



# 深度学习：从理论到实践

## 第一章：深度学习基础（下）



# 深度学习基础（二）

---

- ✓ 回归与分类
- ✓ 梯度下降
- ✓ 信息论

# 回归与分类

---

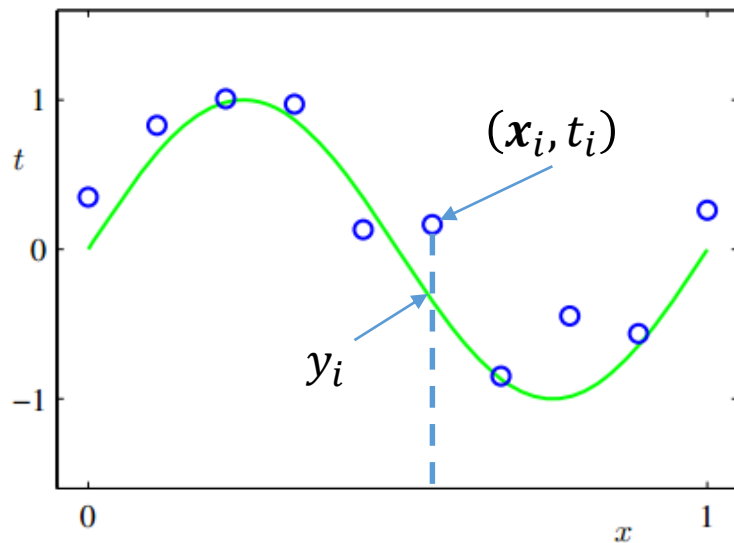
- ✓ 曲线拟合
- ✓ 线性回归
- ✓ logistic回归
- ✓ softmax回归

# 曲线拟合

✓ 回归基本问题：已知样本集 $D$ 的 $n$ 个样本  $(x_i, t_i)_{i=1}^n$ ，望获知自变量 $x$ 与变量 $y$ 的映射关系 $f$ 。

✓ 一个例子：

真实曲线 $\sin(2\pi x)$ ，  
 $n=10$ 个训练样本。  
根据样本点拟合曲线。



## ✓ 多项式函数拟合

$$g(x, \mathbf{w}) = \sum_{j=0}^M w_j x^j = w_0 + w_1 x_1 + \cdots + w_M x^M$$

其中多项式的阶、即最高次数 $M$ 待选择。

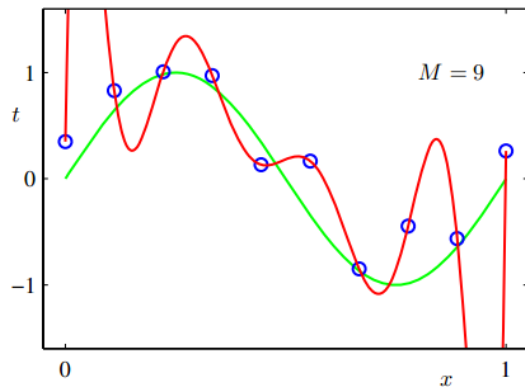
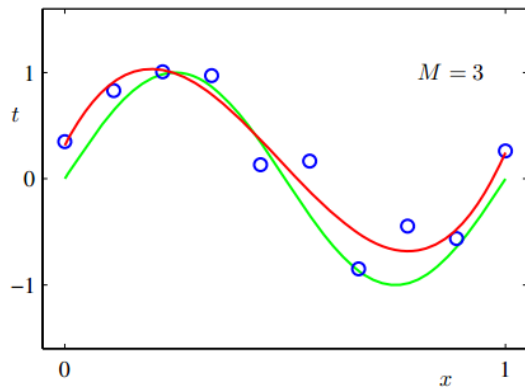
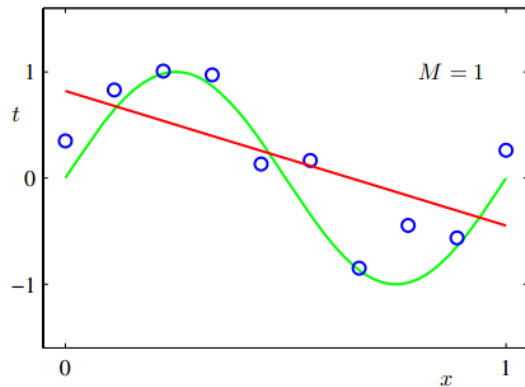
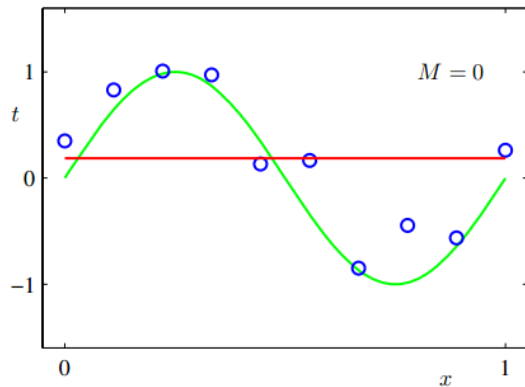
## ✓ 最小二乘拟合损失

