

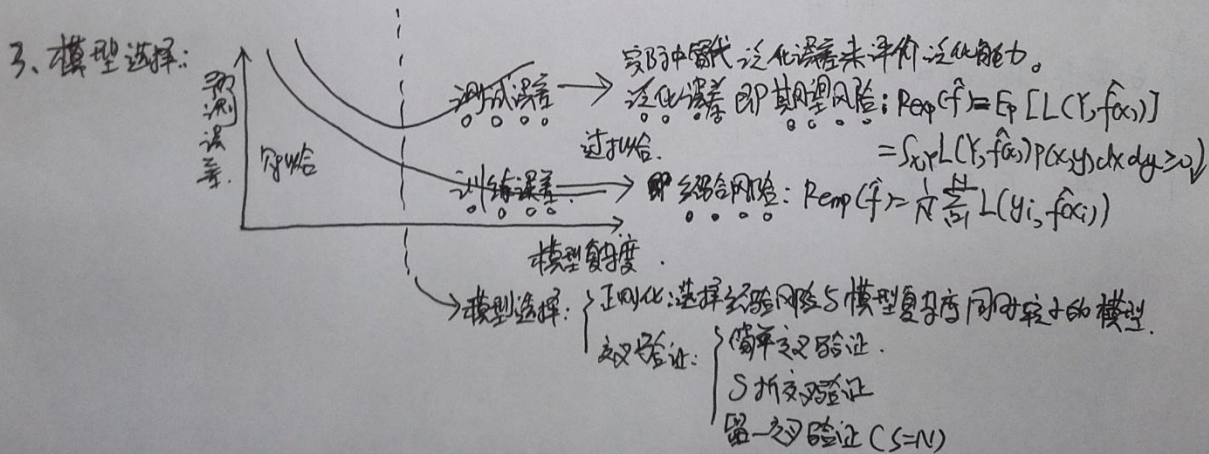
$\mathbb{E}[L(Y, f(x))] = \int_{\mathcal{Y}} \int_{\mathcal{X}} L(y, f(x)) P(X, Y) dx dy$ 即期望风险最小化

经验风险最小化: $\min \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i))$

极大似然估计

结构风险最小化: $\min \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)) + \lambda J(f)$

正则化估计



模型泛化误差与训练误差的关系: $R_{\text{exp}}(\hat{f}) \leq R_{\text{emp}}(\hat{f}) + \mathcal{O}(d, N, \delta)$

其中, $\mathcal{O}(d, N, \delta) = \sqrt{\frac{1}{2N} (\log d + \log \frac{1}{\delta})}$

$N \rightarrow \infty$ 即样本容量 $\rightarrow \infty \Rightarrow$ 泛化误差 $\rightarrow 0$

$d \uparrow$ 即假设空间 $\uparrow \Rightarrow$ 泛化误差 \uparrow

4. 生成法与判别法, 分类问题, 标注问题与回归问题.

生成法: 由 $P(X)$ 间接求 $P(Y|X)$, 如 NB, HMM

判别法: 直接求 $f(x)$ 或 $P(Y|X)$

可还原 $P(X, Y)$
 生成法: 慢
 可存在隐变量

准确率
 简化学习问题.

分类问题: 输入离散/连续, 输出离散 $\begin{cases} Y = f(x) \\ P(Y|X) \end{cases} \rightarrow \text{文本分类}$

标注问题: 输入、输出均为变量序列 $\begin{cases} Y = f(x) \\ P(Y|X) \end{cases} \rightarrow \text{词性标注}$
 如 HMM, CRF

回归问题: 输入、输出连续 $Y = f(x)$ $\rightarrow \text{股价预测}$
 常用平方损失函数 \Rightarrow 最小二乘法

Relevant

NonRelevant

Retrieved. true positives (tp)

false positives (fp)

Not Retrieved false negatives (fn)

true negatives (tn)

