

# 线性回归 Linear Regression



# 目录(CONTENT)



- 01 线性回归 (Linear Regression)
- 02 正则化 (Regularization)
- 03 逻辑斯蒂回归 (Logistic Regression)
- 04 多分类学习 (Multiclass Classification)
- 05 类别不平衡问题 (Class-Imbalance)
- 06 小结(Summary)

# 线性回归





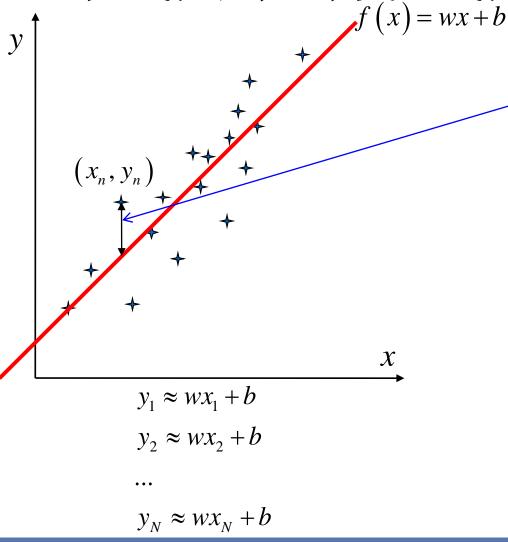
## 线性回归

**Linear Regression** 

# 线性回归-->一元线性回归



## ■ 一元线性回归/单变量线性回归



▶ 回归误差:

$$e_n = y_n - (wx_n + b)$$

> 均方和误差:

$$\sum_{n=1}^{N} e_n^2 = \sum_{n=1}^{N} (y_n - wx_n - b)^2$$
✓ 欧式距离

✓ 能量...

> 最小二乘法

(Least Square Method)

$$(w^*, b^*) = \arg\min_{w, b} \sum_{n=1}^{N} (y_n - wx_n - b)^2$$

# 线性回归-->多元线性回归



W

## ■ 多元线性回归 (Multivariate Linear Regression)

$$y_{1} \approx x_{1}w + b$$
  $y_{1} \approx x_{1,1}w_{1} + x_{1,2}w_{2} + ... + x_{1,D}w_{D} + b$   $y_{2} \approx x_{2}w + b$  ...  $y_{N} \approx x_{N}w + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$  ...  $y_{N} \approx x_{N}w + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$   $y_{N} \approx x_{N,1}w_{1} + x_{N,2}w_{2} + ... + x_{N,D}w_{D} + b$ 

# 线性回归-->最小二乘法故事



1801年,意大利天文学家朱赛普·皮亚齐发现了第一颗小行星谷神星。经过40天的跟踪观测后,由于谷神星运行至太阳背后,使得皮亚齐失去了谷神星的位置。随后全世界的科学家利用皮亚齐的观测数据开始寻找谷神星,但是根据大多数人计算的结果来寻找谷神星都没有结果。时年24岁的高斯也计算了谷神星的轨道。奥地利天文学家海因里希·奥尔伯斯根据高斯计算出来的轨道重新发现了谷神星。

高斯使用的最小二乘法的方法发表于1809年他的著作《天体运动论》中。

法国科学家勒让德于1806年独立发明"最小二乘法",但因不为世人所知而默默无闻。

勒让德曾与高斯为谁最早创立最小二乘法原理发生争执。

# 线性回归-->问题求解/优化



### ■ 问题求解/优化

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

➤ 闭合解/解析解 (Closed-form Solution/ Analytic Solution):

$$J(\mathbf{w}) = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_{2}^{2} = (\mathbf{y} - \mathbf{X}\mathbf{w})^{T} (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 2\mathbf{X}^{T} (\mathbf{y} - \mathbf{X}\mathbf{w}) = 2\mathbf{X}^{T} \mathbf{X}\mathbf{w} - 2\mathbf{X}^{T} \mathbf{y} = \mathbf{0}$$

$$\mathbf{w}^{*} = (\mathbf{X}^{T} \mathbf{X})^{-1} \mathbf{X}^{T} \mathbf{y} \qquad \mathbf{w}^{*} = (\mathbf{X}^{T} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^{T} \mathbf{y}$$

$$\checkmark \mathbb{E} \mathbb{Q} \mathcal{U}$$

$$\psi \mathbb{E} \mathbb{Q} \mathcal{U}$$

➤ 梯度下降法 (Gradient Descent):

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 2\mathbf{X}^{T}(\mathbf{y} - \mathbf{X}\mathbf{w}) \quad \Rightarrow \quad \mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$$

https://www.bilibili.com/video/av41473299?from=search&seid=8346799062692305739



# 线性回归-->最小二乘几何解释

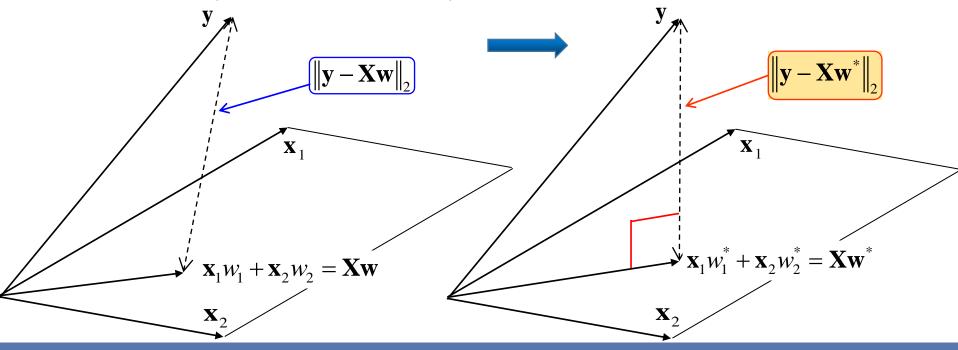


## ■ 子空间X上的<u>正交投影</u>

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{w} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \left( \begin{bmatrix} x_{1,1} \\ \dots \\ x_{N,1} \end{bmatrix} w_1 + \begin{bmatrix} x_{1,2} \\ \dots \\ x_{N,2} \end{bmatrix} w_2 + \dots + \begin{bmatrix} x_{1,D} \\ \dots \\ x_{N,D} \end{bmatrix} w_D \right)$$

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

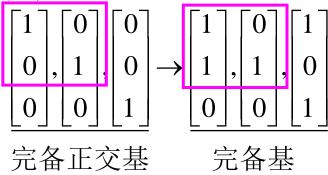
$$= \mathbf{y} - (\mathbf{x}_1 w_1 + \mathbf{x}_2 w_2 + \dots + \mathbf{x}_d w_D)$$