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{ \underbrace{g(x_n, \mathbf{w})}_{\text{回归值}} - \underbrace{t_n}_{\text{样本值}} \}^2$$

# 曲线拟合

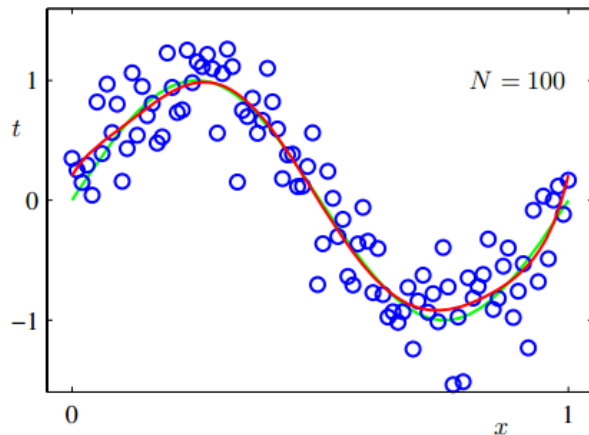
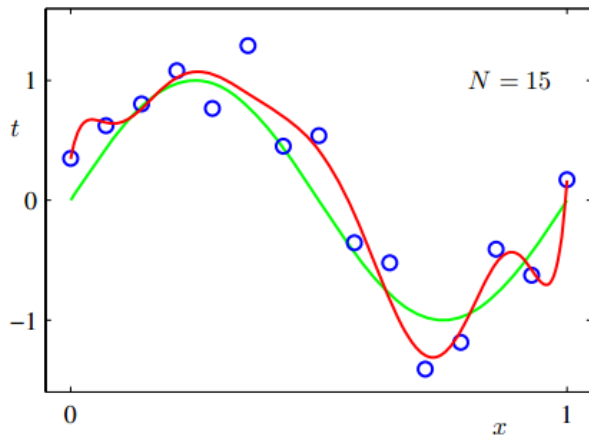
## ✓ 不同阶多项式 拟合的结果

