

越大 RMSE 值越高.随着 k 越大,模型计算成本越高,所以,选择基于 SVD++ 的思想来构建 Skill-LFM 模型.

实验表明,取 $\lambda_1 = 0.002$ 、隐因子 $k = 3$ 、迭代次数 $\text{steps} = 100$ 时,如图 2 所示,技能水平对知识点偏置 $s_{i,k}$ 的正则化参数 λ_2 在 0.005 时表现最好.

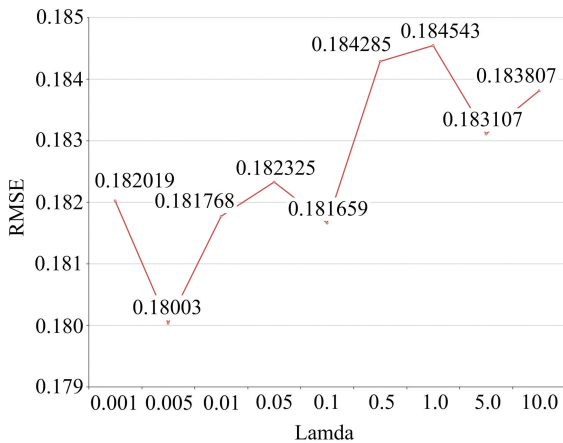


图 2 Skill-LFM 模型 λ_2 对 RMSE 的影响

Fig.2 Impact of λ_2 on RMSE of Skill-LFM

4 结论

本文提出融合用户技能水平的隐语义模型用于知识库系统中知识点的推荐,和传统 ItemCF、NMF 以及基础隐语义模型相比,在测试集上 RMSE 值最低.

应用上,基于技能的隐语义模型可在知识推荐领域推广应用,如呼叫中心融合 CSR 技能的知识推荐、在线课堂融合学生技能的习题推荐等.

方法上,在知识领域除融合用户技能水平这一属性,还可加入更多用户信息和知识点信息来解决冷启动,并提升推荐多样性和准确性.上下文感知技术是知识点推荐方法研究的重点,也是解决新知识点冷启动问题的一个方向.

性能上,SGD 单纯内存训练比较耗时,考虑并行 SGD 技术来支持大规模数据集的计算.

参考文献 (References)

[1] ULLMAN J D. Principles of Database and Knowledge-Base Systems [M]. New York: Computer Science Press, 1988.

[2] ALAVI M, LEIDNER D E. Knowledge management and knowledge management systems: Conceptual foundations and research issues [J]. MIS Quarterly, 2001, 25(1): 107-136.

[3] GIBONEY J S, BROWN S A, LOWRY P B, et al. User acceptance of knowledge-based system recommendations: explanations, arguments, and fit [J]. Decision Support Systems, 2015, 72(C):1-10.

[4] VELÁSQUEZ J D, PALADE V. Building a knowledge base for implementing a web-based computerized recommendation system [J]. International Journal on Artificial Intelligence Tools, 2007, 16(5): 793-828.

[5] KOREN Y, BELL R, VOLINSKY C. Matrix factorization techniques for recommender systems [J]. Computer, 2009, 42(8): 30-37.

[6] SARWAR B, KARYPIS G, KONSTAN J, et al. Application of dimensionality reduction in recommender systems [EB/OL]. ACM WebKDD-2000, [2018-11-17] <http://glaros.dtc.umn.edu/gkhome/node/122>.

[7] SARWAR B, KARYPIS G, KONSTAN J, et al. Item-based collaborative filtering recommendation algorithms [C]// International Conference on World Wide Web. Hong Kong, China: ACM Press, 2001: 285-295.

[8] LINDEN G, SMITH B, YORK J. Amazon.com recommendations: item-to-item collaborative filtering [J]. IEEE Internet Computing, 2003, 7(1):76-80.

[9] ZENG C, XING C X, ZHOU L Z. Survey of personalization technology [J]. Journal of Software, 2002, 13(10): 1952-1961.

[10] RICCI F, ROKACH L, SHAPIRA B, et al. Recommender Systems Handbook [M]. Springer, 2011.

[11] ADOMAVICIUS G, TUZHILIN A. Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions [J]. IEEE Transactions on Knowledge and Data Engineering, 2005, 17(6): 734-749.

[12] BENNETT J, LANNING S, NETFLIX N. The Netflix Prize [C]// KDD Cup and Workshop in Conjunction with KDD, 2009.

[13] MASSA P, AVESANI P. Trust-aware recommender systems [C]// Proceedings of the Conference on Recommender Systems. Minnesota, USA: ACM Press, 2007: 17-24.

[14] GOLDBERG D, NICHOLS D, OKI B M, et al. Using collaborative filtering to weave an information tapestry [J]. Communications of the ACM, 1992, 35(12): 61-70.