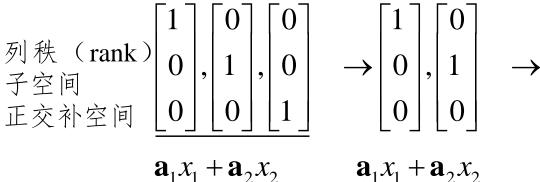


# 线性回归-->最小二乘几何解释

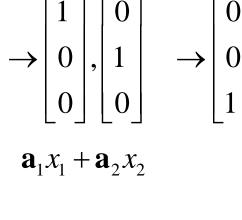


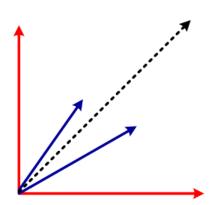
## ■子空间

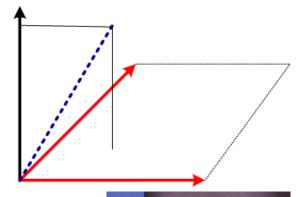




 $+\mathbf{a}_3 x_3$ 







Linear Algebra Lecture #1 W. Gilbert Strang

MIT公开课《线性代数》:

https://www.bilibili.com/video/av24368594

# 线性回归-->最小二乘概率解释



## ■ 高斯噪声假设

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{w} \qquad \Rightarrow \qquad \begin{bmatrix} e_1 \\ \vdots \\ e_N \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} x_{1,1}w_1 + x_{1,2}w_2 + \dots + x_{1,D}w_D \\ \vdots \\ x_{N,1}w_1 + x_{N,2}w_2 + \dots + x_{N,D}w_D \end{bmatrix}$$

> 零均值等方差高斯噪声假设

$$e_n \sim \mathcal{N}(0, \sigma^2) \Rightarrow p(e_n) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{e_n^2}{2\sigma^2}\right)$$

> 最大似然估计

$$p(e_1, e_2, ..., e_N) = \prod_{n=1}^{N} p(e_n) = \left(\frac{1}{\sqrt{2\pi\sigma}}\right)^N \exp\left(-\sum_{n=1}^{N} e_n^2 / 2\sigma^2\right)$$

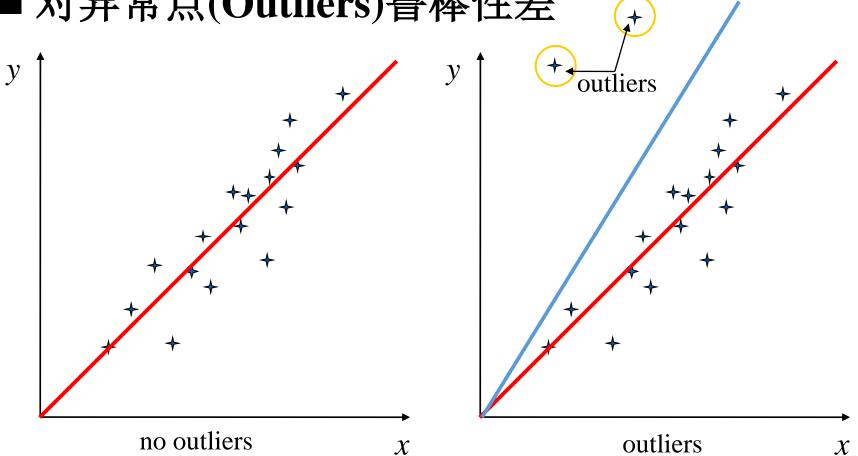
> 概率解释

$$\max_{\mathbf{w}} p(e_1, e_2, ..., e_N; \mathbf{w}) \Leftrightarrow \min_{\mathbf{w}} \sum_{n=1}^{N} e_n^2 \Leftrightarrow \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

# 线性回归-->最小二乘鲁棒性



■ 对异常点(Outliers)鲁棒性差



- ➤ 随机取样一致 (Random Sample Consensus, RANSAC)
- ➤ 鲁棒回归 (Robust Regression)

# 线性回归-->建模非线性



## ■线性的含义

$$y = wx + b$$

线性?



$$y = w_1 x + w_2 x^2 + w_3 x^3 + b$$

线性?



$$x_1 = x$$
,  $x_2 = x^2$ ,  $x_3 = x^3$   
 $y = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$ 

线性?



- 线性并不指对输入变量的线性,而是指对参数空间的线性。也就说对于输入来说,完全可以先对其进行非线性变换,再进行线性组合。从这个角度来说,线性模型完全具有描述非线性的能力。
- ▶ 通用非线性化方法: 核学习方法 (Kernel-based Learning Algorithms)

# 线性回归-->广义线性模型

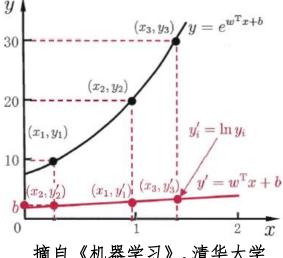


## ■广义线性模型 (Generalized Linear Model)

➤ 对数线性回归 (Log-Linear Regression): <sup>y</sup>

$$y = \exp(wx + b)$$
  $\log y = wx + b$ 

形式上仍是线性回归,实质上已经是求取输入空间到输出空间的非线性函数映射



摘自《机器学习》,清华大学 出版社,周志华著,P56图3.1

➤ 广义线性模型 (Generalized Linear Model):

$$y = g^{-1}(\mathbf{w}^T \mathbf{x}) \qquad y' = g(y) = \mathbf{w}^T \mathbf{x}$$
  
单调可微函数





# 正则化 Regularization



- 一般的目标函数包含两部分:
  - 1. 数据项 (Data Term):

回归/分类的目标,比如:误差尽可能小;

分类尽可能准确等。

2. 正则化项 (Regularization Term):

对参数空间的限制/对解额外属性的追求(比如稀疏)

数据项 
$$O(\mathbf{x}) = D(\mathbf{x}) + \lambda R(\mathbf{x})$$
 正则化项

$$D(\mathbf{x}) \ge 0, R(\mathbf{x}) \ge 0$$



### ■ 岭回归 (Ridge Regression):

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} + \lambda \|\mathbf{x}\|_{2}^{2}$$