高次多项式精确拟合，预测性能却极差，产生**过拟合**。



# 回避过拟合的两类办法

✓ 获取更多的训练样本



✓ 获取关于模式的先验知识，例如  $A \sin(\omega x + t)$

# 曲线拟合优化方法

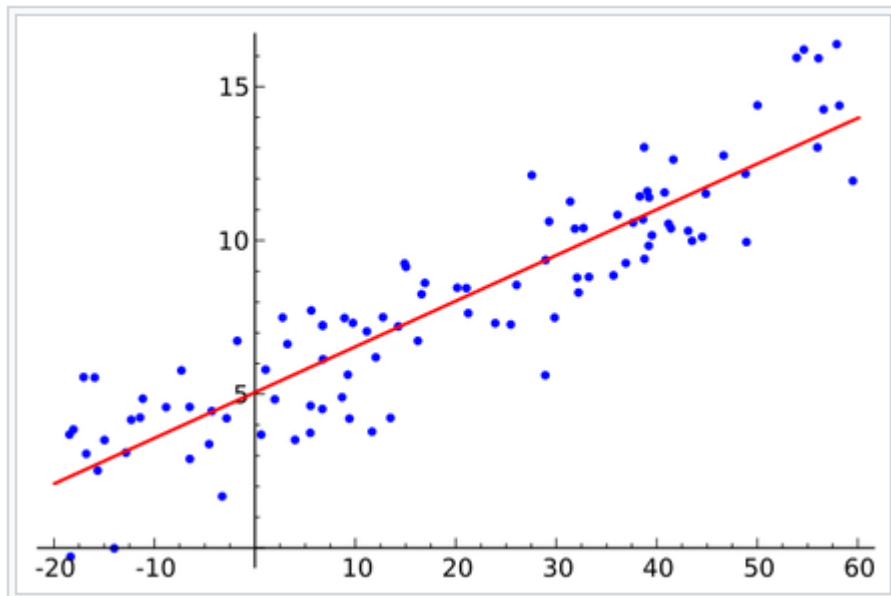
- ✓ 令  $\mathbf{z} = [1, x, x^2, \dots, x^M]^T$  , 拟合函数  $g = \mathbf{z}^T \mathbf{w}$  , 损失函数  $E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{\mathbf{z}_n^T \mathbf{w} - t_n\}^2$
- ✓ 可见, 拟合函数转化为线性函数, 即M维空间的线性回归。

特征变换



# 线性回归

- ✓ 线性回归：变量与自变量之间线性相关性，使用线性回归估计二者间的映射。



## ✓ 线性回归问题

$$E(\mathbf{w}) = \min_{\mathbf{w}} \frac{1}{2} \sum_{n=1}^N \{\mathbf{a}_n^T \mathbf{w} - b_n\}^2 = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{A}\mathbf{w} - \mathbf{b}\|^2 = \min_{\mathbf{w}} \left( \frac{\mathbf{w}^T \mathbf{A}^T \mathbf{A} \mathbf{w}}{2} - \mathbf{w}^T \mathbf{A}^T \mathbf{b} + \frac{\mathbf{b}^T \mathbf{b}}{2} \right)$$

$$, \text{ 其中 } \mathbf{A} = \begin{pmatrix} \mathbf{a}_1^T \\ \vdots \\ \mathbf{a}_N^T \end{pmatrix}, \mathbf{b} = \begin{pmatrix} b_1 \\ \vdots \\ b_N \end{pmatrix}.$$

$$✓ \text{ 最小二乘解: } \nabla_{\mathbf{w}} E(\mathbf{w}) = \mathbf{A}^T \mathbf{A} \mathbf{w} - \mathbf{A}^T \mathbf{b} = 0$$

$$\mathbf{w}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$$

# logistic回归

---

- ✓ 二分类问题，已知训练样本集 $D$ 的 $n$ 个样本  $(\mathbf{x}_i, t_i)_{i=1}^n$ ，其中 $t_i \in \{0, 1\}$ 为类别标签， $\mathbf{x}_i \in \mathcal{R}^d$ 为特征向量。
- ✓ 拟合特征向量到类别标签的回归，常用logistic回归。