[15] BURKE R. Hybrid recommender systems: Survey and experiments [J]. User Modeling and User-Adapted Interaction, 2002, 12(4): 331-370.

[16] KOREN Y. Factorization meets the neighborhood: A multifaceted collaborative filtering model [C]// ACM

- SIGKDD International Conference on Knowledge Discovery and Data Mining. Las Vegas, USA: ACM Press, 2008: 426-434.
- [17] KOREN Y. Collaborative filtering with temporal dynamics[J]. *Communications of the ACM*, 2010, 53(4): 89-97.
- [18] JAMALI M, ESTER M. A matrix factorization technique with trust propagation for recommendation in social networks [C]// *ACM Conference on Recommender Systems*. Barcelona, Spain: ACM Press, 2010: 135-142.
- [19] CHEN T, ZHANG W, LU Q, et al. SVDFeature: A toolkit for feature-based collaborative filtering [J]. *Journal of Machine Learning Research*, 2012, 13(1): 3619-3622.
- [20] OSMANLI O N, SMOYAL HAKKI TOROSLU. Using tag similarity in SVD-based recommendation systems[C]// *International Conference on Application of Information and Communication Technologies*. Baku, Azerbaijan: IEEE Press, 2011: 1-4.
- [21] XU Z, CHANG X, XU F, et al. $L_{1/2}$ regularization: A thresholding representation theory and a fast solver[J]. *IEEE Transactions on Neural Networks & Learning Systems*, 2012, 23(7): 1013-1027.
- [22] RUDER S. An overview of gradient descent optimization algorithms[J]. *Machine Learning*, 2016: arXiv:1609.04747 [cs.LG].
- [23] SALAKHUTDINOV R, MNIH A. Probabilistic matrix factorization[C]// *International Conference on Neural Information Processing Systems*. Vancouver, Canada: Curran Associates Inc., 2007: 1257-1264.
- [24] LEE D D, SEUNG H S. Algorithms for non-negative matrix factorization[C]// *International Conference on Neural Information Processing Systems*. MIT Press, 2000: 535-541.
- [25] SHARIFI Z, REZGHI M, NASIRI M. A new algorithm for solving data sparsity problem based on Non negative matrix factorization in recommender systems[C]// *International Conference on Computer and Knowledge Engineering*. Mashhad, Iran: IEEE Press, 2014: 56-61.
- [26] LUO X, ZHOU M, XIA Y, et al. An efficient non-negative matrix-factorization-based approach to collaborative filtering for recommender systems [J]. *IEEE Transactions on Industrial Informatics*, 2014, 10(2): 1273-1284.
- [27] WANG L, MENG X, ZHANG Y, et al. Applying HOSVD to alleviate the sparsity problem in context-aware recommender systems [J]. *Chinese Journal of Electronics*, 2013, 22(4): 773-778.
- [28] KUTTY S, CHEN L, NAYAK R. A people-to-people recommendation system using tensor space models [C]// *ACM Symposium on Applied Computing*. Trento, Italy: ACM Press, 2012: 187-192.
- [29] KARATZOGLOU A, AMATRIAIN X, BALTRUNAS L, et al. Multiverse recommendation: n -dimensional tensor factorization for context-aware collaborative filtering [C]// *ACM Conference on Recommender Systems*. Barcelona, Spain: ACM Press, 2010: 79-86.
- [30] CREMONESI P, TURRIN R, TURRIN R. Performance of recommender algorithms on top- n recommendation tasks [C]// *ACM Conference on Recommender Systems*. Barcelona, Spain: ACM Press, 2010: 39-46.
- [31] WILLMOTT C J, MATSUURA K. Advantages of the mean absolute error (MAE) over the root mean square error (RMSE) in assessing average model performance [J]. *Climate Research*, 2005, 30(1): 79-82.
- [32] CHAI T, DRAXLER R R. Root mean square error (RMSE) or mean absolute error (MAE)? [J]. *Geoscientific Model Development*, 2014, 7(3): 1247-1250.
- [33] TAK, CS G, SZY I, et al. Matrix factorization and neighbor based algorithms for the netflix prize problem [C]// *ACM Conference on Recommender Systems*. Lausanne, Switzerland: ACM Press, 2008: 267-274.
- [34] OTT P. Incremental Matrix Factorization for Collaborative Filtering[M]// *Science, Technology and Design 01/2008*, Anhalt University of Applied Sciences.
- [35] LEE D D, SEUNG H S. Learning the parts of objects by non-negative matrix factorization[J]. *Nature*, 1999, 401(6755): 788791.