Data Term: 
$$\frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2}$$

*Regularization Term* :  $\|\mathbf{x}\|_{2}^{2}$ 

 $trade-off:\lambda$ 

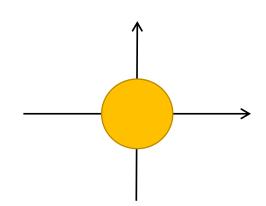
▶ 正则化是对解空间的一种限制

$$\min_{\mathbf{x}} D(\mathbf{x}) + \lambda R(\mathbf{x})$$

$$\Leftrightarrow$$

$$\min_{\mathbf{x}} D(\mathbf{x})$$

$$s.t. R(\mathbf{x}) \leq M$$



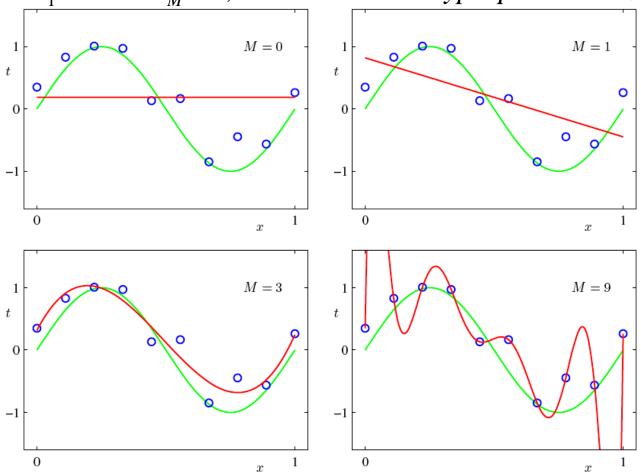


### ■ 举例

最小二乘模型: 
$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2}$$

$$y = w_0 + w_1 x + \dots + w_M x^M;$$







### ■ 举例

Table 1.1 Table of the coefficients w\* for polynomials of various order.

Observe how the typical magnitude of the coefficients increases dramatically as the order of the polynomial increases.

	M = 0	M = 1	M = 6	M = 9
$w_0^{\star}$	0.19	0.82	0.31	0.35
$w_1^{\star}$		-1.27	7.99	232.37
$w_2^{\star}$			-25.43	-5321.83
$w_3^{\star}$			17.37	48568.31
$w_4^{\star}$				-231639.30
$w_5^{\star}$				640042.26
$w_6^{\star}$				-1061800.52
$w_7^{\star}$				1042400.18
$w_8^{\star}$				-557682.99
$w_9^{\star}$				125201.43



#### ■ 举例

>增大样本可以避免过拟合,但是。。。

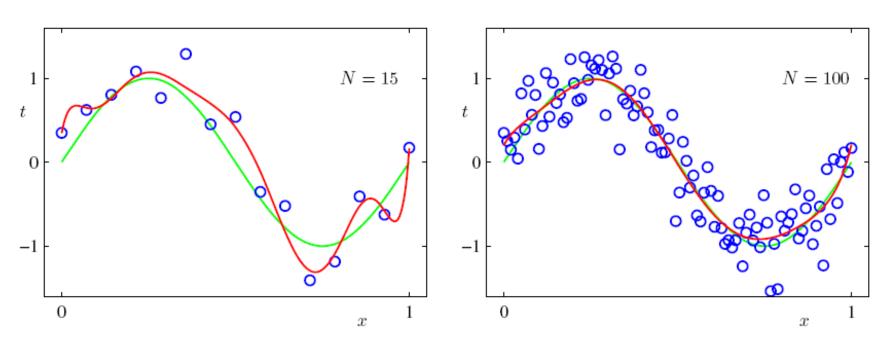


Figure 1.6 Plots of the solutions obtained by minimizing the sum-of-squares error function using the M=9 polynomial for N=15 data points (left plot) and N=100 data points (right plot). We see that increasing the size of the data set reduces the over-fitting problem.



#### ■ 举例

岭回归: 
$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} + \frac{1}{2} \lambda \|\mathbf{x}\|_{2}^{2}$$

#### >正则化,控制模型的复杂度

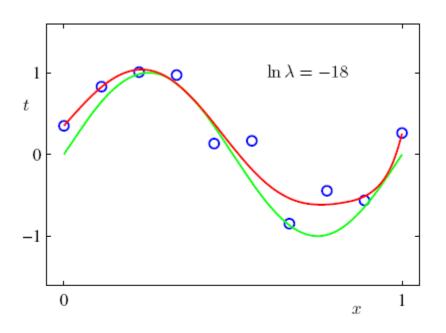


Figure 1.7 Plots of M=9 polynomials fitted to the data function (1.4) for two values of the regularization parameters of no regularizer, i.e.,  $\lambda=0$ , corresponding to  $\ln\lambda=0$ 

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^{\star}$	0.35	0.35	0.13
$w_1^{\star}$	232.37	4.74	-0.05
$w_2^{\star}$	-5321.83	-0.77	-0.06
$w_3^{\star}$	48568.31	-31.97	-0.05
$w_4^{\star}$	-231639.30	-3.89	-0.03
$w_5^{\star}$	640042.26	55.28	-0.02
$w_6^{\star}$	-1061800.52	41.32	-0.01
$w_7^{\star}$	1042400.18	-45.95	-0.00
$w_8^{\star}$	-557682.99	-91.53	0.00
$w_9^{\star}$	125201.43	72.68	0.01



### ■ 不同正则化追求解的不同属性

addition, the second term is a regularization term. Different kinds of regularization lead to different properties of the coding coefficients ( $\ell_1$ -regularization (2(a)) encourages sparsity [29],  $\ell_{2,1}$ -regularization (2(b)) promotes group sparsity, and  $\ell_2$ -regularization (2(c)) does not lead to sparsity but makes the solution simple and stable [35]). After obtaining the coding

 Code a test sample y by a linear representation (2(a), 2(b) or 2(c)),

2(a) 
$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} + \lambda \|\mathbf{x}\|_{1} \text{ (SR [29])}$$
  
2(b)  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} + \lambda \|\mathbf{x}\|_{2,1} \text{ (GSR [17])}$   
2(c)  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} + \lambda \|\mathbf{x}\|_{2}^{2} \text{ (CR [35])}$ 

to obtain the corresponding coding coefficients.

➤ Dong Wang, Huchuan Lu, Minghsuan Yang, Kernel Collaborative Face Recognition, Pattern Recognition, 48(10):3025-3037, 2015.





