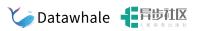


《机器学习公式详解》 (南瓜书)

第3章 二分类线性判别分析

本节主讲: 谢文睿

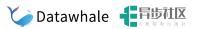
本节大纲



西瓜书对应章节: 3.4

- 1. 算法原理(模型)
- 2. 损失函数推导(策略)
- 3. 拉格朗日乘子法
- $4. 求解<math>\mathbf{w}$ (算法)
- 5. 广义特征值和广义瑞利商

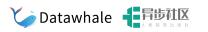
算法原理



从几何的角度,让全体训练样本经过投影后:

- 异类样本的中心尽可能远
- 同类样本的方差尽可能小

损失函数推导



经过投影后, 异类样本的中心尽可能远(并非严格投影):

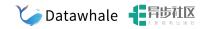
$$\max \| oldsymbol{w}^{ ext{T}} oldsymbol{\mu}_0 - oldsymbol{w}^{ ext{T}} oldsymbol{\mu}_1 \|_2^2$$

$$\max \||oldsymbol{w}| \cdot |oldsymbol{\mu}_0| \cdot \cos heta_0 - |oldsymbol{w}| \cdot |oldsymbol{\mu}_1| \cdot \cos heta_1\|_2^2$$

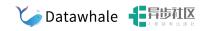
经过投影后,同类样本的方差尽可能小(并非严格方差):

$$\min oldsymbol{w}^{ ext{T}} oldsymbol{\Sigma}_0 oldsymbol{w}$$

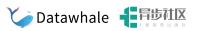
$$egin{aligned} oldsymbol{w}^{\mathrm{T}} oldsymbol{\Sigma}_0 oldsymbol{w} &= oldsymbol{w}^{\mathrm{T}} \left(\sum_{oldsymbol{x} \in X_0} (oldsymbol{x} - oldsymbol{\mu}_0) (oldsymbol{x} - oldsymbol{\mu}_0)^{\mathrm{T}} oldsymbol{w} \ &= \sum_{oldsymbol{x} \in X_0} (oldsymbol{w}^{\mathrm{T}} oldsymbol{x} - oldsymbol{w}^{\mathrm{T}} oldsymbol{\mu}_0) (oldsymbol{x}^{\mathrm{T}} oldsymbol{w} - oldsymbol{\mu}_0^{\mathrm{T}} oldsymbol{w}) \end{aligned}$$



$$egin{aligned} \max J &= rac{\|oldsymbol{w}^{\mathrm{T}}oldsymbol{\mu}_{0} - oldsymbol{w}^{\mathrm{T}}oldsymbol{\mu}_{1}\|_{2}^{2}}{oldsymbol{w}^{\mathrm{T}}oldsymbol{\Sigma}_{0}oldsymbol{w} + oldsymbol{w}^{\mathrm{T}}oldsymbol{\Sigma}_{1}oldsymbol{w}} \ &= rac{\|(oldsymbol{w}^{\mathrm{T}}oldsymbol{\mu}_{0} - oldsymbol{w}^{\mathrm{T}}oldsymbol{\Sigma}_{0} + oldsymbol{\Sigma}_{1})oldsymbol{w}}{oldsymbol{w}^{\mathrm{T}}oldsymbol{\Sigma}_{0} + oldsymbol{\Sigma}_{1})oldsymbol{w}} \ &= rac{\|(oldsymbol{\mu}_{0} - oldsymbol{\mu}_{1})^{\mathrm{T}}oldsymbol{w}\|_{2}^{2}}{oldsymbol{w}^{\mathrm{T}}oldsymbol{\Sigma}_{0} + oldsymbol{\Sigma}_{1})oldsymbol{w}} \ &= rac{\|(oldsymbol{\mu}_{0} - oldsymbol{\mu}_{1})^{\mathrm{T}}oldsymbol{w}}{oldsymbol{w}^{\mathrm{T}}oldsymbol{\Sigma}_{0} + oldsymbol{\Sigma}_{1})oldsymbol{w}} \ &= rac{oldsymbol{w}^{\mathrm{T}}oldsymbol{\mu}_{0} - oldsymbol{\mu}_{1})(oldsymbol{\mu}_{0} - oldsymbol{\mu}_{1})^{\mathrm{T}}oldsymbol{w}}{oldsymbol{w}^{\mathrm{T}}oldsymbol{\Sigma}_{0} + oldsymbol{\Sigma}_{1})oldsymbol{w}} \end{aligned}$$