精确率 Precision = $\frac{tp}{tp + fp}$

召回率 Recall = $\frac{tp}{tp + fn}$

F1值: $\frac{2}{F1} = \frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}$

LDA & PCA

1. LDA的输入带标签, 监督学习, 也称 Fisher \rightarrow 使投影后的方差

对 $y = w^T x$ 有

类别 i 的原始中心点, $m_i = \frac{1}{n_i} \sum_{x \in P_i} x$

类别 i 投影后的中心点, $\tilde{m}_i = w^T m_i \Rightarrow |\tilde{m}_1 - \tilde{m}_2|^2 = w^T (m_1 - m_2)(m_1 - m_2)^T w = w^T S_B w$

投影后的方差 $S_i = \sum_{y \in P_i} (y - \tilde{m}_i)^2 = \sum_{x \in P_i} (w^T x - w^T m_i)^2 = \sum_{x \in P_i} w^T (x - m_i)(x - m_i)^T w = w^T S_i w$

损失函数 $J(w) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{S_1^2 + S_2^2} \uparrow \Rightarrow \begin{cases} \text{类间距离} \uparrow \\ \text{类内距离} \downarrow \end{cases} \quad S_1^2 + S_2^2 = w^T (S_1 + S_2) w = w^T S_w w$

$$= \frac{w^T S_B w}{w^T S_w w}$$

拉格朗日乘子法: $J(w) = \frac{w^T S_B w}{w^T S_w w} \uparrow$ 等价于 $\begin{cases} J(w) = w^T S_B w \uparrow \\ w^T S_w w = 1 \end{cases}$

所以 $C(w) = w^T S_B w - \lambda(w^T S_w w - 1) \Rightarrow \frac{dC}{d\lambda} = 2 S_B w - 2\lambda S_w w = 0 \Rightarrow S_B w = \lambda S_w w \Rightarrow$ 求特征值

对于多分类问题, $\begin{cases} S_B = \sum_{i=1}^C n_i (m_i - m)(m_i - m)^T \\ S_w = \sum_{i=1}^C S_i \end{cases} \Rightarrow S_B w = \lambda S_w w \Rightarrow$ 求特征值

其中 w 对应于特征值中最大特征向量!

2. PCA的输入不带标签, 无监督学习 \rightarrow 使投影后的方差

$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$

假设 u_1 为投影向量, 则投影后的方差 $\frac{1}{N} \sum_{n=1}^N (u_1^T x_n - u_1^T \bar{x})^2 = u_1^T \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T u_1 = u_1^T S u_1$

拉格朗日乘子法: $J(u_1) = u_1^T S u_1 \uparrow$ 所以 $C(u) = u^T S u - \lambda(u^T u - 1) \Rightarrow \frac{dC}{d\lambda} = 2 S u - 2\lambda u = 0$
 $\begin{cases} u^T u = 1 \text{ (已知)} \end{cases}$

$\Rightarrow S u = \lambda u \Rightarrow$ 求特征值, 因为 $\|u\| = 1$ 所以 u 为特征向量, λ 为特征值。若 $J(u) \uparrow \Rightarrow \lambda \uparrow$

所以 PCA 要求当 D 维数据空间投影到低维空间 ($M < D$), 则取前 M 个特征向量构成的投影矩阵, 即能够使得投影后的方差最大的数据。

感知机与支持向量机 (二分类)

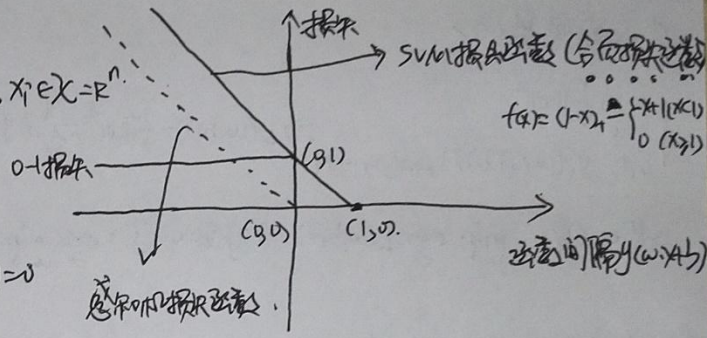
1. 函数间隔与几何间隔

函数间隔: $\hat{\gamma}_i = y_i(w \cdot x_i + b) \Rightarrow$ 超平面函数间隔: $\hat{\gamma} = \min_{i=1,2,\dots,N} \hat{\gamma}_i$
 几何间隔: $\gamma_i = y_i(\frac{w}{\|w\|} \cdot x_i + \frac{b}{\|w\|}) \Rightarrow$ 超平面几何间隔: $\gamma = \frac{\hat{\gamma}}{\|w\|} = \min_{i=1,2,\dots,N} \frac{\hat{\gamma}_i}{\|w\|}$