基于封装式偏差回归的心脏生理和病理年龄估计算法

李勇明^{1,2}, 肖洁¹, 王品¹, 颜芳¹

(1. 重庆大学通信工程学院, 重庆 400044; 2. 重庆大学脑科学协同创新中心, 重庆 400044)

摘要: 研究表明, 年龄与心脏病具有紧密联系, 心脏年龄对于心脏健康状态检测和监测都有着极为重要的作用, 为此提出了心脏生理年龄和病理年龄估计算法, 前者仅基于正常人样本建立回归模型; 后者基于所有类别样本, 通过最小均方误差建立回归模型, 其引入年龄偏差表征病理年龄和实际年龄差异, 通过最大化分类准确率来搜索最优年龄偏差, 从而实现病理年龄估计, 文中采用了公共数据集进行验证, 实验结果表明, 两种估计年龄的分类能力和信息表征能力均优于实际年龄, 显著性差异远小于 0.01. 本文首次提出了心脏病理年龄这一新概念, 有助于为心脏健康状态检测和监测提供有效的标记物.

关键词: 心脏病; 诊断; 心脏生理年龄; 心脏病理年龄; 机器学习

中图分类号: TP391 **文献标识码:** A doi: 10.3969/j.issn.0253-2778.2018.09.011

引用格式: 李勇明, 肖洁, 王品, 等. 基于封装式偏差回归的心脏生理和病理年龄估计算法[J]. 中国科学技术大学学报, 2018, 48(9): 762-769.

LI Yongming, XIAO Jie, WANG Pin, et al. Heart physiological and pathological age estimation based on wrapper deviation regression[J]. Journal of University of Science and Technology of China, 2018, 48(9): 762-769.

Heart physiological and pathological age estimation based on wrapper deviation regression

LI Yongming^{1,2}, XIAO Jie¹, WANG Pin¹, YAN Fang¹

(1. *Communication engineering of Chongqing University, Chongqing 400044, China;*

2. *Collaborative Innovation Center for Brain Science, Chongqing University, Chongqing 400044, China*)

Abstract: Researches show that a person age is highly related to his heart. Heart age is very important for examining and monitoring of the heart's state. Two algorithms for estimating the physiological and pathological age of the heart were proposed based on data mining technique. The first algorithm is based on a regression model for healthy people by using the mean absolute error (MAE), while the latter is based on a regression model for all types of people by considering the age deviation. The optimal age deviation is searched within the range of deviation candidates and is obtained by maximizing the classification accuracy. Based on the optimal age deviation and real age, the heart pathological age is obtained. The public heart dataset is used for verification of the proposed algorithm. Experimental results show that two estimated heart ages are better than the real age, with the apparent significance level the lower than 0.01. Compared with the current heart age estimation algorithm, the heart pathological age estimation algorithm can lead

收稿日期: 2018-05-28; **修回日期:** 2018-09-18

作者简介: 李勇明(通讯作者), 男, 1976年生, 博士/教授. 研究方向: 大数据分析、生物医学信号与信息处理等. E-mail: yongmingli@cqu.edu.cn

to the better classification capability and is more helpful with improving the classification accuracy of heart disease as a marker or feature. Besides, a new concept——heart pathological age is proposed for the first time, and which may help provide an effective marker for monitoring and supervising heart health.

Key words: heart disease; diagnosis; heart physiological age; heart pathological age; machine learning

0 引言

心脏病是一种重要病种,危害严重^[1].许多学者已经发现年龄与心脏病之间存在较紧密的关系^[2-6].Davis 发现非洲裔美国老年妇女更有可能患有心血管疾病(cardiovascular disease, CVD)^[2].Finegold 等发现年龄与缺血性心脏病(ischemic heart disease, IHD)有关,人口老龄化使 IHD 成为主要死亡原因^[3].Koopman 等评估了 205 649 个慢性冠心病(coronary heart disease, CHD)病人的住院率,住院年龄被推向老年,早发冠心病住院率随时间推移呈下降趋势,在老年时住院率增加^[4].Padur 等对 170 例心脏科患者的病历进行回顾性分析,并观察到 51-60 岁年龄组的患者更多^[5].Chang 等发现 IHD 死亡率的年龄效应呈对数线性增长,其中 80 岁 84 岁的 IHD 死亡率分别为城乡 20 岁以上人群的 277 和 161 倍,其结果表明,人口老龄化是 IHD 死亡率迅速上升背后的主要因素^[6].以上这些结果都表明了年龄与心脏病具有明确的相关关系.

进一步研究发现,实际年龄与心脏病的关系并不足以用于心脏病诊断^[7-9].Wu 等发现患心脏病的危险与年龄之间没有显著性关系^[7].Frankel 等发现目前还不足以解释包括年龄在内的传统危险因素与成年期 CHD 的关系^[8].Mendes 等发现有一些研究结果支持欧盟成员国青年群体的 CHD 死亡率可能比老年群体更稳定的假设^[9].以上结果说明,实际年龄不能充分有效地反应心脏老化状态和过程.因此,有必要从心脏医疗数据中定量估计心脏内在年龄(心脏年龄)来取代实际年龄,从而更准确的诊断和监测心脏状态、老化及发病程度.

