5 \cdot x_t^5 + w_4 \cdot x_t^4 + w_3 \cdot x_t^3 + w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

- 计算模型参数，可得

$$w_5^* = 1.184 \times 10^{-5}, w_4^* = -1.22 \times 10^{-3}, w_3^* = 4.64 \times 10^{-2},$$

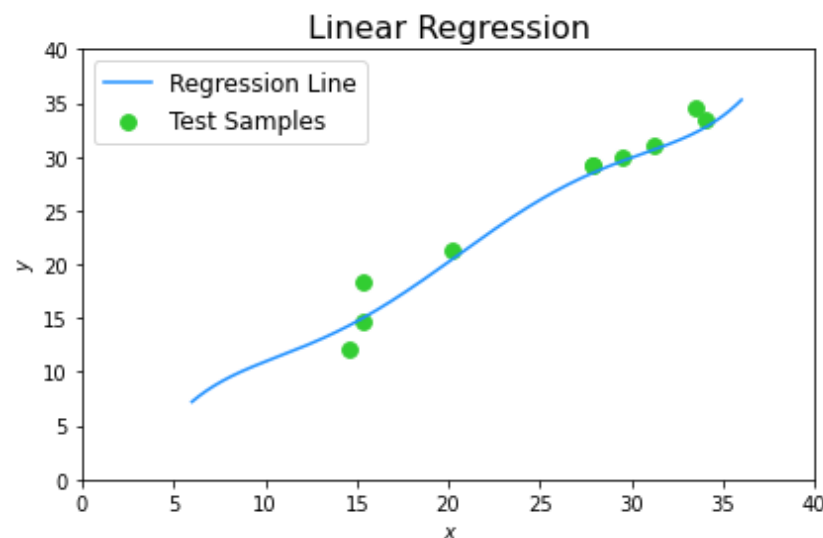
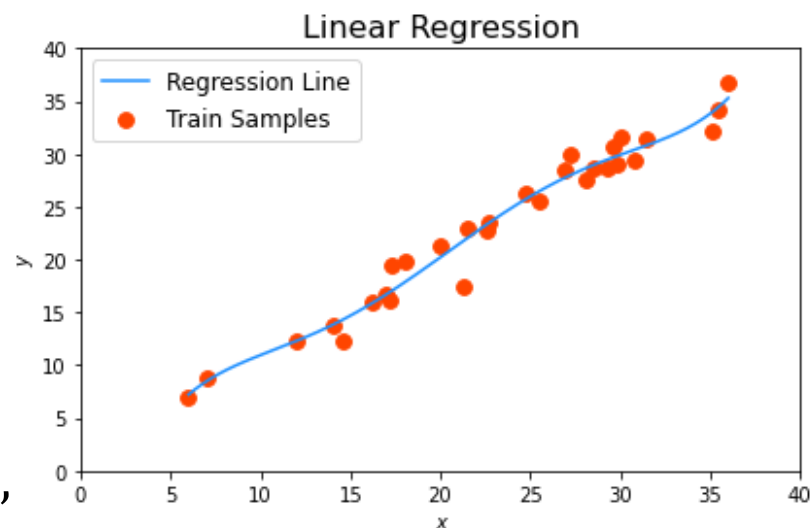
$$w_2^* = -0.080, w_1^* = 6.948, b^* = -14.37$$

