

1) 以多组不同参数值初始化多个神经网络, 按标注方法训练后, 取其中误差最小的作为最终参数

2) 模拟退火, 每一步以一定概率接受比当前更差的结果

3) 种群随机梯度下降

4 其他常见神经网络

1) RBF (Radial Basis Function) 径向基函数

$$\phi(x) = \sum_{i=1}^n W_i p(x, c_i) \quad p(x, c_i) = e^{-\beta_i \|x - c_i\|^2}$$

确定神经元中心  $c_i$ , 再利用BP算法确定  $W_i, \beta_i$

2) ART (adaptive resonance theory) 自适应谐振理论, 竞争型学习  
winner-take-all

3) SOM (self-organizing map) 自组织映射: 每个输出层神经元含4个连接权重, 自身携带的权向量之间的距离, 距离最近的神经元成为竞争获胜者, beat matching unit 最佳匹配单元

## §6 支持向量机

1. 感知器模型

$$y = \text{sgn}(\theta \cdot x) = \begin{cases} 1, & \theta \cdot x > 0 \\ -1, & \theta \cdot x < 0 \end{cases} \quad \begin{matrix} y \theta x > 0 & \text{正确分类} \\ y \theta x < 0 & \text{错误分类} \end{matrix}$$

Loss function, 期望使所有样本到超平面的距离之和最小

$$L = \sum_{i=1}^k \frac{y_i \theta \cdot x_i}{\|\theta\|_2} \rightarrow L = - \sum_{i=1}^k y_i \theta \cdot x_i \quad \frac{\partial L(\theta)}{\partial \theta} = - \sum_{i=1}^k y_i x_i$$

线性可分 SVM: 让离超平面最近的点尽可能远离这个点 超平面 (separating hyperplane)  
非线性可分 SVM

支持向量 (support vector) 离分割超平面最近的点

间隔 (margin): 支持向量表点到分割超平面的距离

$$W^T x + b = 0$$

W - 2 倍

b - 1 倍

$$r = \frac{|W^T x + b|}{\|W\|}$$

$$r = \frac{2}{\|W\|} (\text{margin})$$

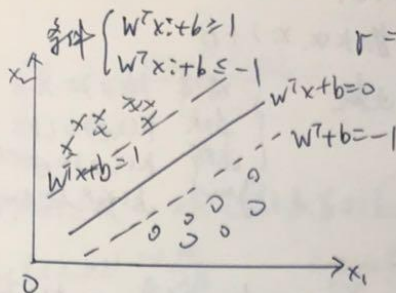
最大化问题

$$\max_{w, b} \frac{2}{\|W\|}$$

等价于

$$\min_{w, b} \frac{1}{2} \|W\|^2$$

$$y_i (W^T x_i + b) \geq 1$$



2. 对偶化

特性: 对偶问题的对偶是原问题; 无论原问题是凸, 对偶问题皆为凸优化;  
—— 可给出原问题的下界; 当满足一定条件, 原问题与对偶问题是完美匹配的。

拉格朗日乘子法:  $L(w, b, \alpha) = \frac{1}{2} \|W\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i (W^T x_i + b))$

对  $w, b$  求导为 0.  $w = \sum_{i=1}^n \alpha_i y_i x_i$   $b = \sum_{i=1}^n \alpha_i y_i$

$$\Rightarrow \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$f(w) = W^T x + b = \sum_{i=1}^n \alpha_i y_i x_i^T x + b$$

约束:  $\begin{cases} \alpha_i \geq 0 \\ y_i f(x_i) - 1 \geq 0 \\ \alpha_i (y_i f(x_i) - 1) = 0 \end{cases}$

优化这变量: SMO (Sequential Min Q) 算法  
→ 选择两个条件中违反最大之量, 使时间最短  
最快, 使所选两变量所对应样本间隔最大

3. 核函数

非线性可分 SVM 将原始 n 维空间映射到高维空间

$$f(x) = W^T \phi(x) + b \quad \min_{w, b} \frac{1}{2} \|W\|^2$$

$$\Rightarrow \max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \phi(x_i)^T \phi(x_j)$$

$$k(x_i, x_j) = (\phi(x_i), \phi(x_j)) = \phi(x_i)^T \phi(x_j)$$

不必计算高维空间中的内积

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j k(x_i, x_j)$$

$$\Rightarrow f(w) = \sum_{i=1}^n \alpha_i y_i k(x, x_i) + b$$

核函数必须是正定核函数

线性	$k(x, z) = x \cdot z$
多项	$k(x, z) = (x \cdot z + r)^d$
高斯	$k(x, z) = e^{-r \ x - z\ ^2}$
Sigmoid	$k(x, z) = \tanh(r x \cdot z + r)$

4. 软间隔与正则化

SVM 线性可分数据. 但对 outlier 敏感  $\rightarrow$  软间隔. 松弛因子  $\xi_i$

松弛因子  $\xi_i$  越大  $\rightarrow$  离超平面越远

损失函数:  $\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$

核函数  $L(w, b, \xi, \beta, \mu) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i + \sum_{i=1}^n \beta_i [1 - \xi_i - y_i (w^T x_i + b)] - \sum_{i=1}^n \mu_i \xi_i$

$\min_{w, b, \xi, \beta, \mu} L(w, b, \xi, \beta, \mu) \iff \max_{\beta, \mu} \min_{w, b, \xi} L(w, b, \xi, \beta, \mu)$

对  $w, b, \xi$  求导  $\Rightarrow \mu(\beta) = \sum_{i=1}^n \beta_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \beta_i \beta_j y_i y_j x_i^T x_j$

硬间隔  $w^*, b^* \quad f(x) = \text{sgn}(w^* x + b^*)$

软间隔 LR. (LR 损失函数为 sigmoid)

加入松弛因子后 模型更 robust. 泛化能力更强.

sk-learn 中 SVC (核函数), LinearSVC (线性), OneClassSVM (异常检测)

支持向量也可作回归. SVR 不常用

有的最优化问题: { 无约束条件: 梯度下降法, 牛顿法, 全牛顿下降法...  
有约束条件: 拉格朗日乘子法  
不光滑的: KKT