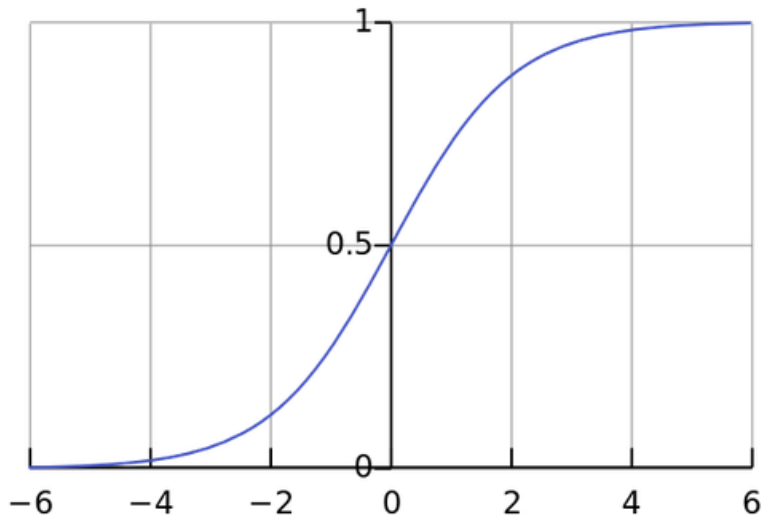
# logistic函数

✓ 标准logistic函数

$$f(x) = \frac{1}{1 + \exp(-x)}$$

✓  $f(-x) = 1 - f(x),$

✓  $f'(x) = f(x)(1 - f(x))$



# logistic回归

- ✓ logistic回归的回归函数为

$$g(\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})},$$

其中 $\mathbf{w}$  为回归参数。

- ✓ logistic回归，样本 $(\mathbf{x}_i, t_i)$ 的概率密度

$$P(\mathbf{x}_i, t_i; \mathbf{w}) = \begin{cases} g(\mathbf{x}_i), & t_i = 1 \\ 1 - g(\mathbf{x}_i), & t_i = 0 \end{cases} = g(\mathbf{x}_i)^{t_i} \cdot (1 - g(\mathbf{x}_i))^{1-t_i}$$

# logistic回归

- ✓ 极大似然估计 $\mathbf{w}$
- ✓ 似然函数 $p(D|\mathbf{w}) = \prod_{i=1}^n g(\mathbf{x}_i)^{t_i} (1 - g(\mathbf{x}_i))^{1-t_i}$ ,  
取负对数似然作为代价函数有

$$E(\mathbf{w}) = - \left[ \sum_{i=1}^n t_i \ln g(\mathbf{x}_i) + (1 - t_i) \ln(1 - g(\mathbf{x}_i)) \right]$$

# logistic回归: 优化

- ✓ 该优化问题采用Newton-Raphson迭代优化 ,

$$\mathbf{w}^{new} = \mathbf{w}^{old} - \mathbf{H}^{-1} \nabla E(\mathbf{w}),$$

其中 $\mathbf{H}$ 为 $E(\mathbf{w})$ 关于 $\mathbf{w}$ 的二阶导数矩阵。

- ✓  $\nabla E(\mathbf{w}) = \sum_{i=1}^n (y_i - t_i) \mathbf{x}_i = \mathbf{X}^T (\mathbf{y} - \mathbf{t})$
- ✓  $\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{i=1}^n y_i (1 - y_i) \mathbf{x}_i \mathbf{x}_i^T = \mathbf{X}^T \mathbf{R} \mathbf{X}$

# logistic回归：举例

- ✓ 一门考试，20位考生花费0~6小时备考。现在希望获悉备考时长与是否通过考试的关系。

0.50	0.75	1.00	1.25	1.50
0	0	0	0	0
1.75	1.75	2.00	2.25	2.50
0	1	0	1	0
2.75	3.00	3.25	3.50	4.00
1	0	1	0	1
4.25	4.50	4.75	5.00	5.50
1	1	1	1	1