- 模型误差：

➤ 训练集均方误差： $1.797 < 1.878$

➤ 测试集均方误差： $2.396 > 2.053$

过拟合



过拟合

- 复杂的模型可以更好地拟合训练数据
- 但未必会在测试数据上获得更好的效果

$$y = w \cdot x_t + b$$

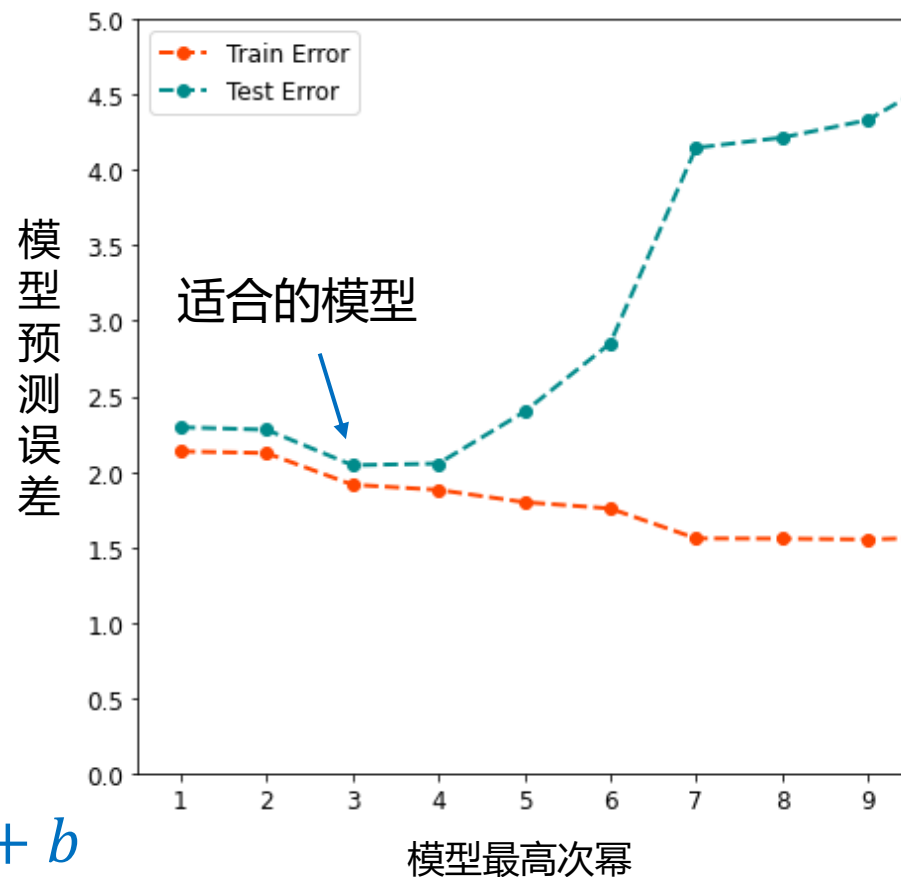
$$y = w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

$$y = w_3 \cdot x_t^3 + w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

$$y = w_4 \cdot x_t^4 + w_3 \cdot x_t^3 + w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

$$y = w_5 \cdot x_t^5 + w_4 \cdot x_t^4 + w_3 \cdot x_t^3 + w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

... ..



多变量线性回归

- 回到第一步：确定模型空间

- 除前一天的日均气温外，考虑与前一天的相对湿度、风速、气压是否相关

$$y = w_1 \cdot x_t + w_2 \cdot x_t^2 + w_3 \cdot x_t^3 + w_4 \cdot x_h + w_5 \cdot x_w + w_6 \cdot x_p + b$$

线性回归模型： $f(\mathbf{x}) = \mathbf{wx} + b$

数据： (\mathbf{x}_i, y_i)

	日均气温 (mean temp)	相对湿度 (humidity)	风速 (wind speed)	气压 (pressure)
	7.40	92.00	2.980	1017.80
\mathbf{x}_i	$x_{i,t}$	$x_{i,h}$	$x_{i,w}$	$x_{i,p}$

增加其他特征变量

- 回归模型:

$$y = w_1 \cdot x_t + w_2 \cdot x_t^2 + w_3 \cdot x_t^3 \\ + w_4 \cdot x_h + w_5 \cdot x_w + w_6 \cdot x_p + b$$