■ 正则化项可以看作参数的先验(从贝叶斯估计角度)

观测数据: y

参数/系数: x

$$p(\mathbf{x} | \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{y})} = \frac{p(\mathbf{y} | \mathbf{x}) p(\mathbf{x})}{p(\mathbf{y})} \propto p(\mathbf{y} | \mathbf{x}) p(\mathbf{x})$$

Posteriori

Likelihood prior

(1)MAP (Maximum a Posteriori):

$$\mathbf{x}^* = \arg\max_{\mathbf{x}} p(\mathbf{x} | \mathbf{y}) = \arg\max_{\mathbf{x}} p(\mathbf{x}, \mathbf{y}) = \arg\max_{\mathbf{x}} p(\mathbf{y} | \mathbf{x}) p(\mathbf{x})$$

(2) MLE (Maximum Likelihood Estimation):

$$\mathbf{x}^* = \arg\max_{\mathbf{x}} \, p\left(\mathbf{y} \mid \mathbf{x}\right)$$

 $ps: p(\mathbf{x}) = const, uniform prior$ 

The coefficient x can be obtained by maximizing the posteriori probability  $p(\mathbf{x}|\mathbf{y})$ , which is also equivalent to maximizing the joint likelihood probability  $p(\mathbf{x},\mathbf{y})$ . Assume there is a uniform prior, the coefficient x is estimated by  $\hat{\mathbf{x}} = \arg\max_{\mathbf{x}} p(\mathbf{y}|\mathbf{x}) = \arg\max_{\mathbf{x}} p(\mathbf{e})$ , which is the maximum likelihood estimation (MLE).



■ MLE(Maximum Likelihood Estimation)与OLS(Ordinary Least-Squares):

$$y = Ax$$

 $error / residual : \mathbf{e} = \mathbf{y} - \mathbf{A}\mathbf{x}$ 

Assume that  $\mathbf{e} = [e_1, e_2, ..., e_N]$  follow i.i.d Gaussian distribution  $N(0, \sigma^2)$ 

 $\mathbf{x} = [x_1, x_2, ..., x_d]$  follow <u>some</u> uniform prior.

$$p(\mathbf{x}|\mathbf{y}) \propto p(\mathbf{y}|\mathbf{x}) p(\mathbf{x})$$

$$= p(\mathbf{e}) = \prod_{n=1}^{N} p(e_n)$$

$$= \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma}} \exp \left[ -\frac{e_n^2}{2\sigma^2} \right]$$

$$= \left\{ \prod_{n=1}^{N} \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{\sigma^{2}}\right) \right\} \exp\left(-\frac{1}{2} \left\|\mathbf{e}\right\|_{2}^{2}\right)$$

$$\max_{\mathbf{x}} p(\mathbf{x} | \mathbf{y}) \Leftrightarrow \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2}$$



#### ■ MAP(Maximum A Posterior)与RR(Ridge Regression):

 $\propto \exp \left[ -\frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} - \frac{\lambda}{2} \|\mathbf{x}\|_{2}^{2} \right] \qquad \left( ps : \lambda = \frac{\sigma^{2}}{\sigma^{2}} \right)$ 

$$\begin{aligned} \mathbf{y} &= \mathbf{A}\mathbf{x} & error / residual : \mathbf{e} &= \mathbf{y} - \mathbf{A}\mathbf{x} \\ Assume that \ \mathbf{e} &= \left[e_1, e_2, ..., e_N\right] follow i.i.d \ Gaussian \ distribution \ N\left(0, \sigma^2\right) \\ \mathbf{x} &= \left[x_1, x_2, ..., x_D\right] also \ follow i.i.d \ Gaussian \ distribution \ N\left(0, \sigma_1^2\right) \\ p\left(\mathbf{x} \mid \mathbf{y}\right) &\propto p\left(\mathbf{y} \mid \mathbf{x}\right) p\left(\mathbf{x}\right) = p\left(\mathbf{e}\right) p\left(\mathbf{x}\right) = \left[\prod_{n=1}^N p\left(e_n\right)\right] \left[\prod_{d=1}^D p\left(x_d\right)\right] \\ &= \left\{\prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{e_n^2}{2\sigma^2}\right]\right\} \left\{\prod_{d=1}^D \frac{1}{\sqrt{2\pi}\sigma_1} \exp\left[-\frac{x_d^2}{2\sigma_1^2}\right]\right\} \\ &\propto \left\{\prod_{n=1}^N \exp\left[-\frac{e_n^2}{2\sigma^2}\right]\right\} \left\{\prod_{d=1}^D \exp\left[-\frac{x_d^2}{2\sigma_1^2}\right]\right\} \\ &= \exp\left(-\frac{\left\|\mathbf{e}\right\|_2^2}{2\sigma^2}\right) \exp\left(-\frac{\left\|\mathbf{x}\right\|_2^2}{2\sigma_1^2}\right) = \exp\left(\frac{1}{\sigma^2}\right) \exp\left[-\frac{\left\|\mathbf{e}\right\|_2^2}{2\sigma_1^2}\right] \left\|\mathbf{x}\right\|_2^2 \end{aligned}$$



### ■ Ordinary Least-Squares v.s. Ridge Regression:

#### 1. Objective Function:

$$OLS: \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2}$$

$$RR: \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2} + \frac{\lambda}{2} \|\mathbf{x}\|_{2}^{2}$$

#### 2.Solution:

$$OLS: \mathbf{x} = \left(\mathbf{A}^T \mathbf{A}\right)^{-1} \mathbf{A}^T \mathbf{y}$$

$$RR: \mathbf{x} = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I})^{-1} \mathbf{A}^T \mathbf{y}$$

#### 3. Probability:

OLS: MLE

RR: MAP

- ➤ 吴恩达《机器学习》-正则化 https://www.bilibili.com/video/ av55276229
- ► L1&L2正则化详解

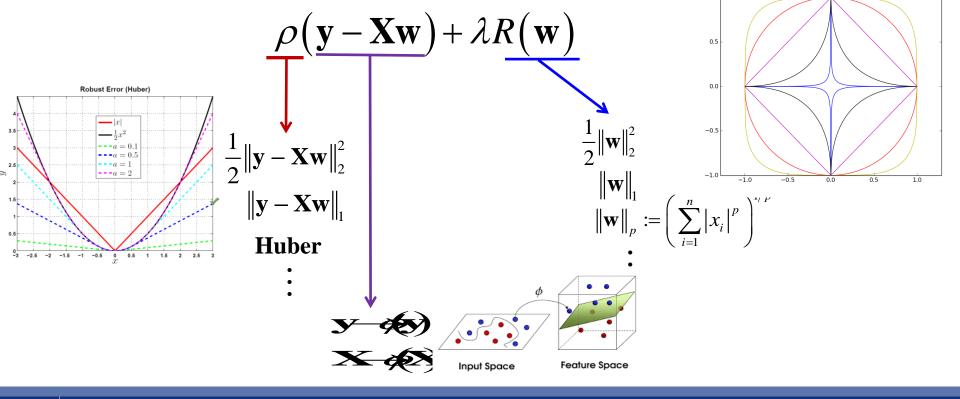
  https://www.bilibili.com/video/
  av77106463?from=search&sei
  d=4369320229005019988
- ► 什么是L1 L2正则化?