2. 感知机与支持向量机的输入、输出

输入: $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}, x_i \in X = \mathbb{R}^n, y_i \in Y = \{-1, +1\}, i=1,2,\dots,N$

输出: $\begin{cases} \text{分离超平面: } w^* \cdot x + b^* = 0 \\ \text{决策函数: } f(x) = \text{sign}(w^* \cdot x + b^*) = 0 \end{cases}$



3. 感知机学习: 误分类最小, 求分离超平面, 解不唯一。

误分类点

$$\min_{w,b} \sum_{x_i \in M} y_i (\frac{w}{\|w\|} \cdot x_i + \frac{b}{\|w\|}) \Rightarrow \min_{w,b} L(w,b), \text{ 且 } L(w,b) = - \sum_{x_i \in M} y_i (w \cdot x_i + b) \xrightarrow{\text{梯度下降法}} \begin{cases} w^* = \sum_{i=1}^N \alpha_i y_i x_i \\ b^* = \sum_{i=1}^N \alpha_i y_i \end{cases}$$

梯度下降法学习过程

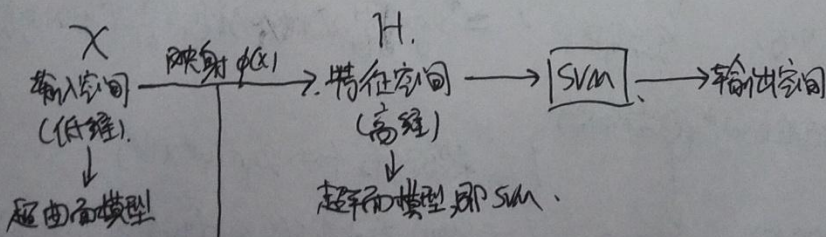
初始化 w, b

选取 $(x_i, y_i) \rightarrow \begin{cases} \text{如果 } y_i(w \cdot x_i + b) > 0, \text{ 不更新 } w, b \\ \text{如果 } y_i(w \cdot x_i + b) \leq 0, \text{ 则 } \begin{cases} w \leftarrow w + \eta y_i x_i \\ b \leftarrow b + \eta y_i \end{cases} \end{cases}$

直至训练集没有误分类点

超平面几何间隔

4. 支持向量机: 间隔最大化, 求分离超平面, 解唯一。



核函数 $K(x, z) = \phi(x) \cdot \phi(z)$, 本质是内积

一般给定输入空间和核函数 K , 不必定义 ϕ , 由 K 推导 ϕ , 但 ϕ 不唯一。

比如输入空间 \mathbb{R}^2 , $K(x, z) = (x \cdot z)^2$

则 $(x \cdot z)^2 = (x^{(1)}, x^{(2)} \cdot z^{(1)}, z^{(2)})^2$

$= (x^{(1)} z^{(1)} + x^{(2)} z^{(2)})^2$

$= (x^{(1)} z^{(1)})^2 + 2 x^{(1)} z^{(1)} x^{(2)} z^{(2)} + (z^{(1)} z^{(2)})^2$

对特征空间,

若 $\mathbb{R}^3 \Rightarrow \phi(x) = (x^2, x^2, x^2)$

若 $\mathbb{R}^4 \Rightarrow \phi(x) = (x^2, x^2, x^2, x^2)$

支持向量机 { 数据线性可分 \Rightarrow 硬间隔最大化 } 线性支持向量机
 { 数据近似线性可分 \Rightarrow 软间隔最大化 }
 { 数据线性不可分 \Rightarrow 核函数及软间隔最大化 }

超平面间隔最大化

$$\max_{w,b} \min_{i=1,2,\dots,N} y_i \left(\frac{w}{\|w\|} \cdot x_i + \frac{b}{\|w\|} \right) \Rightarrow \begin{cases} \max_{w,b} \gamma \\ \text{s.t. } y_i (w \cdot x_i + b) \geq \gamma, i=1,2,\dots,N \end{cases}$$