- 计算模型参数, 可得

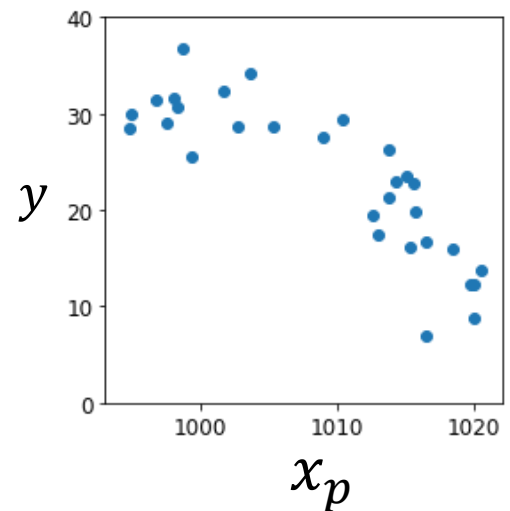
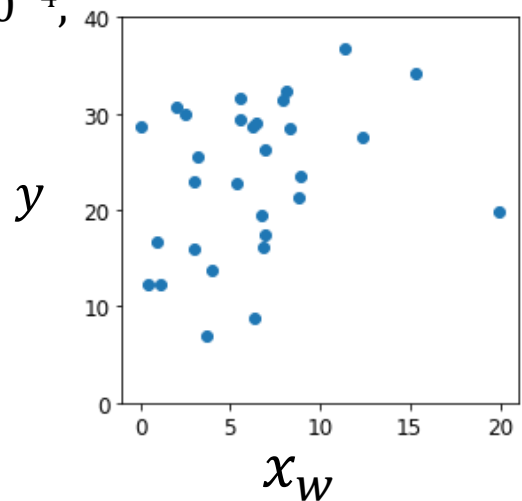
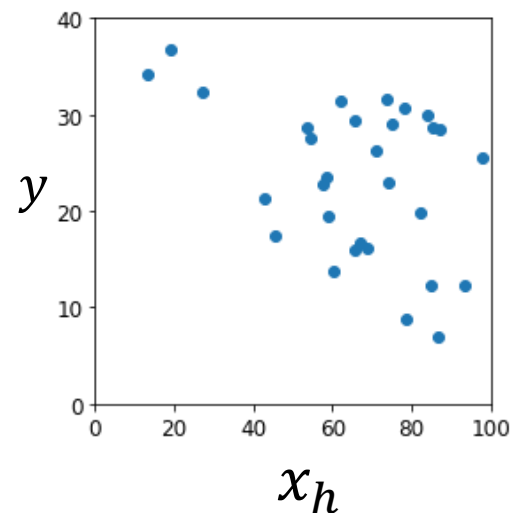
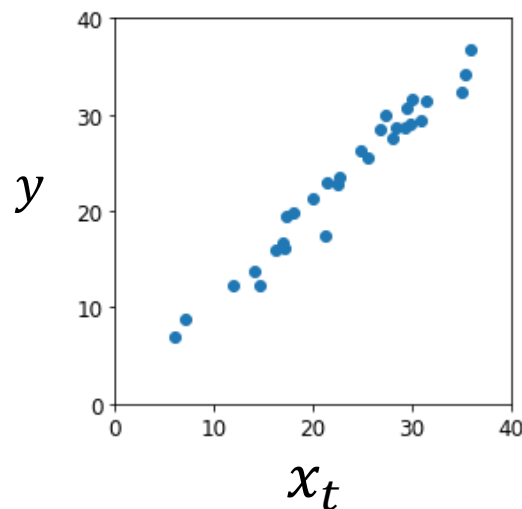
$$w_6^* = -0.011, w_5^* = 0.010, w_4^* = 1.18 \times 10^{-2}, w_3^* = -2.58 \times 10^{-4}, \\ w_2^* = 1.39 \times 10^{-2}, w_1^* = 0.667, b^* = 109.5$$

- 模型误差:

➤ 训练集均方误差: $1.553 < 1.913$

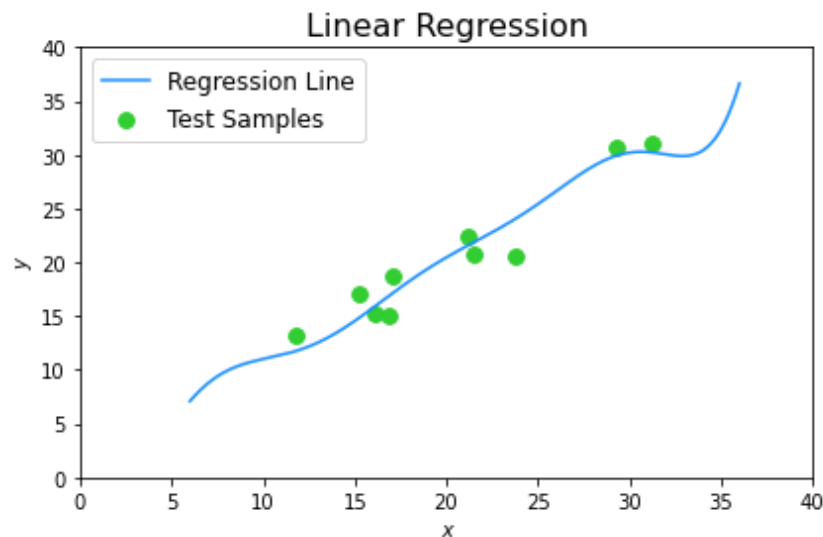
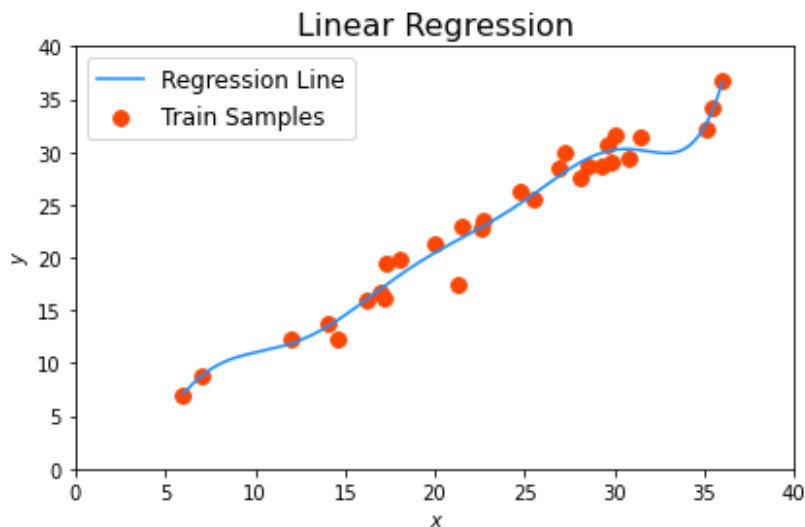
➤ 测试集均方误差: $2.278 > 2.042$

过拟合



正则化的基本思想

- 奥卡姆剃刀(Occam's razor)原理
 - 选择能够很好地解释已知数据但更简单的模型。
- 简单的函数更为平滑，也就不容易发生过拟合的问题。



$$y = w_{10} \cdot x_t^{10} + w_9 \cdot x_t^9 + \cdots + w_2 \cdot x_t^2 + w_1 \cdot x_t + b$$