目前,几乎无公开文献研究心脏年龄的定量估计算法.现有的评分法、问卷法等具有明显的主观性,难以客观、量化和推广.迄今为止,在其他一些领域中,学者们已经展开了年龄估计的定量研究,取得了显著成效.这包括:面部年龄估计^[10-11]、牙齿年龄估计^[12]、新生儿或胎龄估计^[13]、骨龄估计^[14-15]、脑年龄估计^[16-18]和创伤年龄估计^[19]等.这些结果都表明,年龄定量估计是可行和有效的,其中机器学习是

挖掘年龄信息的有效手段.

上述大多数研究都没有涉及如何将年龄定量估计用于疾病诊断.近年来,一些学者开始着手相关研究.Frankel 等讨论了如何估计用于阿尔茨海默病分类的脑年龄^[20-21].Löwe 等讨论了如何估计用于 2 型糖尿病分类的脑年龄^[22].Dias 等讨论了如何估计用于牙周病分类的牙齿年龄^[23].Moyses 等讨论了如何估计用于阿尔茨海默病分类的面部年龄^[24].上述结果表明,年龄定量估计有助于疾病分类,定量估计的年龄能成为有效的诊断标记物.

目前尚未发现利用机器学习方法估计心脏年龄并用于心脏病分类诊断的工作.基于以上分析,本文基于回归方法提出了两种用于心脏病分类的心脏年龄估计算法.第一种是心脏生理年龄估计(HeartAge_Estima)算法,其思路是基于正常人样本建立回归模型,训练标签为实际年龄,通过最小化估计年龄和实际年龄差异进行模型训练,其误差函数为估计年龄与实际年龄之间的差异,即平均绝对误差(mean absolute error, MAE).用 HeartAge_Estima 算法估计的心脏年龄称为心脏生理年龄.第二种是心脏病理年龄估计(Path_HeartAge_Estima)算法,其思路是基于所有类别样本建立回归模型,引入年龄偏差表征病理年龄和实际年龄差异,训练标签不是实际年龄,而是实际年龄加偏差,通过最大化分类准确率来搜索最优年龄偏差.由于偏差与心脏病不同状态相对应,因此这些偏差可以定量估计心脏衰老或病变程度,从而有助于心脏病分类诊断.用 Path_HeartAge_Estima 算法估计的年龄称为心脏病理年龄.

1 基于封装式偏差回归的心脏年龄估计算法

1.1 心脏生理年龄估计算法

回归模型的适应度函数是算法的核心,本文采用 MAE 表示估计年龄与实际年龄的误差,MAE 定义为

$$MAE = \frac{1}{N} \sum_{j=1}^N |y_j - \hat{y}_j| \quad (1)$$

式中, N 为样本数目, y_j 为实际年龄, \hat{y}_j 为估计年龄, 则 HeartAge_Estima 算法的适应度函数表示为

$$F_1 = \arg \left[\min \left(\frac{1}{N} \sum_{j=1}^N |y_j - \hat{y}_j| \right) \right] \quad (2)$$

1.2 心脏病理年龄估计算法

HeartAge_Estima 算法只考虑了正常人, 即在训练模型时没有考虑正常人与心脏病人的差异. 为此, 本文进一步提出了 Path_HeartAge_Estima 算法, 其训练模型的适应度函数基于反映估计年龄分类能力的分类评价准则设计. 本文同时介绍了两种流行的分类评价准则——可分度距离准则和相关性准则, 其值与分类能力或回归能力成正比^[25-27]. Path_HeartAge_Estima 算法的适应度函数表示为

$$F_{21} |_{de^1, de^2, \dots, de^c} = \arg \left[\max \left(\frac{\sum_{i=1}^c P_i \overline{(y_i(de^1, de^2, \dots, de^c) - y(de^1, de^2, \dots, de^c))}^2}{\sum_{i=1}^c P_i \frac{1}{N_i} \sum_{k=1}^{N_i} (y_{ik}(de^1, de^2, \dots, de^c) - y_i(de^1, de^2, \dots, de^c))^2} \right) \right]_{\min(\|\hat{y}_j - (y_j + de^q)\|, \alpha)} \quad (5)$$

式中, P_i 表示第 i 类样本占总样本的比例, N_i 表示第 i 类样本的数量, $y_{ik}(de^1, de^2, \dots, de^c)$ 表示第 i 类第 k 个样本的年龄估计值, $\overline{y_i(de^1, de^2, \dots, de^c)}$ 表示第 i 类样本年龄估计值的均值, $\overline{y(de^1, de^2, \dots, de^c)}$ 表示所有样本年龄估计值的均值.