其中 $\gamma = \min_{i=1,2,\dots,N} y_i (w \cdot x_i + b)$

\Rightarrow 分情况讨论

① 硬间隔最大化

$$\begin{cases} \min_{w,b} \frac{1}{2} \|w\|^2 \\ \text{s.t. } y_i (w \cdot x_i + b) \geq 1, i=1,2,\dots,N \end{cases} \Rightarrow L(w,b) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^N \alpha_i y_i (w \cdot x_i + b) + \sum_{i=1}^N \alpha_i$$

\Rightarrow 原始问题: $\min_{w,b} \max_{\alpha} L(w,b,\alpha) \Rightarrow$ 对偶问题: $\max_{\alpha} \min_{w,b} L(w,b,\alpha) \Rightarrow$

$$\begin{cases} w^* = \sum_{i=1}^N \alpha_i^* y_i x_i \\ b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i (x_i \cdot x_j) \end{cases}$$

$\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ 且存在 $\alpha_j^* > 0$ (支持向量)

② 软间隔最大化 \rightarrow 惩罚因子

$$\begin{cases} \min_{w,b,\xi} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t. } y_i (w \cdot x_i + b) \geq 1 - \xi_i \\ \xi_i \geq 0, i=1,2,\dots,N \end{cases} \Rightarrow L(w,b,\xi,\alpha,u) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [y_i (w \cdot x_i + b) - 1 + \xi_i] - \sum_{i=1}^N u_i \xi_i$$

$\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ 且存在 $0 < \alpha_j^* < C$ (支持向量)

\Rightarrow 原始问题: $\min_{w,b,\xi} \max_{\alpha,u} L(w,b,\xi,\alpha,u) \Rightarrow$ 对偶问题: $\max_{\alpha,u} \min_{w,b,\xi} L(w,b,\xi,\alpha,u) \Rightarrow$

$$\begin{cases} w^* = \sum_{i=1}^N \alpha_i^* y_i x_i \\ b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i (x_i \cdot x_j) \end{cases}$$

③ 核函数及软间隔最大化, 核函数 $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$

$x_i \cdot x_j \rightarrow \phi(x_i) \cdot \phi(x_j) = K(x_i, x_j)$, 核函数使 SVM 学习隐式地在特征空间进行。

外解: $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ 且存在 $0 < \alpha_j^* < C$ (支持向量)

$$\begin{cases} w^* = \sum_{i=1}^N \alpha_i^* y_i x_i \\ b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i (x_i \cdot x_j) \end{cases} \Rightarrow \text{决策函数 } f(x) = \text{sign}(w^* \cdot x + b^*)$$

$$= \text{sign} \left[\sum_{i=1}^N \alpha_i^* y_i K(x_i, x) + \left(y_j - \sum_{i=1}^N \alpha_i^* y_i (x_i \cdot x_j) \right) \right]$$

惩罚因子

现在解: $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_N^*)^T$ 且存在 $0 < \alpha_j^* < C$ (支持向量)

$$\begin{cases} w^* = \sum_{i=1}^N \alpha_i^* y_i x_i \\ b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i K(x_i, x_j) \end{cases} \Rightarrow \text{决策函数 } f(x) = \text{sign}(w^* \cdot x + b^*)$$

$$= \text{sign} \left[\sum_{i=1}^N \alpha_i^* y_i K(x_i, x) + \left(y_j - \sum_{i=1}^N \alpha_i^* y_i K(x_i, x_j) \right) \right]$$

核函数 核函数

常用核函数 { 多项式核函数
高斯核函数
字符串核函数

朴素贝叶斯法 (多类分类) \rightarrow 生成模型, 因为 $P(Y|X) = \frac{P(X,Y)}{P(X)}$

1. 朴素贝叶斯法的学习与分类.