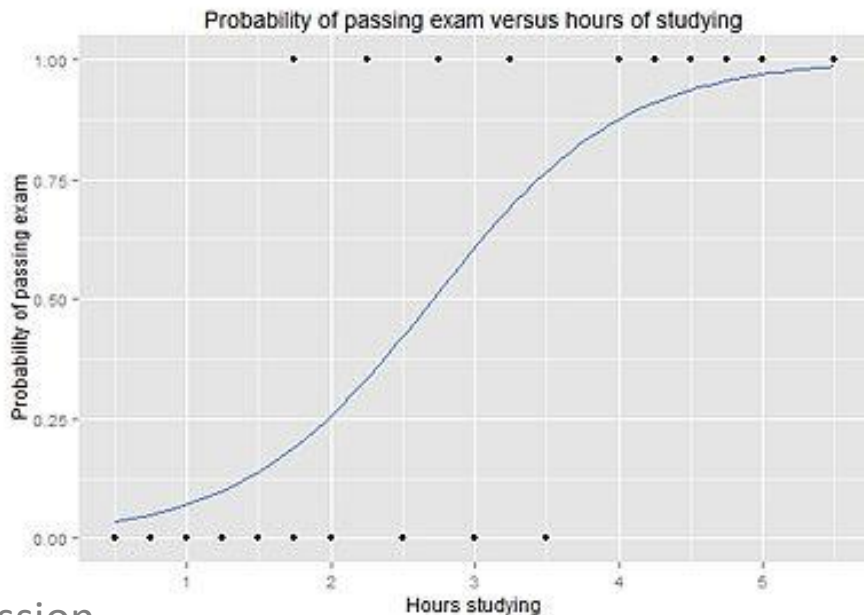


# logistic回归：举例

✓ 解释变量仅仅为1维的学习时间，回归参数为2维向量。

✓ 通过考试的概率为

✓ 
$$\frac{1}{1 + \exp\left(-\left(1.5046 \cdot \text{时长} - 4.0777\right)\right)}$$



- ✓ softmax回归用于多类分类问题。
- ✓  $K$ 类问题中，不采用标签定义 $1:K$ ，而使用标签向量：

$$\mathbf{t} = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} 1 \\ \\ k \\ \\ K \end{matrix}$$

0 – 1的 $K$ 维向量，若属于 $k$ 类，则向量的 $k$ 分量为1，其他分量均为0。

# softmax回归

✓ softmax回归的回归函数为

$$p(C_k|x) = \frac{\exp(\mathbf{w}_k^T \mathbf{x})}{\sum_j \exp(\mathbf{w}_j^T \mathbf{x})}$$

其中 $\mathbf{w}_k$  为第 $k$ 类的回归参数。

✓ softmax回归，某个 $\mathbf{x}_i$ 样本的概率：

$$\prod_{k=1}^K p(C_k|\mathbf{x}_i)^{t_{ik}}$$

其中， $\mathbf{t}_i = (t_{i1}, \dots, t_{ik}, \dots, t_{iK})^T$ 为 $\mathbf{x}$ 的标签向量

- ✓ 极大似然估计回归参数  $\mathbf{w}_1, \dots, \mathbf{w}_K$
- ✓ 似然函数  $p(D|\mathbf{w}_1, \dots, \mathbf{w}_K) = \prod_i \prod_{k=1}^K p(\mathcal{C}_k|\mathbf{x}_i)^{t_{ik}}$  , 取负对数似然

$$E(\mathbf{w}_1, \dots, \mathbf{w}_K) = - \sum_i \sum_{k=1}^K t_{ik} \ln p(\mathcal{C}_k|\mathbf{x}_i)$$

