

基于 libsvm

- 总样本数**28056**，其中正样本**2796**，负样本**25260**。
- 随机取**5000**个样本训练，其余测试。
- 样本归一化，在训练样本上，求出每个维度的均值和方差，在训练和测试样本上同时归一化。

$$\text{new}X = \frac{X - \text{mean}(X)}{\text{std}(X)}$$

- 高斯核
- **5-fold cross validation**，在 **CScale = [2⁻⁵, 2¹⁵]**; **gamma = [2⁻¹⁵, 2³]** 上遍历求识别率的最大值。

上述**C**和**gamma**的区间设置参见**LIBSVM**自带的介绍：

a practical guide to support vector classification

- 训练参数设置 **svmtrain(yTraining, xTraining, cmd)**

cmd参数如下：

Svm - type

(1) -s 0 "-s svm_type : set type of SVM (default 0)\n"

" 0 -- C-SVC (multi-class classification)\n"

" 1 -- nu-SVC (multi-class classification)\n"

" 2 -- one-class SVM\n"

" 3 -- epsilon-SVR (regression)\n"

" 4 -- nu-SVR (regression)\n"

(2) -t 2

kernel - type

"-t kernel_type : set type of kernel function (default 2)\n"

" 0 -- linear: u*v\n"

" 1 -- polynomial: (gamma*u'*v + coef0)^degree\n"

" 2 -- radial basis function: exp(-gamma*|u-v|^2)\n"

" 3 -- sigmoid: tanh(gamma*u'*v + coef0)\n"

" 4 -- precomputed kernel (kernel values in training_instance_matrix)\n"

$\rightarrow K(x, x_2)$

(3) -c CVALUE

"-c cost : set the parameter C of C-SVC, epsilon-SVR, and nu-SVR (default 1)\n"

输入矩阵 (5000 个样本)

$$\begin{bmatrix} K(x_1, x_1), \dots, K(x_1, x_{5000}) \\ \vdots \\ K(x_{5000}, x_1), \dots, K(x_{5000}, x_{5000}) \end{bmatrix}$$

$K_{5000 \times 5000}$

(4) -g gammaValue

"-g gamma : set gamma in kernel function (default 1/num_features)\n"

(5) -v 5

"-v n : n-fold cross validation mode\n"

- 训练后获得的参数

(1) C = 16, gamma = 0.0825

(2) 支持向量 (即alpha不为0的向量) : 358个 (162个正样本, 196个负样本)

(3) b = 6.2863

对于一个测试样本 x :

$$\text{若 } \sum_{i=1}^n \alpha_i y_i K(x_i, x) + b \geq 0 \quad y = 1$$

经验：支持向量在总样本 20% ~ 30%

否则 $y = -1$

交叉验证 (Cross Validation)

5000个样本分为5组 每组1000个



① (a b c d) 训练 e 测试

⋮

C_5^4 训练 $5 - C_5^4$ 测试

总样本 28056 正样本 2796 负样本 25260

随机取 5000 个样本训练 其余测试

样本归一化 在训练样本上 求出每个维度的均值和方差
在训练和测试样本上同时归一化

$$\text{new } X = \frac{X - \text{mean}(X)}{\text{std}(X)}$$

高斯核

5-fold cross validation 在

$$C \text{ scale} = [2^{-5}, 2^{15}]$$

$$\gamma = [2^{-15}, 2^3]$$

混淆矩阵

		预测	
		正样本	负样本
实 正样本		TP	FN
实 负样本		FP	TN

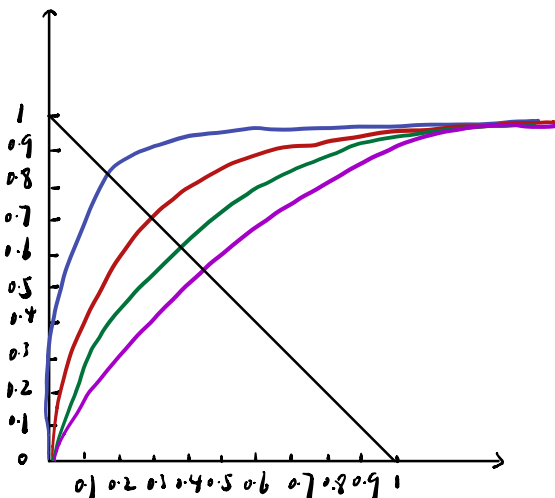
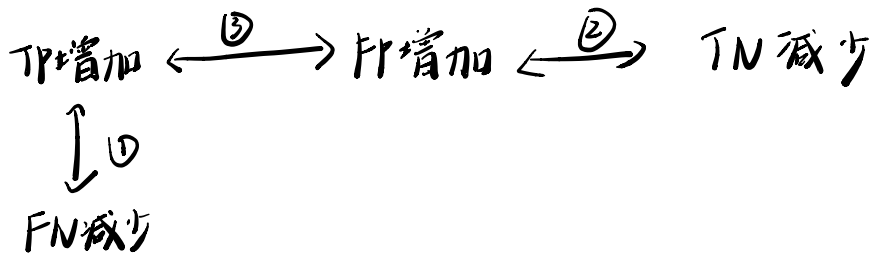
最好不要发生! $FP=0$ TP 越好 TP

ROC曲线 (Receiver Operating Character)

$$TP + FN = 1 \quad ①$$

$$FP + TN = 1 \quad ②$$

对同一个系统来说 若 TP 增加 FP 也增加 ③



等错误率 (Equal Error Rate, EER)

是两类错误 FP 和 FN 相等时的错误率

可以直观地表示系统性能