Bayes 定理

条件独立性假设 (使模型中特征的概率独立)

已知: 先验概率 $P(Y)$, 求后验概率 $P(Y|X) = \frac{P(X,Y)}{P(X)} = \frac{P(Y)P(X|Y)}{P(X)} = \frac{P(Y) \prod_j P(x^{(j)}|Y)}{P(X)}$

具体: $\begin{cases} P(Y=c_k) \\ P(X=x|Y=c_k) \end{cases} \quad P(Y=c_k|X=x) = \frac{P(Y=c_k) \prod_j P(X^{(j)}=x^{(j)}|Y=c_k)}{P(X)}$

因为后验概率最大化等价于在选择 0-1 损失函数时的期望风险最小化

所以求得后验概率最大的输出 $y = \arg \max_{c_k} P(Y=c_k) \prod_j P(X^{(j)}=x^{(j)}|Y=c_k)$

2. 朴素贝叶斯法的参数估计.

极大似然估计 (等价于经验风险最小化)

1. 先验概率的极大似然估计: $P(Y=c_k) = \frac{\sum_{i=1}^N I(y_i=c_k)}{N}$
2. 对于给定的属性 $x = (x^{(1)}, x^{(2)}, \dots, x^{(m)})^T$, 计算
条件概率的极大似然估计: $P(X^{(j)}=x^{(j)}|Y=c_k) = \frac{\sum_{i=1}^N I(x_i^{(j)}=x^{(j)}, y_i=c_k)}{\sum_{i=1}^N I(y_i=c_k)}$
3. 确定实例 x 的类: $y = \arg \max_{c_k} P(Y=c_k) \prod_{j=1}^n P(X^{(j)}=x^{(j)}|Y=c_k)$

贝叶斯估计 (等价于结构风险最小化)

1. 先验概率的贝叶斯估计: $P_\lambda(Y=c_k) = \frac{\sum_{i=1}^N I(y_i=c_k) + \lambda}{N + (K)\lambda}$
 $\rightarrow K$ 为类标记可能的取值个数
2. 条件概率的贝叶斯估计
 $P_\lambda(X^{(j)}=x^{(j)}|Y=c_k) = \frac{\sum_{i=1}^N I(x_i^{(j)}=x^{(j)}, y_i=c_k) + \lambda}{\sum_{i=1}^N I(y_i=c_k) + (S_j)\lambda}$
 $\rightarrow S_j$ 为第 j 个特征可能的取值个数
3. 同上.

最大熵模型与 Logistic 回归 (连续分类)

1. 最大熵原理认为, 熵最大的模型是最好的模型. 是概率模型学习或估计的一个准则!

变量 X 的熵: $H(X) = -\sum_{i=1}^n P_i \log P_i$ 或 $H(P) = -\sum_{i=1}^n P_i \log P_i$

给定 X, Y 的联合熵: $H(X, Y) = -\sum_{i=1}^n P_i \cdot H(Y|X=x_i) = -\sum_{i=1}^n P_i \left[-\sum_{j=1}^m P(y_j|x_i) \log P(y_j|x_i) \right]$
 $= -\sum_{i=1}^n \sum_{j=1}^m P_i \cdot P(y_j|x_i) \log P(y_j|x_i)$
 或 $H(P) = -\sum_{i=1}^n \sum_{j=1}^m P_i \cdot P(y_j|x_i) \log P(y_j|x_i)$

模型 $P(Y|X)$ 的熵: $H(P) = -\sum_{x,y} \tilde{P}(x) \cdot P(y|x) \log P(y|x)$

最大熵模型 $\begin{cases} \max_{P \in C} H(P) = -\sum_{x,y} \tilde{P}(x) P(y|x) \log P(y|x) \\ \text{s.t. } E_P(f_i) = E_{\tilde{P}}(f_i), i=1,2,\dots,n \text{ 即 } \sum_{x,y} \tilde{P}(x) P(y|x) f_i(x,y) = \sum_{x,y} \tilde{P}(x,y) f_i(x,y) \\ \sum_y P(y|x) = 1 \end{cases}$

\Leftrightarrow

$\begin{cases} \min_{P \in C} -H(P) = \sum_{x,y} \tilde{P}(x) P(y|x) \log P(y|x) \\ \text{s.t. } E_P(f_i) - E_{\tilde{P}}(f_i) = 0, i=1,2,\dots,n \\ \sum_y P(y|x) = 1 \end{cases}$

$L(P, w) = -H(P) + w_0 \left(1 - \sum_y P(y|x)\right) + \sum_{i=1}^n w_i (E_{\tilde{P}}(f_i) - E_P(f_i))$