# 梯度下降(Gradient Descent)

---



- ✓ 基本理论
- ✓ 初始化与步长
- ✓ 深度学习应用

# 梯度下降基本理论

## ✓ 最小化问题

$$\min_x f(x)$$

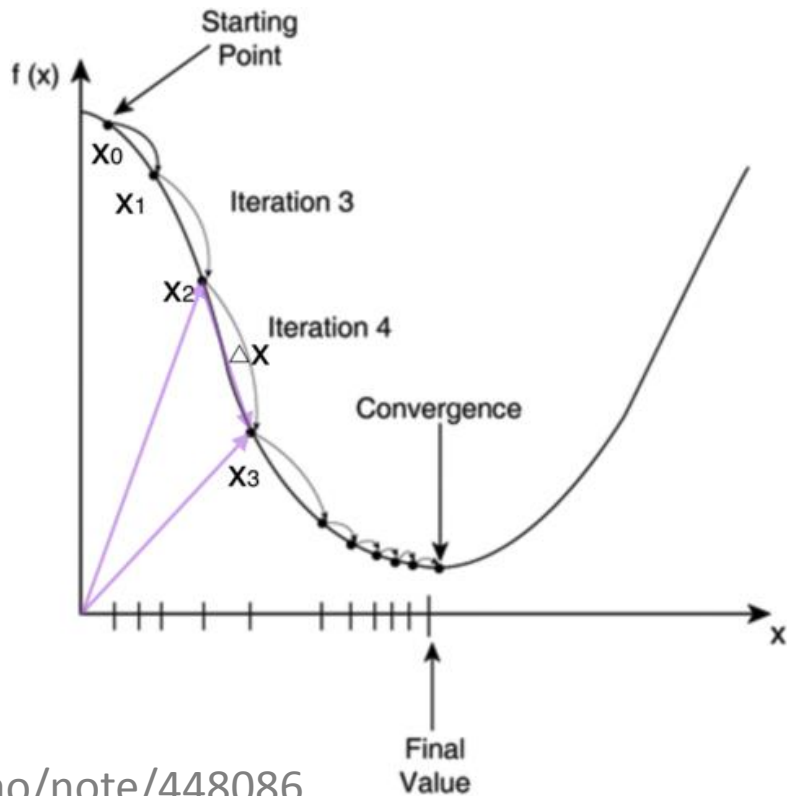
## ✓ 负梯度方向为函数值下降最快的方向

$$\boldsymbol{x}^{new} = \boldsymbol{x}^{old} - \lambda \nabla_x f(\boldsymbol{x})$$

其中， $\nabla_x f(\boldsymbol{x})$ 为函数在 $\boldsymbol{x}$ 处的梯度， $\lambda$ 为步长、亦称学习率。