  https://www.bilibili.com/video/
  av16009446?from=search&sei
  d=4369320229005019988

# 线性回归归纳



$$\mathbf{e} = \mathbf{y} - \mathbf{X}\mathbf{w} \qquad \Rightarrow \qquad \begin{bmatrix} e_1 \\ \vdots \\ e_N \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} - \begin{bmatrix} x_{1,1}w_1 + x_{1,2}w_2 + \dots + x_{1,D}w_D \\ \vdots \\ x_{N,1}w_1 + x_{N,2}w_2 + \dots + x_{N,D}w_D \end{bmatrix}$$



# 逻辑斯蒂回归





## 逻辑斯蒂回归

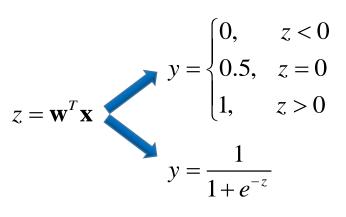
**Logistic Regression** 

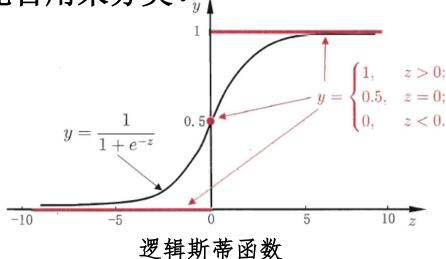
# 逻辑斯蒂回归-->引言



# ■逻辑斯蒂回归 (Logistic Regression)

> 线性回归的任务是预测,能否用来分类?"





> 广义线性回归特列

《机器学习》,清华大学出版社,周志华著,P58

# 逻辑斯蒂回归-->概率解释

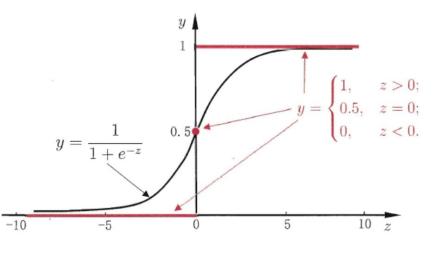


## ■ 概率解释

$$>$$
 分类概率  $y = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})}$ 

$$p_{1}(\mathbf{x}; \mathbf{w}) = p(y = 1 | \mathbf{x}; \mathbf{w}) = \frac{\exp(\mathbf{w}^{T} \mathbf{x})}{1 + \exp(\mathbf{w}^{T} \mathbf{x})}$$

$$p_{0}(\mathbf{x}; \mathbf{w}) = p(y = 0 | \mathbf{x}; \mathbf{w}) = \frac{1}{1 + \exp(\mathbf{w}^{T} \mathbf{x})}$$



### > 贝叶斯决策

$$\log \frac{p(y=1|\mathbf{x};\mathbf{w})}{p(y=0|\mathbf{x};\mathbf{w})} = \mathbf{w}^T \mathbf{x} \qquad \Rightarrow \begin{cases} \mathbf{w}^T \mathbf{x} > 0 & \to y=1 \\ \mathbf{w}^T \mathbf{x} < 0 & \to y=0 \end{cases}$$

### > 判决模型

- ✔ 从生成模型和判决模型角度,逻辑斯蒂回归算法属于判决模型;
- ✓ 直接对分类的可能性进行建模,无需事先假设数据分布,避免数据分布假设不准确所带来的问题;
- ✓ 并且逻辑斯蒂回归不仅输出分类类别,而且能得到概率近似。

# 逻辑斯蒂回归-->极大似然推导



### ■ 极大似然推导

▶ 伯努利分布

$y_n$	1	0
$p_n$	$p^1(\mathbf{x}_n;\mathbf{w})$	$p^0(\mathbf{x}_n;\mathbf{w})$

$$p^{1}(\mathbf{x}_{n};\mathbf{w})+p^{0}(\mathbf{x}_{n};\mathbf{w})=1$$

➤ 最大似然估计 (i.i.d, independent and identically distributed)

$$p(\mathbf{x}_{1}, \mathbf{x}_{2}, ..., \mathbf{x}_{N}; \mathbf{w}) = \prod_{n=1}^{N} \left[ p^{1}(\mathbf{x}_{n}; \mathbf{w}) \right]^{y_{n}} \left[ p^{0}(\mathbf{x}_{n}; \mathbf{w}) \right]^{1-y_{n}}$$

$$\max_{\mathbf{w}} p(\mathbf{x}_{1}, \mathbf{x}_{2}, ..., \mathbf{x}_{N}; \mathbf{w}) \Leftrightarrow \min_{\mathbf{w}} \left[ -\log p(\mathbf{x}_{1}, \mathbf{x}_{2}, ..., \mathbf{x}_{N}; \mathbf{w}) \right]$$

> 目标函数

$$L(\mathbf{w}) = -\log p(\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_N; \mathbf{w})$$

$$= -\sum_{n=1}^{N} \left[ y_n \log p^1(\mathbf{x}_n; \mathbf{w}) + (1 - y_n) \log p^0(\mathbf{x}_n; \mathbf{w}) \right]$$

$$= \sum_{n=1}^{N} \left[ -y_n \mathbf{w}^T \mathbf{x}_n + \log \left( 1 + \exp \left( \mathbf{w}^T \mathbf{x}_n \right) \right) \right]$$

# 逻辑斯蒂回归-->问题求解/优化 @ 大连强二大学



### 问题求解/优化

$$\mathbf{w}^* = \arg\min_{\mathbf{w}} L(\mathbf{w})$$

$$L(\mathbf{w}) = \sum_{n=1}^{N} \left[ -y_n \mathbf{w}^T \mathbf{x}_n + \log(1 + \exp(\mathbf{w}^T \mathbf{x}_n)) \right]$$