$P_w(y|x) = \frac{1}{Z_w(x)} \exp\left(\sum_{i=1}^n w_i f_i(x,y)\right)$
 其中 $Z_w(x) = \sum_y \exp\left(\sum_{i=1}^n w_i f_i(x,y)\right)$

原问题是: $\min_{P \in C} \max_w L(P, w)$

对偶问题是: $\max_w \min_{P \in C} L(P, w) \Rightarrow$

$\left. \begin{aligned} &\text{① 求解 } \min_{P \in C} L(P, w) \text{ 即 } L(P, w) \text{ 关于 } P \text{ 的最小值} \\ &\text{即求最大熵模型公式 } P_w(y|x) \text{ 有 } \min_{P \in C} L(P, w) = L(P_w, w) \\ &\text{② 求解 } \max_w L(P_w, w) \text{ 即 } L(P_w, w) \text{ 关于 } w \text{ 的最大值} \\ &\text{即最大熵模型的学习, 得到 } w^* \end{aligned} \right\} \Rightarrow \text{最大熵模型 } P_w(y|x)$

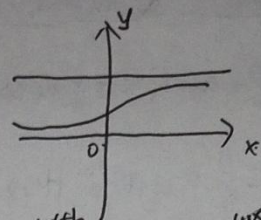
此步等价于 $P_w(y|x)$ 的最大似然估计 $L_P(P_w) = \log \prod_{x,y} P_w(y|x)^{\tilde{P}(x,y)} = \sum_{x,y} \tilde{P}(x,y) \log P_w(y|x)$

因此最大熵模型的学习归结为最大似然估计问题!

求解最大似然估计问题, 得 w^* $\left\{ \begin{aligned} &\text{迭代优化法} \\ &\text{梯度下降法} \\ &\text{牛顿法或拟牛顿法} \end{aligned} \right.$

2. Logistic回归

Logistic分布 $f(x) = \frac{1}{1 + e^{-\theta x}}$ $\begin{cases} x \rightarrow +\infty \Rightarrow f(x) \rightarrow 1 \\ x \rightarrow -\infty \Rightarrow f(x) \rightarrow 0 \end{cases}$



假设 $P(Y|X)$ 服从 Logistic 分布, 有 $\begin{cases} z(x) = P(Y=1|X) = \frac{e^{w_0 + w_1 x}}{1 + e^{w_0 + w_1 x}} \rightarrow \frac{e^{w_1 x}}{1 + e^{w_1 x}} \Rightarrow \frac{z(x)}{1 - z(x)} = e^{w_1 x} \\ 1 - z(x) = P(Y=0|X) = \frac{1}{1 + e^{w_0 + w_1 x}} \rightarrow \frac{1}{1 + e^{w_1 x}} \end{cases}$

学习问题采用极大似然估计, 有 $L(w) = \log \prod_{i=1}^n [z(x_i)]^{y_i} [1 - z(x_i)]^{1 - y_i}$

$$= \sum_{i=1}^n (y_i \log z(x_i) + (1 - y_i) \log [1 - z(x_i)])$$

$$= \sum_{i=1}^n (y_i \log \frac{z(x_i)}{1 - z(x_i)} + \log [1 - z(x_i)])$$

$$= \sum_{i=1}^n [y_i (w \cdot x_i) - \log (1 + e^{w \cdot x_i})]$$

求使 $L(w)$ 取得极大值的 $w \Rightarrow$ $\begin{cases} \text{梯度下降法} \\ \text{牛顿法或拟牛顿法} \end{cases}$

梯度下降法: $\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial w} \cdot \frac{\partial w}{\partial w_i}$

$$= \sum_{i=1}^n (y_i \cdot x_i - \frac{1}{1 + e^{w \cdot x_i}} \cdot e^{w \cdot x_i} \cdot x_i) \cdot \frac{\partial w}{\partial w_i}$$

$$= \sum_{i=1}^n (y_i \cdot x_i - z(x_i) \cdot x_i) \cdot \frac{\partial w}{\partial w_i} = \sum_{i=1}^n (y_i - z(x_i)) x_i \cdot \frac{\partial w}{\partial w_i}$$

$$= (y_i - z(x_i)) x_i$$

$\therefore w_i \leftarrow w_i + (y_i - z(x_i)) x_i$

k近邻法 (分类)