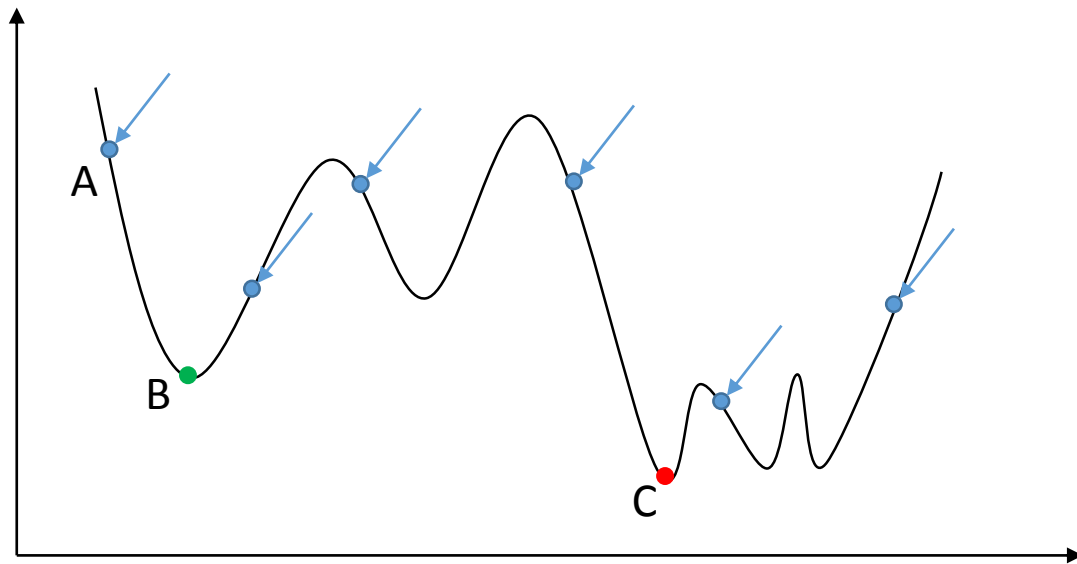
# 梯度下降基本理论

## ✓ 迭代更新



# 初始化

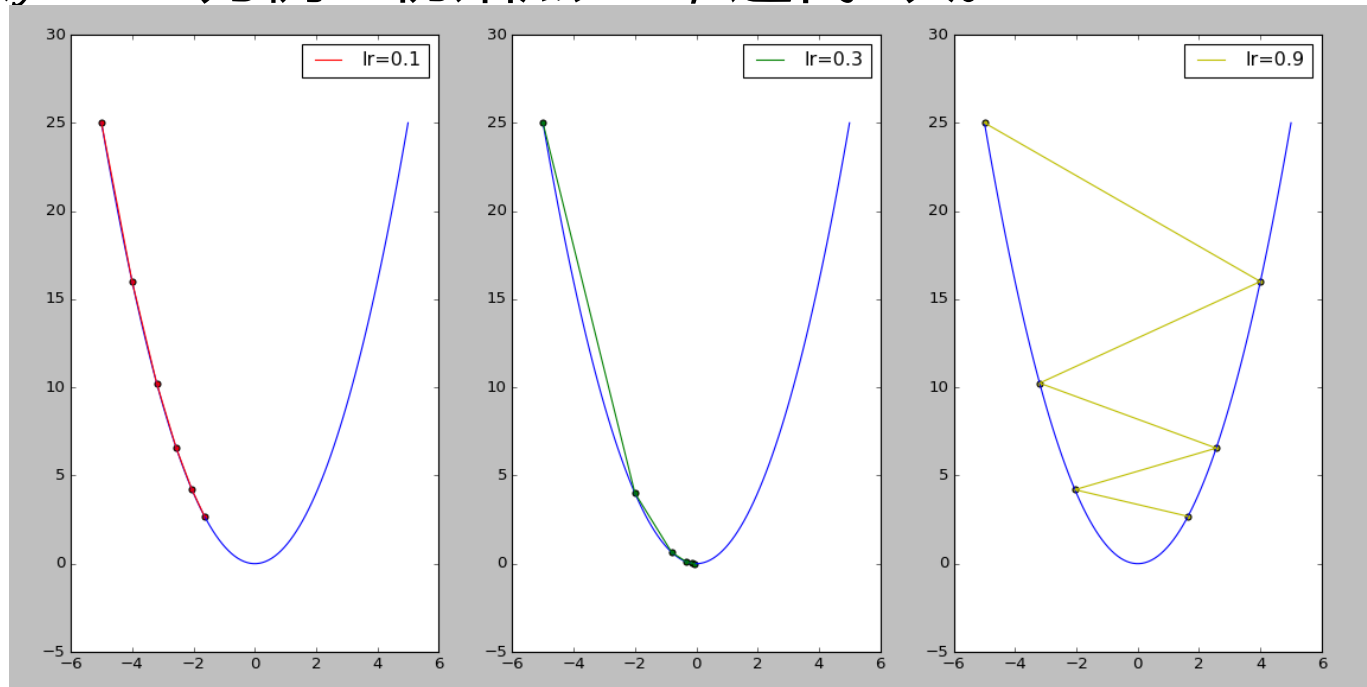
- ✓ 函数存在多个局部极小值点。迭代前需要选取恰当的初始点。





# 步长选取

✓ 以 $y = x^2$ 为例。初始点-5，迭代5次。



✓  $\theta$  参数估计问题：假定样本集  $D$  含有  $n$  个样本  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ 。假定损失函数为  $J(\theta; D, f)$ ，其中  $f(\mathbf{x}; \theta)$  为待估计参数的映射。

✓ 线性回归  $J(\theta; D, f) = \frac{1}{2n} \sum_{i=1}^n (f(\mathbf{x}_i; \theta) - y_i)^2$

✓ logistic 回归  $J(\theta; D, f) = -\frac{1}{n} [\sum_{i=1}^n y_i \ln f(\mathbf{x}_i; \theta) + (1 - y_i) \ln(1 - f(\mathbf{x}_i; \theta))]$