▶ 一阶和二阶导数

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -\sum_{n=1}^{N} \left( y_n - \frac{\exp(\mathbf{w}^T \mathbf{x}_n)}{1 + \exp(\mathbf{w}^T \mathbf{x}_n)} \right) \mathbf{x}_n = -\sum_{n=1}^{N} \left[ y_n - p^1(\mathbf{x}_n; \mathbf{w}) \right] \mathbf{x}_n$$

$$\frac{\partial^2 L(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = \sum_{n=1}^{N} \frac{\exp(\mathbf{w}^T \mathbf{x}_n)}{\left(1 + \exp(\mathbf{w}^T \mathbf{x}_n)\right)^2} \mathbf{x}_n \mathbf{x}_n^T = \sum_{n=1}^{N} p^1(\mathbf{x}_n; \mathbf{w}) p^0(\mathbf{x}_n; \mathbf{w}) \mathbf{x}_n \mathbf{x}_n^T$$

> 迭代求解

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \, \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$$

梯度下降法

**Gradient Descent Method** 

$$\mathbf{w} \leftarrow \mathbf{w} - \left(\frac{\partial^2 L(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T}\right)^{-1} \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}}$$
  
牛顿法

**Newton Method** 

# 逻辑斯蒂回归-->多类分类



## ■ 多类分类

> 多元逻辑回归 (Multi-nomial Logistic Regression)

${\cal Y}_n$	1	2	 C
$p_{n}$	$p^1(\mathbf{x}_n;\mathbf{w}_1)$	$p^2(\mathbf{x}_n;\mathbf{w}_2)$	 $p^{C}\left(\mathbf{x}_{n};\mathbf{w}_{C}\right)$

$$p^{c}\left(\mathbf{x}_{n};\mathbf{w}_{c}\right) = \frac{\exp\left(\mathbf{w}_{c}^{T}\mathbf{x}_{n}\right)}{1 + \sum_{k=1}^{C-1} \exp\left(\mathbf{w}_{k}^{T}\mathbf{x}_{n}\right)}, \quad c = 1, 2, ..., C-1; \quad p^{C}\left(\mathbf{x}_{n};\mathbf{w}_{C}\right) = \frac{1}{1 + \sum_{k=1}^{C-1} \exp\left(\mathbf{w}_{k}^{T}\mathbf{x}_{n}\right)}$$

### > Softmax回归

$$p^{c}(\mathbf{x}_{n}; \mathbf{w}_{c}) = \frac{\exp(\mathbf{w}_{c}^{T} \mathbf{x}_{n})}{\sum_{k=1}^{C} \exp(\mathbf{w}_{k}^{T} \mathbf{x}_{n})}, c = 1, 2, ..., C$$

$$\mathbf{W} = \left[\mathbf{w}_1, \mathbf{w}_2, ..., \mathbf{w}_C\right]$$

$$\mathbf{y}_n = \left[ \delta(y_n = 1), \delta(y_n = 2), ..., \delta(y_n = C) \right]^T$$

$$h(\mathbf{W}^T\mathbf{x}) = \frac{\exp(\mathbf{W}^T\mathbf{x})}{\mathbf{1}^T \exp(\mathbf{W}^T\mathbf{x})}$$

$$L(\mathbf{W}) = -\frac{1}{N} \sum_{n=1}^{N} (\mathbf{y}_n)^T \log h(\mathbf{W}^T \mathbf{x}_n)$$

交叉熵损失函数

(Cross Entropy Loss)

• https://www.bilibili.com/video/av71259996

## 多分类学习





## 多分类学习

## Multiclass Classification

# 多分类学习-->引言



## ■二分类-->多分类

- ▶ 通常开始讨论分类问题时大多讨论两类问题 (如线性回归、SVM、Adaboost等),但是实际问题往往是多类的。
- ➤ 线性回归中可以将输出由一个值([0,1]或[-1,1])扩展一个指示向量(Indicator Vector)来处理多类问题:

$$\mathbf{y}_n = \left[\delta(y_n = 1), \delta(y_n = 2), ..., \delta(y_n = C)\right]^T$$



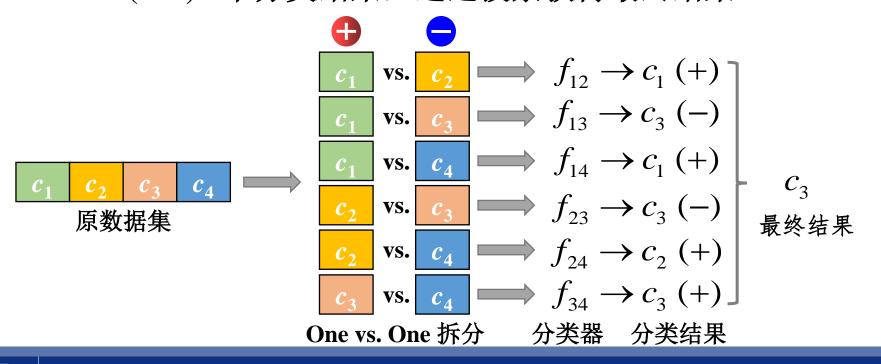
## ■通用方案

- > 存在利用二分类学习器解决多分类问题的通用方案?
  - ✓ 一对一 (One vs. One, OvO)
  - ✓ 一对其余 (One vs. Rest, OvR)
  - ✓ 多对多 (Many vs. Many, MvM)

# 多分类学习-->一对一



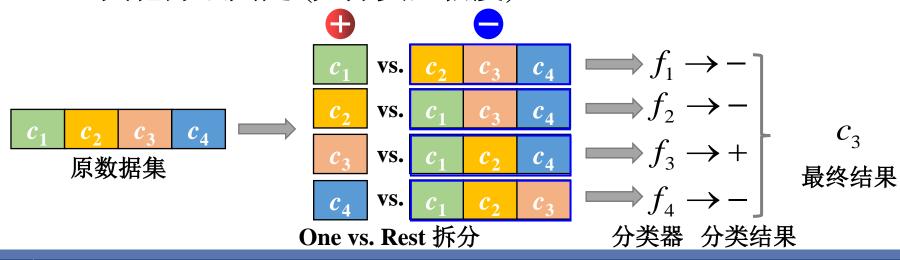
- 一对一 (One vs. One, OvO)
  - > <u>训练</u>:将C类分类问题拆分为C(C-1)/2个二分类问题,训练C(C-1)/2个分类器 (每个分类器区分两个类别);
  - $\ge$  <u>测试</u>:利用所有分类器对待测试样本进行分类,获得 C(C-1)/2个分类结果,通过投票获得最终结果。