拉格朗日乘子法



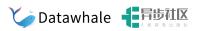
对于仅含等式约束的优化问题(解释下为啥是min):

$$egin{array}{ll} \min_{oldsymbol{x}} & f(oldsymbol{x}) \ \mathrm{s.t.} & h_i(oldsymbol{x}) = 0 & i = 1, 2, ..., n \end{array}$$

其中自变量 $\mathbf{x} \in \mathbb{R}^n$, $f(\mathbf{x})$ 和 $h_i(\mathbf{x})$ 均有连续的一阶偏导数。首先列出其拉格朗日函数:

$$L(oldsymbol{x},oldsymbol{\lambda}) = f(oldsymbol{x}) + \sum_{i=1}^n \lambda_i h_i(oldsymbol{x})$$

其中 $\lambda = (\lambda_1, \lambda_2, ..., \lambda_n)^{\mathrm{T}}$ 为拉格朗日乘子。然后对拉格朗日函数关于 \boldsymbol{x} 求偏导,并令导数等于 $\boldsymbol{0}$ 再搭配约束条件 $h_i(\boldsymbol{x}) = 0$ 解出 \boldsymbol{x} ,求解出的所有 \boldsymbol{x} 即为上述优化问题的所有可能【极值点】



$$egin{array}{ll} \min_{oldsymbol{w}} & -oldsymbol{w}^{\mathrm{T}} \mathbf{S}_b oldsymbol{w} \ \mathrm{s.t.} & oldsymbol{w}^{\mathrm{T}} \mathbf{S}_w oldsymbol{w} = 1 \Leftrightarrow oldsymbol{w}^{\mathrm{T}} \mathbf{S}_w oldsymbol{w} - 1 = 0 \end{array}$$

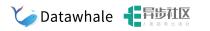
由拉格朗日乘子法可得拉格朗日函数为

$$L(oldsymbol{w},\lambda) = -oldsymbol{w}^{\mathrm{T}}\mathbf{S}_boldsymbol{w} + \lambda(oldsymbol{w}^{\mathrm{T}}\mathbf{S}_woldsymbol{w} - 1)$$

对w求偏导可得

$$egin{aligned} rac{\partial L(oldsymbol{w},\lambda)}{\partial oldsymbol{w}} &= -rac{\partial (oldsymbol{w}^{\mathrm{T}}\mathbf{S}_boldsymbol{w})}{\partial oldsymbol{w}} + \lambda rac{\partial (oldsymbol{w}^{\mathrm{T}}\mathbf{S}_woldsymbol{w} - 1)}{\partial oldsymbol{w}} \ &= -(\mathbf{S}_b + \mathbf{S}_b^{\mathrm{T}})oldsymbol{w} + \lambda (\mathbf{S}_w + \mathbf{S}_w^{\mathrm{T}})oldsymbol{w} \end{aligned}$$

由于
$$\mathbf{S}_b = \mathbf{S}_b^{\mathrm{T}}, \mathbf{S}_w = \mathbf{S}_w^{\mathrm{T}}$$
,所以



$$rac{\partial L(oldsymbol{w},\lambda)}{\partial oldsymbol{w}} = -2 \mathbf{S}_b oldsymbol{w} + 2 \lambda \mathbf{S}_w oldsymbol{w}$$

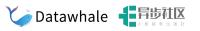
令上式等于0即可得

$$egin{align} -2\mathbf{S}_boldsymbol{w}+2\lambda\mathbf{S}_woldsymbol{w}&=0\ \mathbf{S}_boldsymbol{w}&=\lambda\mathbf{S}_woldsymbol{w}\ &(oldsymbol{\mu}_0-oldsymbol{\mu}_1)(oldsymbol{\mu}_0-oldsymbol{\mu}_1)^{\mathrm{T}}oldsymbol{w}&=\lambda\mathbf{S}_woldsymbol{w} \end{aligned}$$

若令
$$(\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)^{\mathrm{T}} \boldsymbol{w} = \gamma$$
,则

$$egin{aligned} oldsymbol{\gamma}(oldsymbol{\mu}_0 - oldsymbol{\mu}_1) &= \lambda \mathbf{S}_w oldsymbol{w} \ oldsymbol{w} &= rac{\gamma}{\lambda} \mathbf{S}_w^{-1} (oldsymbol{\mu}_0 - oldsymbol{\mu}_1) \end{aligned}$$

由于最终要求解的 \boldsymbol{w} 不关心其大小,只关心其方向,所以 $\frac{\gamma}{\lambda}$ 这个常数项可以任意取值, 西瓜书中所说的"不妨令 $\mathbf{S}_b \boldsymbol{w} = \lambda (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1)$ "就等价于令 $\gamma = \lambda$,进而使得 $\frac{\gamma}{\lambda} = 1$,此 时求解出的 \boldsymbol{w} 即为公式(3.39),另外,此处并未严格按照上述拉格朗日乘子法再刻意考 虑等式约束条件 $\boldsymbol{w}^{\mathrm{T}} \mathbf{S}_w \boldsymbol{w} - 1 = 0$ 也是因为不在乎 \boldsymbol{w} 的大小。



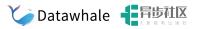
此时用拉格朗日乘子法求出来的极值点w一定是最小值点吗?