1. k近邻算法 { 输入: $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$; 实例特征向量 x
 输出: 实例 x 所属的 y

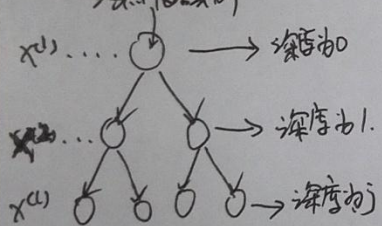
(1) 根据给定的距离度量, 在训练集 T 中找出与 x 最近的 k 个点, 记作 $N_k(x)$

(2) 在 $N_k(x)$ 中根据分类决策规则 (如多数表决) 决定 x 的类别 $y = \arg \max_j \sum_{x_i \in N_k(x)} I(y_i = c_j)$

三要素: { 距离度量 $\Rightarrow L_p(x_i, x_j) = (\sum_{l=1}^n |x_i^{(l)} - x_j^{(l)}|^p)^{1/p}$, 欧氏距离 $\hookrightarrow L_2(x_i, x_j) = (\sum_{l=1}^n |x_i^{(l)} - x_j^{(l)}|^2)^{1/2}$
 k值的选择 \Rightarrow $k \downarrow$, 其最近邻 \Rightarrow 模型复杂, 过拟合 \Rightarrow 交叉验证求 k
 $k \uparrow \Rightarrow$ 模型简单, 欠拟合.
 分类决策规则 \Rightarrow 多数表决等价于在训练集-1损失函数下的经验风险最小化.

2. k近邻搜索: { 线性扫描 (相当于全局搜索)
 kd树 (相当于局部搜索): 适用于训练集 \Rightarrow 空间维数 k (与邻居中 k 含义不同)

构造kd树: { 输入: $T = \{x_1, x_2, \dots, x_n\}$, 其中 $x_i = (x_i^{(1)}, x_i^{(2)}, \dots, x_i^{(k)})^T$
 输出: kd树



对于 $x^{(1)}, x^{(2)}, \dots, x^{(n)}$, 分别对其第 k 维坐标的中位数 (平衡) 进行超平面划分, 将落在划分超平面上的实例点, 保存在对应结点。直到两个子区域都没有实例存在时停止。

其中 $1 \leq j \leq k$: $k = j \pmod k + 1$.
 kd树最深为 k 层 (因为 k 维空间)

用kd树的最近邻搜索: { 输入: kd树; 目标点 x
 输出: x 的最近邻

从根结点出发, 递归向下访问kd树, 直到子结点为止, 以此作为“当前最近点”。

开始局部搜索: 回溯, 更新不满足“当前最近点”
 再看是否与另子结点相交, { 相交 \rightarrow 扫描.
 不相交 \rightarrow 回溯.
 直到根结点结束, 返回“当前最近点”。

决策树 (分类问题)

1. 决策树: 树形结构。内部结点: 对应特征或属性。叶结点对应类。

K类
特征A用于划分数据集D, 故将D划分为m个集合

数据集特征A

特征选择:

$$\text{信息增益: } g(D, A) = H(D) - H(D|A) = \left[-\sum_{k=1}^K \frac{|C_k|}{|D|} \log \frac{|C_k|}{|D|} \right] - \left[-\sum_{i=1}^m \frac{|D_i|}{|D|} H(D_i) \right]$$

$$= \left[-\sum_{k=1}^K \frac{|C_k|}{|D|} \log \frac{|C_k|}{|D|} \right] - \left[-\sum_{i=1}^m \frac{|D_i|}{|D|} \cdot \sum_{k=1}^K \frac{|D_{ik}|}{|D_i|} \log \frac{|D_{ik}|}{|D_i|} \right]$$

→ D中属于类C_k的样本集合 |C_k|

→ D_i中属于类C_k的样本集合 |D_{ik}|

信息增益比: $g_R(D, A) = \frac{g(D, A)}{H(D)} = \frac{g(D, A)}{-\sum_{i=1}^m \frac{|D_i|}{|D|} \log \frac{|D_i|}{|D|}}$

决策树学习

决策树的生成 (局部最优):