# 多分类学习-->一对其余/一对多



- 一对其余 (One vs. Rest, OvR)
  - ▶ <u>训练</u>:将C类分类问题拆分为C个二分类问题,训练C个 分类器(每次将一类作为正类,其余作为负类);
  - ▶ <u>测试</u>: 利用所有分类器对待测试样本进行分类,获得C 个分类结果。若结果仅有一个结果为正类,则对应的类 别标记作为最终分类结果。若存在多个正类,则需要用 其他方法判定 (如分类置信度)。

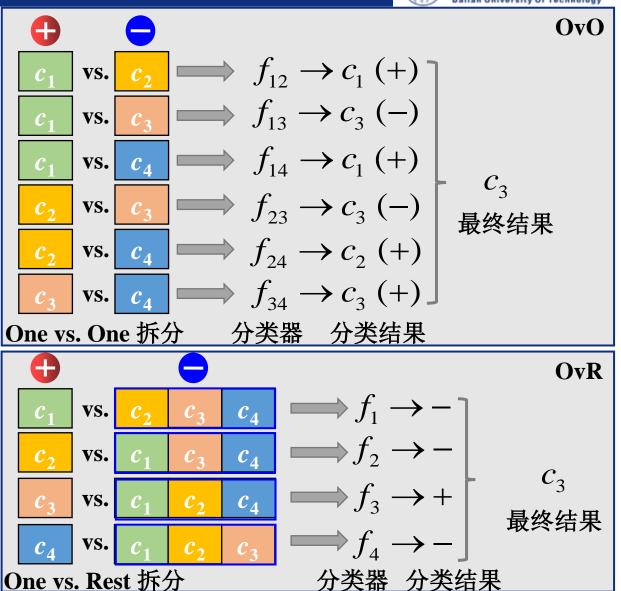


## 多分类学习-->OvO vs. OvR



- ▶ OvO的存储开销和测试 时间开销通常比OvR大:
   OvR只需训练C个分类 器,而OvO需训练C(C-1)/2个分类器。
- ➤ 类别多时,OvO的训练时间开销通常比OvR小: 时间开销通常比OvR小: 训练时,OvR的每个分 类器均使用全部训练样 本,而OvO的每个分类 器仅用到两个类样本;
- 预测性能差不多:至于 预测性能,则取决于具 体的数据分布,在多数 情形下两者差不多。

 $c_1$   $c_2$   $c_3$   $c_4$ 



# 多分类学习-->多对多

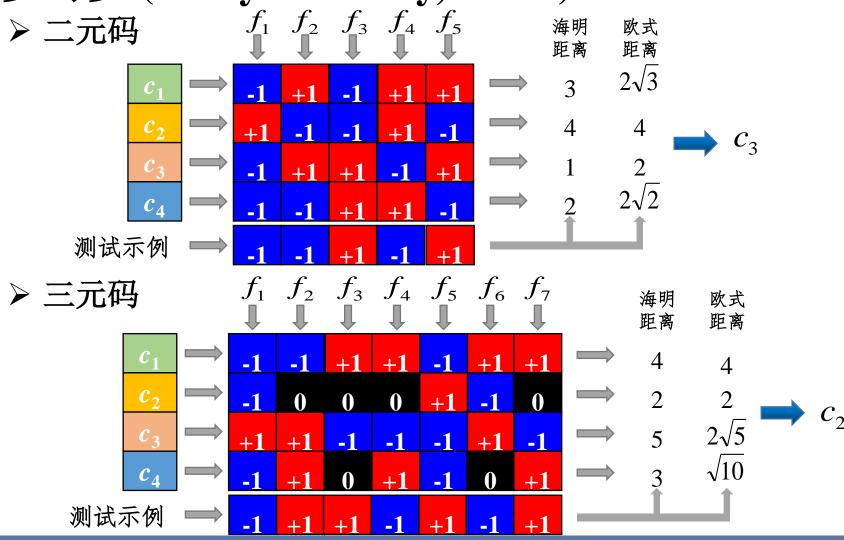


- 多对多 (Many vs. Many, MvM)
  - ▶ 基本思想:每次将若干类作为正类,其余类作为负类。 纠错输出编码 (Error-Correcting Output Codes, ECOC);
  - 编码:对C个类做M次划分,每次划分将一部分类划为正类,另一部分划为负类,从而形成一个二分类训练集。这样一共有M个训练集,可以训练出M个分类器;
  - 解码: M个分类器分别对测试样本进行分类,这些分类标记组成一个编码,将这个分类编码与每个类别各自的编码进行比较,返回距离最小的类作为最终预测结果。

# 多分类学习-->多对多



■ 多对多 (Many vs. Many, MvM)



# 类别不平衡问题





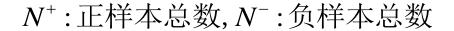
# 类别不平衡问题 Class-Imbalance

# 类别不平衡问题-->引言

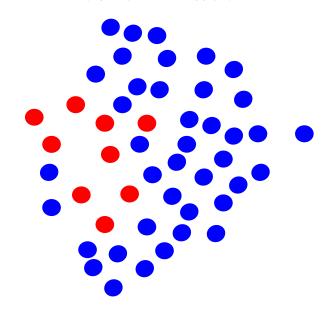


## ■ 类别不平衡 (Class-Imbalance)

- > 分类任务中不同类别的训练样本数目差别很大的情况
  - ✓ 以二分类为例
  - ✓ 假设正样本很少,负样本很多
  - ✓ 类别严重不平衡时,分类器学习 将过分关注数目多的类别
- ➤ 再缩放策略 (rescaling)
  - ✓ 阈值移动 (threshold-moving)
  - ✓ 过采样 (oversampling)
  - ✓ 欠采样 (undersampling)
- ➤ 阈值移动 (threshold-moving)





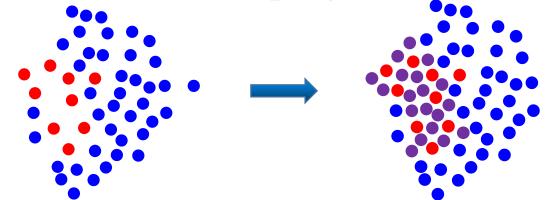


$$\frac{p}{1-p} > \frac{N^+}{N^-}$$
  $\Longrightarrow$   $\frac{p'}{1-p'} = \frac{p}{1-p} \frac{N^-}{N^+} > 1$ 

# 类别不平衡问题-->引言

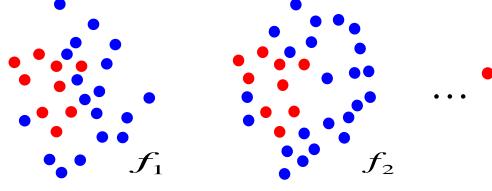


- 类别不平衡 (Class-Imbalance)
  - > 过采样 (oversampling):



- ✓ 样本复制
- ✓ 样本插值
- ✓ 样本生成 (GAN)
  - [1] Chawla N V, Bowyer K W, et al. **SMOTE: Synthetic Minority Over-Sampling Technique**. *JAIR*, 2002.