对于相关性准则, 训练模型的适应度函数表示为

$$F_{22} = \arg \left[\max(\text{corr}(\hat{y}_j, y_j^t)) \right]_{\min(\|\hat{y}_j - (y_j + de^q)\|, \alpha)} \quad (6)$$

式中, \hat{y}_j 为回归模型估计年龄, y_j^t 为样本类别标签, α 表示范数类型, 一般取值为 1 或者 2, 本文取值为 1 表示 1 范数. 本文采用皮尔逊相关系数表示相关性准则, 则式(6)可表示为

$$F_{22} |_{de^1, de^2, \dots, de^c} = \arg \left[\max \left(\frac{c_{el}}{\sqrt{c_{ee} c_{ll}}} \right) \right]_{\min(\|\hat{y}_j - (y_j + de^q)\|, \alpha)} \quad (7)$$

式中, $c_{el} = \frac{1}{N-1} \sum_{j=1}^N (y_j^e - \overline{y^e})(y_j^l - \overline{y^l})$, N 为样本总数, y_j^l 表示第 j 个样本类别标签, $\overline{y^l}$ 表示类别标签均值, y_j^e 表示第 j 个样本年龄估计值, $\overline{y^e}$ 表示

$$F_2 = \arg \left[\max(\text{eval}_{\text{ACC}}) \right]_{\min(\|\hat{y}_j - (y_j + de^q)\|, \alpha)} \quad (3)$$

式中, eval_{ACC} 是基于分类评价准则的间接分类准确率, de 为实际年龄与病理年龄之间的偏差, q 表示第 q 类样本. de^q 与心脏病不同状态对应, 因此基于式(3)估计的心脏年龄更有利于提高心脏病的分类准确率. 对于可分度距离准则, 训练模型的适应度函数表示为

$$F_{21} = \arg \left[\max \left(\frac{S_b}{S_w} \right) \right]_{\min(\|\hat{y}_j - (y_j + de^q)\|, \alpha)} \quad (4)$$

式中, S_b 为类间距离, S_w 为内类距离, α 表示范数类型, 一般取值为 1 或者 2, 本文取值为 1 表示 1 范数. S_b 和 S_w 是估计年龄 \hat{y}_j 的函数, 假设数据类别数目为 c , 第 q 类的偏差是 de^q , 则式(4)可以表示为

年龄估计值的均值.

Path_HeartAge_Estima 算法流程如图 1 所示. 本文以二类分类为例讨论, de^1 表示类别 1 的偏差, 范围为 $[de_{\min}^1, de_{\max}^1]$, de^2 表示类别 2 的偏差, 范围为 $[de_{\min}^2, de_{\max}^2]$. 假设 $A_{\text{ge_class1}}$ 表示类别为 1 的样本的实际年龄, $A_{\text{ge_class2}}$ 表示类别为 2 的样本的实际年龄, 则回归模型的训练标签分别为 $A_{\text{ge_class1}} + de^1$ 与 $A_{\text{ge_class2}} + de^2$ 而不是 $A_{\text{ge_class1}}$ 与 $A_{\text{ge_class2}}$.

Path_HeartAge_Estima 算法伪代码如下:

输入: 训练样本特征矩阵 T_r 、验证样本特征矩阵 V 、测试样本特征矩阵 T_e ;

初始化: 当前偏差组合 (de^1, de^2) ;

过程:

(1) 针对当前偏差组合 (de^1, de^2) , 修改 T_r 的年龄标签, 得到新的训练样本特征矩阵 T_{r1} ;

(2) 将 T_{r1} 送入 SVR 进行训练, 获得训练后的 SVR 模型, 保存 SVR 模型参数;

(3) 基于训练后的 SVR 模型和验证样本特征矩阵 V , 输出验证样本的年龄估计值, 从而计算适应度值 $F_2 |_{de^1, de^2}$, 保存该适应度值及对应的偏差组合 (de^1, de^2) ;

(4) 比较偏差组合 (de^1, de^2) 和 $[de_{\min}^1, de_{\max}^1]$ 及 $[de_{\min}^2, de_{\max}^2]$, 判断遍历是否结束; 若是, 则进入第(5)步; 否则, 则返回第(1)步;