正则化

- 对于有 d 个特征的线性回归模型 $y = \sum_{j=1}^d w_j x_j + b$

- 损失函数为

$$L(w, b) = \sum_{i=1}^n \left(y_i - \sum_{j=1}^d w_j x_{i,j} - b \right)^2$$

- 加入正则化项的损失函数

λ 为超参数, λ 值越大,
模型抗扰动能力越强

- L1正则化

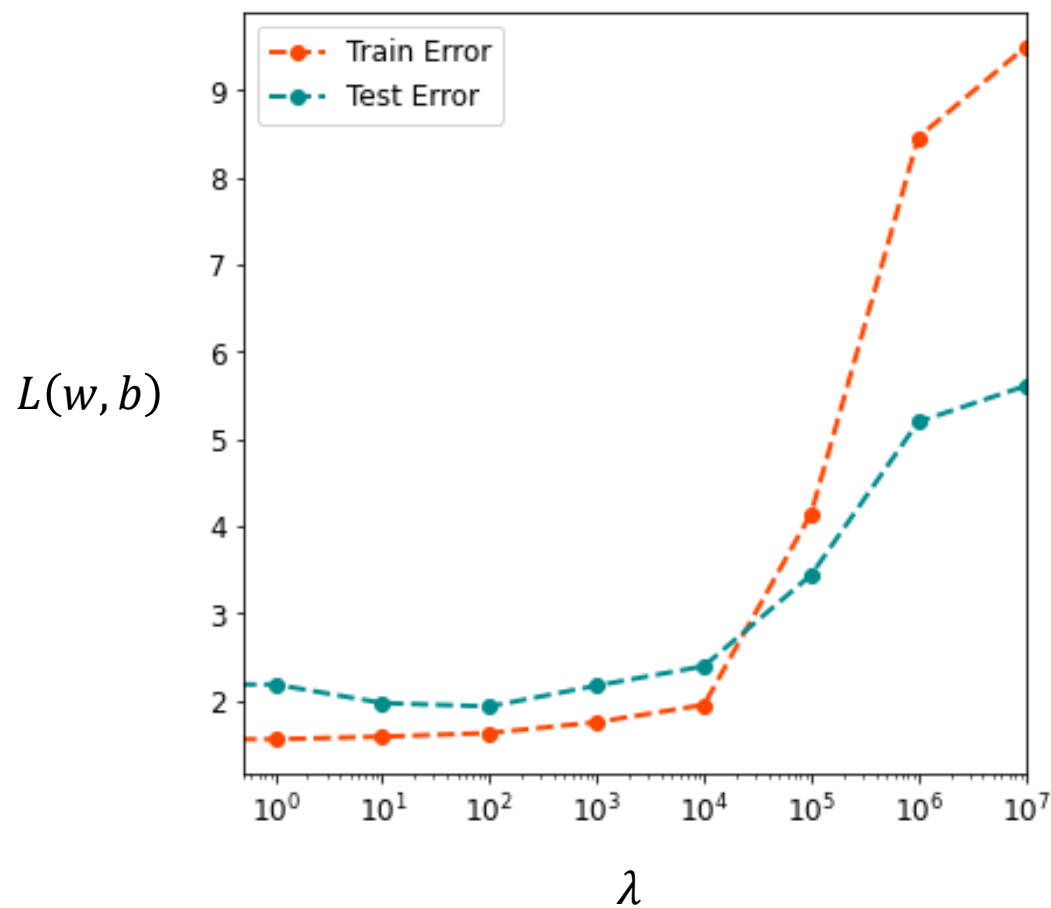
$$L(w, b) = \sum_{i=1}^n \left(y_i - \sum_{j=1}^d w_j x_{i,j} - b \right)^2 + \lambda \sum_{j=1}^d |w_j|$$

- L2正则化

$$L(w, b) = \sum_{i=1}^n \left(y_i - \sum_{j=1}^d w_j x_{i,j} - b \right)^2 + \lambda \sum_{j=1}^d |w_j|^2$$

正则化

- 对于模型: $y = w_1 \cdot x_t + w_2 \cdot x_t^2 + w_3 \cdot x_t^3 + w_4 \cdot x_h + w_5 \cdot x_w + w_6 \cdot x_p + b$
- 损失函数中加入 L2 正则化项
 - $\lambda = 0$ 时:
 - ✓ 训练集均方误差: 1.553
 - ✓ 测试集均方误差: 2.278
 - $\lambda = 100$ 时:
 - ✓ 训练集均方误差: 1.627
 - ✓ 测试集均方误差: 1.932



基于最小二乘法的线性回归模型

基于最小二乘（Least Square）法的线性回归模型

- 将数据样本用 $n \times (d + 1)$ 大小的矩阵 \mathbf{X} 表示：

$$\mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & \cdots & x_{1,d} \\ 1 & x_{2,1} & \cdots & x_{2,d} \\ 1 & \cdots & \cdots & \cdots \\ 1 & x_{n,1} & \cdots & x_{n,d} \end{pmatrix}$$

损失函数的向量形式

- 数据集的目标值序列也可写成向量形式 $\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix}$
- 为了简洁起见，将待求解参数 \mathbf{w} 和 b 合并为 $\mathbf{W} = \begin{pmatrix} b \\ \mathbf{w} \end{pmatrix}$
- 于是，线性回归的损失函数可重写如下：

$$L(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 = \frac{1}{2} (\mathbf{XW} - \mathbf{y})^T (\mathbf{XW} - \mathbf{y})$$