答:是的,因为 $-\mathbf{w}^{\mathrm{T}}\mathbf{S}_{b}\mathbf{w}=-\|\mathbf{w}^{\mathrm{T}}\boldsymbol{\mu}_{0}-\mathbf{w}^{\mathrm{T}}\boldsymbol{\mu}_{1}\|_{2}^{2}\leqslant0$,所以目标函数最大值为

0,且一定存在最小值,所以求出来的极值点只要代入目标函数不为0则一定是最小值

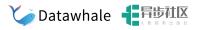
点。

广义特征值



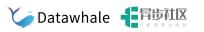
设 \mathbf{A} , \mathbf{B} 为n阶方阵,若存在数 λ ,使得方程 $\mathbf{A}x = \lambda \mathbf{B}x$ 存在非零解,则称 λ 为 \mathbf{A} 相对于 \mathbf{B} 的广义特征值,x为 \mathbf{A} 相对于 \mathbf{B} 的属于广义特征值 λ 的特征向量。特别地,当 $\mathbf{B} = \mathbf{I}$ (单位矩阵)时,广义特征值问题退化为标准特征值问题。

广义瑞利商



设 \mathbf{A} , \mathbf{B} 为n阶厄米(Hermitian)矩阵,且 \mathbf{B} 正定,称 $R(\boldsymbol{x}) = \frac{\boldsymbol{x}^{\mathrm{H}} \mathbf{A} \boldsymbol{x}}{\boldsymbol{x}^{\mathrm{H}} \mathbf{B} \boldsymbol{x}} (\boldsymbol{x} \neq \boldsymbol{0})$ 为 \mathbf{A} 相对于 \mathbf{B} 的广义瑞利商。特别地,当 $\mathbf{B} = \mathbf{I}$ (单位矩阵)时,广义瑞利商退化为瑞利商。

广义瑞利商



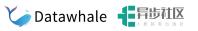
广义瑞利商的性质:设 $\lambda_i, \boldsymbol{x}_i (i=1,2,...,n)$ 为**A**相对于**B**的广义特征值和特征向量,且 $\lambda_1 \leqslant \lambda_2 \leqslant ... \leqslant \lambda_n$ 。

$$\min_{oldsymbol{x}
eq oldsymbol{0}} R(oldsymbol{x}) = rac{oldsymbol{x}^{ ext{H}} oldsymbol{A} oldsymbol{x}}{oldsymbol{x}^{ ext{H}} oldsymbol{B} oldsymbol{x}} = \lambda_1, oldsymbol{x}^* = oldsymbol{x}_1$$

$$\max_{oldsymbol{x}
eq oldsymbol{0}} R(oldsymbol{x}) = rac{oldsymbol{x}^{ ext{H}} oldsymbol{A} oldsymbol{x}}{oldsymbol{x}^{ ext{H}} oldsymbol{B} oldsymbol{x}} = \lambda_n, oldsymbol{x}^* = oldsymbol{x}_n$$

【证明】:当固定 $\mathbf{x}^{\mathrm{H}}\mathbf{B}\mathbf{x}=1$ 时,使用拉格朗日乘子法可推得 $\mathbf{A}\mathbf{x}=\lambda\mathbf{B}\mathbf{x}$ 这样一个广义特征值问题,因此 \mathbf{x} 所有可能的解即为 $\mathbf{x}_i(i=1,2,...,n)$ 这n个广义特征向量,将其分别代入 $R(\mathbf{x})$ 即可推得上述结论。

预告



下一节: 决策树

西瓜书对应章节: 4.1、4.2

结束语



欢迎加入【南瓜书读者交流群】,我们将在群里进行答疑、勘误、本次直播回放、本次直播PPT发放、下次直播通知等最新资源发放和活动通知。加入步骤:

- 1. 关注公众号【Datawhale】,发送【南瓜书】三个字获取机器人"小豚"的微信二维码
- 2. 添加"小豚"为微信好友,然后对"小豚"发送【南瓜书】三个字即可自动邀请进群、