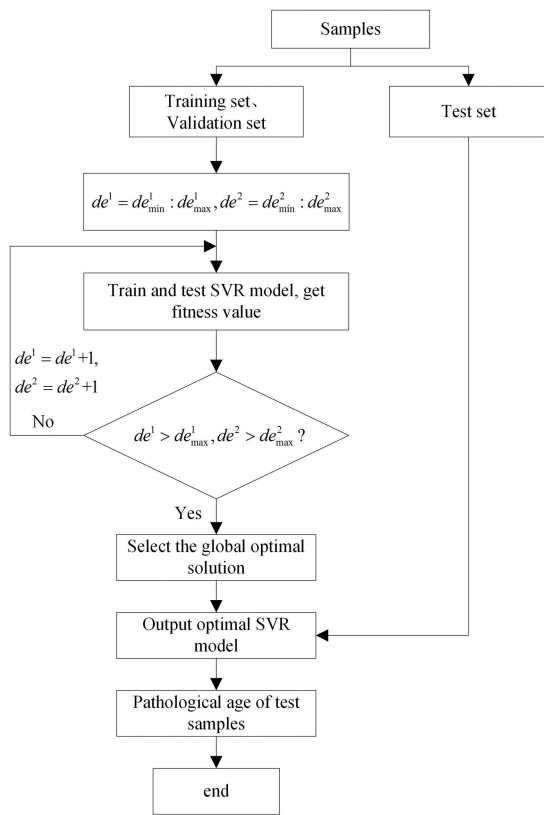


图 1 Path_HeartAge_Estima 算法流程图

Fig.1 The flowchart of Path_HeartAge_Estima algorithm

(5) 比较保存的 $F_2 |_{de^1, de^2}$, 选出最大值 F_{2_max} , 及其对应的最优偏差组合 (de^1_{opt}, de^2_{opt}) 和最优 SVR 模型 SVR_{opt} ;

(6) 将测试样本特征矩阵 T_e 输入到 SVR_{opt} 得到测试样本心脏病理年龄;

输出: 最优偏差组合 (de^1_{opt}, de^2_{opt}) 和测试样本心脏病理年龄.

2 实验结果与分析

2.1 实验条件

为了验证两种心脏年龄估计算法的有效性, 本文选择了心脏病公共数据集 (<http://archive.ics.uci.edu/ml/index.php>) 进行实验. 样本大致分为两类, 分别为正常对照组 (normal control, NC) 与心脏病病人 (heart disease, HD). 通过 hold out 交叉验证法获得 30 组数据, 每组数据由训练集、验证集与测试集 3 部分组成. 本文实验基于 Win10 64 位操作系统平台, 通过 Matlab 2016b 完成数据处理. 回归模型采用支持向量回归机 (support vector regression, SVR), 其核函数采用线性核函数, 参数设置为默认值^[28]. HeartAge_Estima 算法与 Path_HeartAge_Estima 算法是一种框架算法而非具体算

法, 通过采用不同的评价准则或回归模型, 可以产生一系列年龄估计算法. 本文中, 这两种算法在同一条件下进行比较. 基于每种算法, 实验进行 30 次, 后面讨论的所有结果都是 30 次实验结果的统计平均结果.

2.2 心脏生理年龄估计结果

这部分实验实现了基于线性核的 SVR 回归模型. 表 1 比较了无特征压缩 (no feature selection, No FS), P_value (pearson correlation coefficient) 特征选择, PCA (principle component analysis) 特征压缩三种情况下估计年龄的 MAE.

表 1 基于 HeartAge_Estima 算法估计年龄的 MAE

Tab.1 MAE of estimated heart age by HeartAge_Estima algorithm

MAE	No FS	P_value	PCA
Mean	6.173 5	6.745 7	5.941 2
Std	0.488 2	0.583 5	0.486 5

表 1 结果表明, 估计年龄在最好的情况下, MAE 可达 5.94. 在采用特征压缩算法的情况下, 年龄检测更接近于实际年龄, 且 PCA 比 P_value 效果更好.

2.3 心脏病理年龄估计结果

根据初步经验, 年龄偏差 (de^1, de^2) 应该不会太大, 两类样本年龄偏差范围 $[de^1_{min}, de^1_{max}]$ 与 $[de^2_{min}, de^2_{max}]$ 设置为 $[-10, 10]$ 比较合适. 为了系统地比较 Path_HeartAge_Estima 算法与 HeartAge_Estima 算法的差异, 本文通过不同评价准则 (可分度距离准则, 相关性准则)、不同特征压缩算法 (No FS, P_value, PCA) 进行了实验. 对于 HeartAge_Estima 算法, 在无特征压缩情况下进行实验并计算 MAE. 年龄偏差 (de^1, de^2) 估计实验结果如表 2 所示. 为方便描述, P_value 与 PCA 统称为特征选择 (feature selection, FS).