最小二乘法的算法流程

- 我们的目标是寻找一组最优的参数 $\mathbf{W}^* = \begin{pmatrix} b^* \\ \mathbf{w}_* \end{pmatrix}$, 能够最小化损失函数:

$$\operatorname{argmin}_{\mathbf{W}} L(\mathbf{W}) \qquad L(\mathbf{W}) = \frac{1}{2} (\mathbf{X}\mathbf{W} - \mathbf{y})^T (\mathbf{X}\mathbf{W} - \mathbf{y})$$

- 当 $X^T X$ 可逆时, 对其求导并求驻点可以直接求得 \mathbf{W} 的最优解:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}} &= \frac{1}{2} \frac{\partial}{\partial \mathbf{W}} (\mathbf{W}^T X^T X \mathbf{W} - \mathbf{W}^T X^T \mathbf{y} - \mathbf{y}^T X \mathbf{W} + \mathbf{y}^T \mathbf{y}) \\ &= \frac{1}{2} (2X^T X \mathbf{W} - X^T \mathbf{y} - X^T \mathbf{y}) \\ &= X^T X \mathbf{W} - X^T \mathbf{y} \end{aligned}$$

- \mathbf{W} 的解析解: $\mathbf{W} = (X^T X)^{-1} X^T \mathbf{y}$

最小二乘法的应用示例

- 单变量线性回归模型: $y = w \cdot x_t + b$
- 数据表示为:

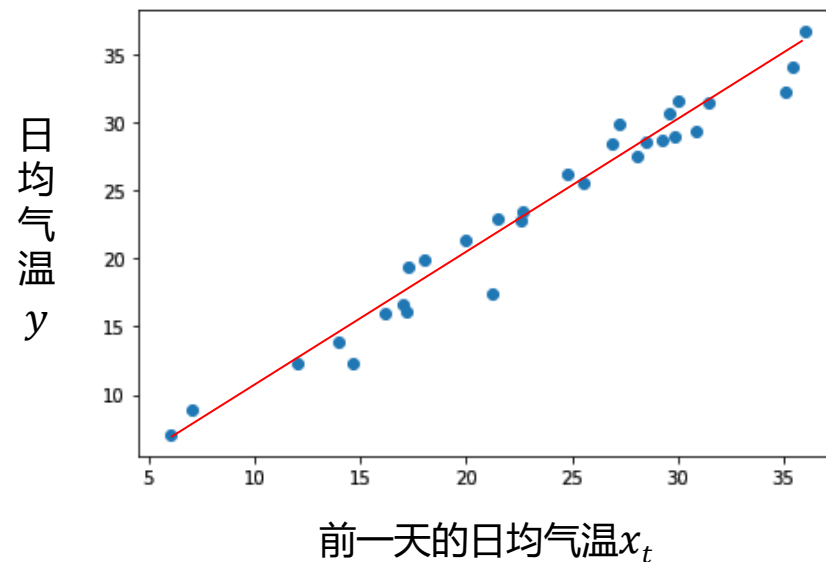
$$X = \begin{pmatrix} 1 & x_{1,t} \\ 1 & x_{2,t} \\ \dots & \dots \\ 1 & x_{30,t} \end{pmatrix}$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_{30} \end{pmatrix}$$

$$\mathbf{w} = \begin{pmatrix} b \\ w \end{pmatrix}$$

- 可得

$$\mathbf{w} = (X^T X)^{-1} X^T \mathbf{y} = \begin{pmatrix} 0.980 \\ 0.965 \end{pmatrix}$$



最小二乘法的问题

- 当这 n 个自变量不是互相独立，而是存在着一些线性关系时，此时 $X^T X$ 不可逆，此时得到的解为病态解，不能作为学习到的最优参数
- 当 $X^T X$ 不可逆时，可采用梯度下降（Gradient Descent）方法求解