# 批量梯度下降(Batch GD)

---

- ✓ 梯度下降： $\theta \leftarrow \theta - \lambda \nabla_{\theta} J(\theta; D, f)$
- ✓ 线性回归： $\theta \leftarrow \theta - \frac{\lambda}{n} \sum_{i=1}^n (f'_{\theta}(x_i; \theta) - y_i)$
- ✓ 每次迭代，使用所有的训练样本，迭代速度受到样本量的影响。尤其样本特别多时，迭代速度较慢

# 随机梯度下降(Stochastic GD)



- ✓ 为回避一次使用所有样本，随机梯度下降每次迭代随机选取单个样本而决定迭代方向。
- ✓ 线性回归： $\theta \leftarrow \theta - \lambda(f'_{\theta}(x_i; \theta) - y_i)$
- ✓ 单样本用于单次迭代，计算代价相对BGD较小。同时单个样本的迭代方向，也许能应对多极值点，也可能导致迭代路径的曲折

# 最小批量梯度下降(MiniBatch GD)



- ✓ 折中批量和随机梯度下降，每次迭代随机选取m个样本
- ✓ 线性回归： $\theta \leftarrow \theta - \frac{\lambda}{m} \sum_{i \in \mathcal{I}_m} (f'_{\theta}(x_i; \theta) - y_i)$
- ✓ 常用于深度学习的参数估计

- ✓ 信息熵
- ✓ 相对熵
- ✓ 互信息
- ✓ 交叉熵及深度学习的应用

- ✓ 给定概率密度函数 $p(x)$ ，定义该函数的信息熵

$$H(p) = H[x] = - \int p(x) \ln p(x) dx$$

- ✓ 信息熵描述了分布的混乱程度。均匀分布是使得信息熵最大的概率分布。单点的冲击响应函数对应的信息熵最小。

# 相对熵

- ✓ 给定两个概率密度函数 $p(x)$ 和 $q(x)$ ，描述二者之间的差异（距离），定义相对熵

$$\begin{aligned} KL(p||q) &= - \int p(x) \ln q(x) dx - \left( - \int p(x) \ln p(x) dx \right) \\ &= - \int p(x) \ln \left\{ \frac{q(x)}{p(x)} \right\} dx \end{aligned}$$

- ✓ 对任意概率分布， $KL(p||q) \geq 0$ ，等号当且仅当 $p = q$ 。



- ✓ 对于两个随机变量 $x, y$ ，定义二者之间的互信息

$$\begin{aligned} I[x, y] &= KL(p(x, y) || p(x)p(y)) \\ &= - \iint p(x, y) \ln \left( \frac{p(x)p(y)}{p(x, y)} \right) dx dy \end{aligned}$$

- ✓ 若 $x, y$ 相互独立，则互信息为0，二者相互无关
- ✓  $I[x, y] = H[x] - H[x|y] = H[y] - H[y|x]$ 。

- ✓ 给定两个概率密度函数 $p(x)$ 和 $q(x)$ ，定义 $p(x)$ 关于 $q(x)$ 的交叉熵

$$\begin{aligned} H(p, q) &= E_p(-\ln q) = - \int p(x) \ln q(x) dx \\ &= H(p) + KL(p||q) \end{aligned}$$

- ✓ 交叉熵作为logistic、softmax回归的代价函数，常应用神经网络的输出层。

✓ 交叉熵  $H(p, q) = H(p) + KL(p||q)$

✓ logistic回归的对数似然函数

$$E(\mathbf{w}) = - \left[ \sum_{i=1}^n t_i \ln g(\mathbf{x}_i) + (1 - t_i) \ln(1 - g(\mathbf{x}_i)) \right]$$

✓ logistic、softmax回归，常应用神经网络的分类层。

- ✓ Pattern Recognition and Machine Learning , Christopher M. Bishop , Springer , 2007
- ✓ Pattern Classification , Richard O. Duda , Peter E. Hart and David G. Stork , Wiley-Interscience , 2000

# 在线问答

---

Q&A



课程地址

《点击此处添加演讲主题演讲主题演讲主题演讲主题》



**感谢各位聆听 !**  
Thanks for Listening