对于 NC, 心脏病理年龄与实际年龄之间的差异为 de^1 ; 对于 HD, 心脏病理年龄与实际年龄之间的差异为 de^2 . 表 2 结果表明, 年龄偏差 de^1 总是小于 de^2 , 并且在所有情况下 de^1 与 de^2 之间的差异都很大. 例如, 在无特征压缩和可分度距离准则情况下, de^1 为 -7.7, de^2 为 7.667. 除此之外, de^1 总是比 0 小, de^2 总是比 0 大. 通过 p 值可以看出估计的病理年龄与实际年龄之间具有统计学上的显著性差

异 ($p < 0.01$), 因此本文提出的 Path_HeartAge_Estima 算法在统计学意义上来讲是可靠的. 以上结论表明, 正常人的病理年龄总是低于其实际年龄, 而心脏病人的病理年龄总是高于实际年龄, 表明病理年龄具有较好的分类能力. 由表 2 可知, 在基于 P_value 与 PCA 特征压缩的情况下, 依然具有以上相似结论, 从 de^1 与 de^2 的差异看, P_value 略高于 PCA. 从表 2 还可以看出, 基于 HeartAge_Estima 算

法估计生理年龄, NC 的 MAE 为 0.375, HD 的 MAE 为 4.038, 且与实际年龄均具有统计学上的显著性差异. 此外, 通过差异来区分不同类别的样本, 病理年龄优于生理年龄. 在无特征压缩的情况下, 前者是 3.663(4.038-0.375), 后者是 15.367(7.667-(-7.7)) 或者 15.366(7.733-(-7.633)), 基于 P_value 与 PCA 的特征压缩的情况下的结果与上述结论基本一致.

表 2 年龄偏差 (de^1, de^2) 估计结果

Tab.2 Estimation results of age deviations

experiment methods		NC_HD		
		distance	correlation	
No FS	without age detection	(0,0)	(0,0)	
	HeartAge_Estima	(0.375,4.038)	(0.375,4.038)	
	significant difference (p value)	(<0.001, <0.001)	(<0.001, <0.001)	
	Path_HeartAge_Estima	(-7.7,7.667)	(-7.633,7.733)	
	significant difference (p value)	(<0.001, <0.001)	(<0.001, <0.001)	
P_value	without age detection	(0,0)	(0,0)	
	Path_HeartAge_Estima	(-7.667, 8.033)	(-7.833,7.867)	
	significant difference (p value)	(<0.001, <0.001)	(<0.001, <0.001)	
FS	without age detection	(0,0)	(0,0)	
	PCA	Path_HeartAge_Estima	(-7.133,7)	(-7.033,7.1)
	significant difference (p value)	(<0.001, <0.001)	(<0.001, <0.001)	

2.4 心脏年龄可分度分析

该部分比较了仅考虑实际年龄, HeartAge_Estima 算法与 Path_HeartAge_Estima 算法 3 种情况下的年龄可分度. 3 种特征压缩算法 (No FS、P_value 特征选择、PCA 特征压缩) 与两种评价准则 (可分度距离准则, 相关准则) 在不同组合情况下进行了分组实验. 基于每种特征压缩算法, 根据估计年龄计算可分度距离值和相关性值. 在不同组合情况下得到的可分度距离值和相关性值的结果如表 3 所示.

从表 3 可以看出, 在不同组合情况下, 相比于仅考虑实际年龄, HeartAge_Estima 算法与 Path_HeartAge_Estima 算法具有更好的可分度距离值和相关性值. 因此仅考虑实际年龄是不够的, 心脏年龄

估计非常有必要. 此外, Path_HeartAge_Estima 算法具有比 HeartAge_Estima 算法更好的可分度距离值和相关性值 (见粗体), 特别是在基于 P_value 特征选择和可分度距离准则情况下尤为明显. Path_HeartAge_Estima 算法与 HeartAge_Estima 算法之间, 在基于 P_value 特征选择和可分度距离准则的情况下, 均值差距最大, 分别为 0.738 6 和 0.407 2; 在基于 PCA 特征压缩和可分度距离准则的情况下, 均值分别为 0.648 4 和 0.381 9. 由于可分度距离值或相关性值可以间接表征分类能力, 因此心脏生理年龄和病理年龄相比于实际年龄, HD 的分类能力得到了显著提高. 表 2 和表 3 均表明基于 P_value 的特征选择和可分度距离准则组合的 Path_HeartAge_Estima 算法对 HD 的分类最有帮助.