> 欠采样 (undersampling)



- 集成学习  $f = \frac{1}{M} \sum_{i=1}^{M} f_{M}$
- ✓ EasyEnsemble<sup>[2]</sup>
- ✓ BalanceCascade<sup>[2]</sup>
- [2] Xu-Ying Liu, Jianxin Wu, Zhi-Hua Zhou. **Exploratory Undersampling for Class-Imbalance Learning**. *IEEE TSMCB*, 2009.

## 类别不平衡问题-->加权损失



## ■ 加权损失函数 (Weighted Loss Function)

> 一般损失函数

$$L(\mathbf{X}; \mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} L(\mathbf{x}_n, y_n; \mathbf{w}) = \sum_{n=1}^{N} \frac{1}{N} L(\mathbf{x}_n, y_n; \mathbf{w})$$
$$= \sum_{n=1}^{N^+} \frac{1}{N} L(\mathbf{x}_n, y_n; \mathbf{w}) + \sum_{n=1}^{N^-} \frac{1}{N} L(\mathbf{x}_n, y_n; \mathbf{w})$$

N:样本总数,  $N^{+}:$ 正样本总数,  $N^{-}:$ 负样本总数  $(N = N^{+} + N^{-})$ 

> 加权损失函数

$$L(\mathbf{X}; \mathbf{w}) = \sum_{n=1}^{N} \rho_n L(\mathbf{x}_n, y_n; \mathbf{w})$$

$$= \sum_{n=1}^{N^+} \rho_n^+ L(\mathbf{x}_n, y_n; \mathbf{w}) + \sum_{n=1}^{N^-} \rho_n^- L(\mathbf{x}_n, y_n; \mathbf{w})$$

$$\rho^+ = \frac{N^-}{N}, \rho^- = \frac{N^+}{N} \qquad \Longrightarrow \qquad \frac{\rho^+}{\rho^-} = \frac{N^-}{N^+}$$

## 类别不平衡问题-->Focal Loss

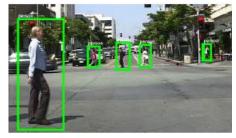


### ■ Focal Loss<sup>[1]</sup>

> 交叉熵损失

$$CE(p_t) = -\log(p_t)$$

$$p_{t} = \begin{cases} p & \text{if } y = 1\\ 1 - p & \text{otherwise} \end{cases}$$



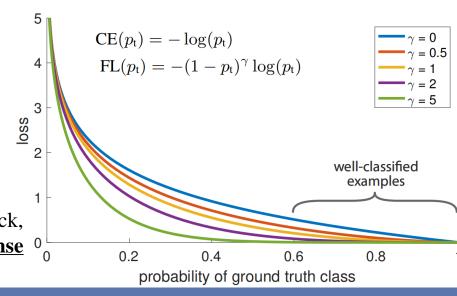
$$L = \frac{1}{N} \sum_{n=1}^{N} CE(p_n^t) = -\frac{1}{N} \sum_{n=1}^{N} y_n \log(p_n) + (1 - y_n) \log(1 - p_n)$$

➤ Focal Loss<sup>[1]</sup> ← 解决目标检测问题中训练集正负样本极度不平衡情况

$$FL(p_t) = -\underline{(1-p_t)^{\gamma}} \log(p_t)$$

- ✓ 自适应样本加权
- ✓ 容易分类样本权重降低  $p_t \to 1 \Rightarrow (1 p_t)^{\gamma} \to 0$
- ✓ 促进分类器学习更关注 困难样本

[1] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, Piotr Dollar. **Focal Loss for Dense** Object Detection, *ICCV*, 2017.

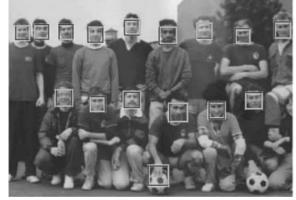


# 类别不平衡问题-->级联分类器



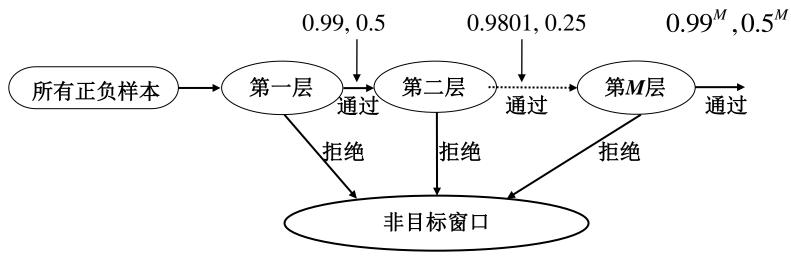
## ■ 级联分类器 (Cascaded Classifier)[1]

- > 人脸检测问题
  - ✓ 正样本 (人脸)
  - ✓ 负样本(非人脸)
  - ✓ 正样本远少于负样本



[1] Paul A. Viola, Michael J. Jones: **Robust Real- Time Face Detection**. *IJCV*, 2004.

➤ 设每层的正确率(TPR)为99%, 误检率(FPR)为50%



ightharpoonup 假设M=10,则可保证在正确率达到90.4%,而误检率为0.1%





## 小结

**Summary** 

# 小结



### ■小结

- ▶线性回归
  - ✓最小二乘法(目标函数,优化求解,几何意义,概率意义)
  - ✓最小二乘法存在问题 (Outliers), 线性理解, 广义线性回归
- ▶正则化
  - ✓ 正则化->过拟合
- >逻辑斯蒂回归
  - ✓从回归到分类,逻辑斯蒂函数,概率解释
  - ✓目标函数,优化求解,多类逻辑斯蒂回归
- >多分类学习
  - ✓直接建立多类模型(指示向量)
  - ✓通用策略(一对一,一对其余,多对多)
- >类别不平衡问题
  - ✓再缩放(阈值移动,欠采样,过采样) 级联分类器
  - ✓加权损失函数(正负样本重加权, Focal Loss)
  - ✓代价敏感损失函数(参见最小风险贝叶斯决策)