- ID3算法: 信息增益选择特征
- C4.5算法: 信息增益比选择特征

⇒ 从根结点开始, 递归地产生决策树。

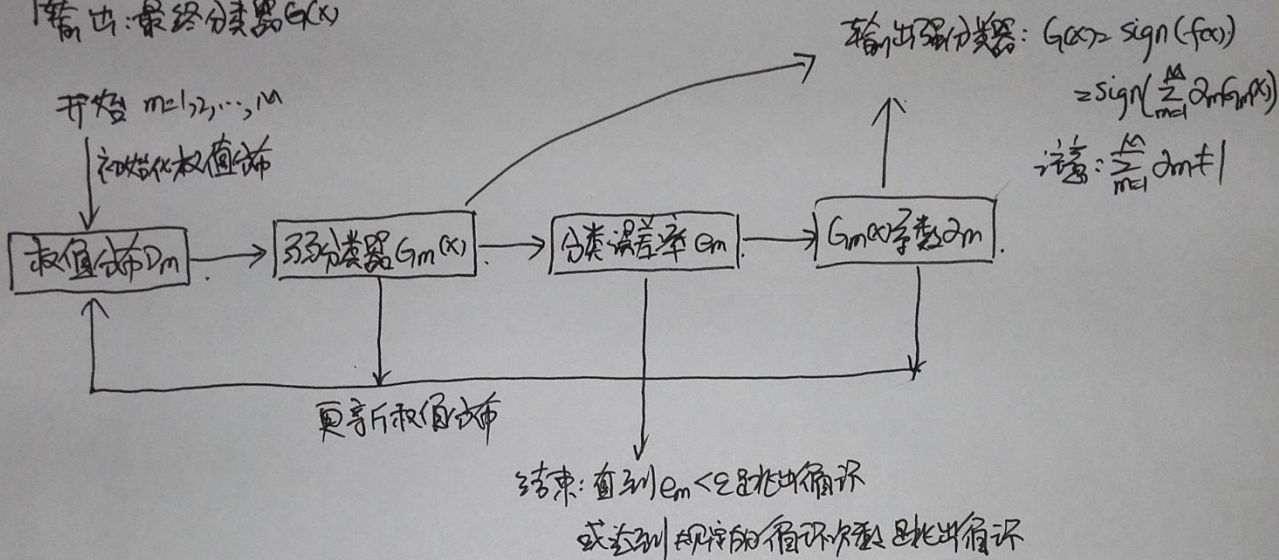
决策树的剪枝 (全局最优): 解决过拟合问题。

AdaBoost算法 (= 分类)

对于二分类问题:

输入: $T = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, $x_i \in X \subseteq \mathbb{R}^d$, $y_i \in Y = \{-1, +1\}$

输出: 最终分类器 $G(x)$



其中, $e_m = P(G_m(x_i) \neq y_i) = \sum_{i=1}^N w_{mi} I(G_m(x_i) \neq y_i)$, w_{mi} —— 第 m 轮中第 i 个实例的权重, $\sum_{i=1}^N w_{mi} = 1$

$\alpha_m = \frac{1}{2} \log \frac{1 - e_m}{e_m} \rightarrow e_m \downarrow \Rightarrow \alpha_m \uparrow$, 分类误差率小的 $G_m(x)$ 作用大.

$w_{m+1,i} = \frac{w_{mi}}{Z_m} \exp(-\alpha_m y_i G_m(x_i))$, $Z_m = \sum_{i=1}^N w_{mi} \exp(-\alpha_m y_i G_m(x_i))$

$= \begin{cases} \frac{w_{mi}}{Z_m} e^{\alpha_m} & (\text{当 } G_m(x_i) = y_i \text{ 时}) \rightarrow \alpha_m \uparrow \Rightarrow w_{m+1,i} \downarrow, \text{被 } G_m(x) \text{ 正确分类的样本权重} \downarrow \\ \frac{w_{mi}}{Z_m} e^{-\alpha_m} & (\text{当 } G_m(x_i) \neq y_i \text{ 时}) \rightarrow \alpha_m \uparrow \Rightarrow w_{m+1,i} \uparrow, \text{被 } G_m(x) \text{ 错误分类的样本权重} \uparrow \end{cases}$