表 3 不同组合情况下适应度函数值
Tab.3 Fitness values under different conditions

experiment methods		NC_HD					
		without age detection		HeartAge_Estima		Path_HeartAge_Estima	
		Mean	Std	Mean	Std	Mean	Std
No FS	distance	0.080 4	0.029 0	0.367 6	0.134 4	0.660 4	0.196 1
	correlation	0.267 5	0.048 3	0.507 5	0.069 7	0.621 2	0.060 9
FS	P_value	0.080 4	0.029 0	0.407 2	0.212 0	0.738 6	0.223 2
	correlation	0.267 5	0.048 3	0.516 4	0.094 5	0.642 6	0.056 8
PCA	distance	0.080 4	0.029 0	0.381 9	0.160 5	0.648 4	0.168 3
	correlation	0.267 5	0.048 3	0.510 9	0.081 2	0.620 0	0.053 0

图 2 为不同组合情况下估计年龄的箱型图,不同行表示不同特征压缩算法下的年龄估计结果;不同列表示不同年龄估计算法,分别为不考虑年龄检测、HeartAge_Estima 算法与 Path_HeartAge_

Estima 算法.由于可分度距离准则与相关性准则具有一致的结论,因此只画出了可分度距离准则情况下的估计年龄的箱型图.

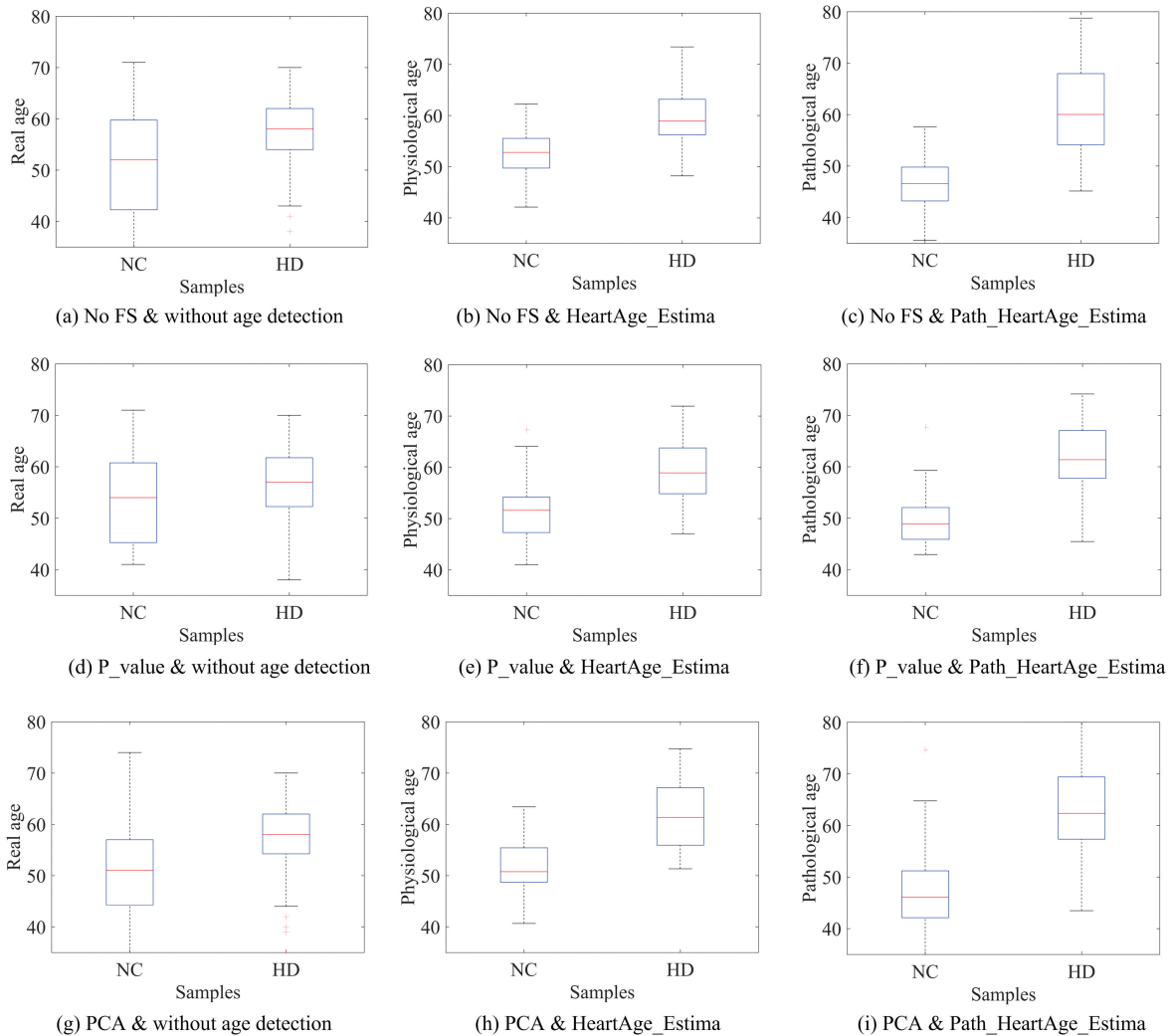


图 2 估计年龄分类能力图形比较

Fig.2 Comparisons of estimated heart ages