

Context-aware and Time-aware Attention-based Model for Disease Risk Prediction with Interpretability



合肥工業大學

柯水洲

2022.11.27

— 研究背景

An Interpretable Fast Model for Predicting The Risk of Heart Failure

Xianli Zhang * Buyue Qian † Xiaoyu Li* Jishang Wei‡ Yingqian Zheng *
Lingyun Song † Qinghua Zheng †

2019ISBM(CCF-B)

Context-aware and Time-aware Attention-based Model for Disease Risk Prediction with Interpretability

增加了一个attention模块
数据集、对比试验
可解释性方面的验证

2022TKDE (early access)

Xianli Zhang, Buyue Qian, Yang Li, Shilei Cao, and Ian Davidson,

使用电子健康记录Electronic Health Records (EHRs)数据实现疾病风险预测
disease risk prediction (DRP)。

局限性: Accuracy、Efficiency、Interpretability

方法: 双路输入 (文+时间) ; 线性变化 (NIPS2016)

二 研究方法

数据
介绍

Heart Failure、Diabetes、
Chronic Kidney Disease

TABLE I
STATISTICS OF THREE COHORTS FOR RISK PREDICTION

Datasets	HF	DIA	CKD
# of cases	808	1, 095	2, 206
# of controls	6, 464	3, 285	6, 618
# of visits	102, 237	68, 946	147, 892
Avg.# of visits per patient	14.05	15.74	16.76
# of unique ICD-9 codes	4518	3944	4892
Avg.# of codes per visit	2.31	2.38	2.38
Avg.# of codes per patient	32.42	37.51	39.94

心力衰竭、糖尿病和慢性肾脏病;对照组

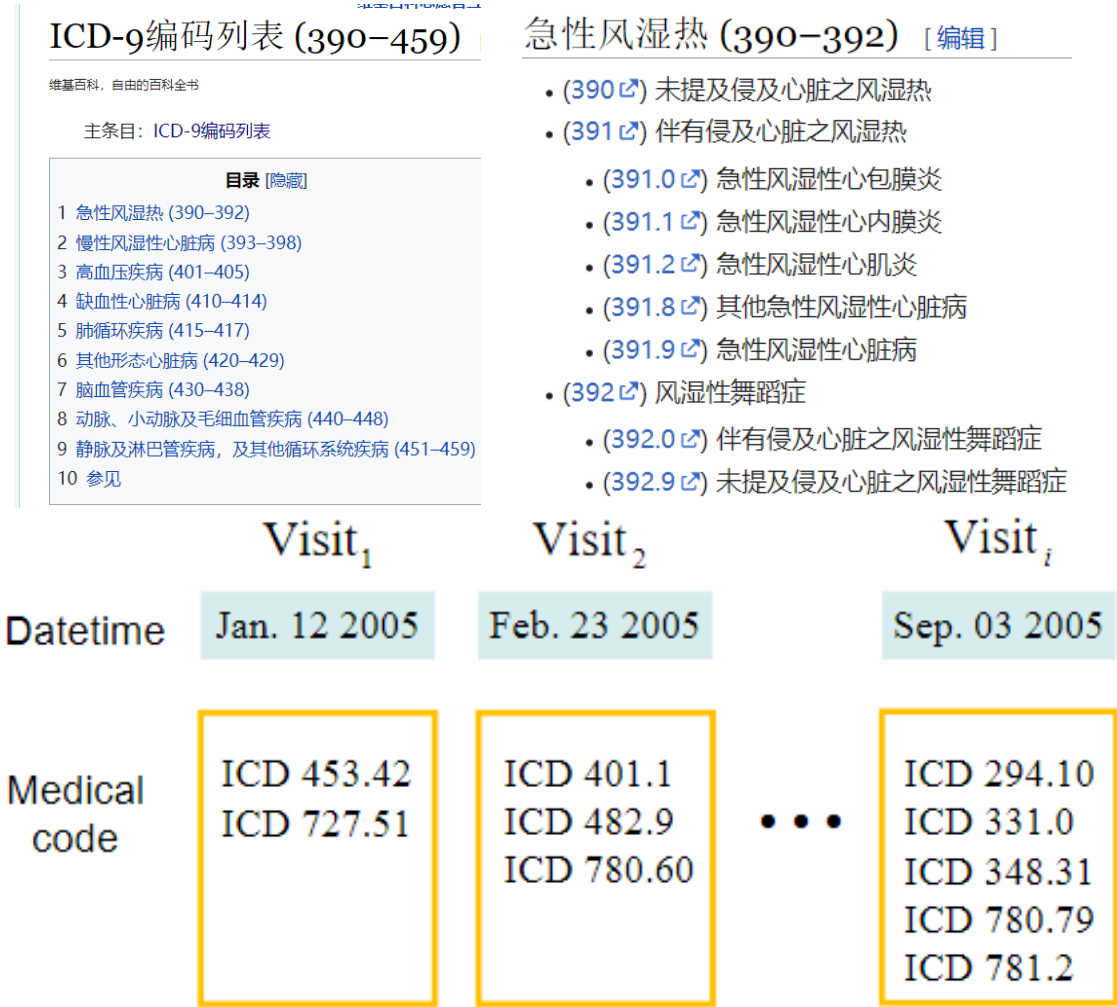


Fig. 1. An example of a patient's EHR data.

二 研究方法

网络结构

E_c : 文本序列输入

E_t : 对应时间序列输入

$$T' = T - \max(T) + 1$$

$$E_t = \text{Conv}(T')$$

$$E_c \in \mathbb{R}^{S \times d_m}$$

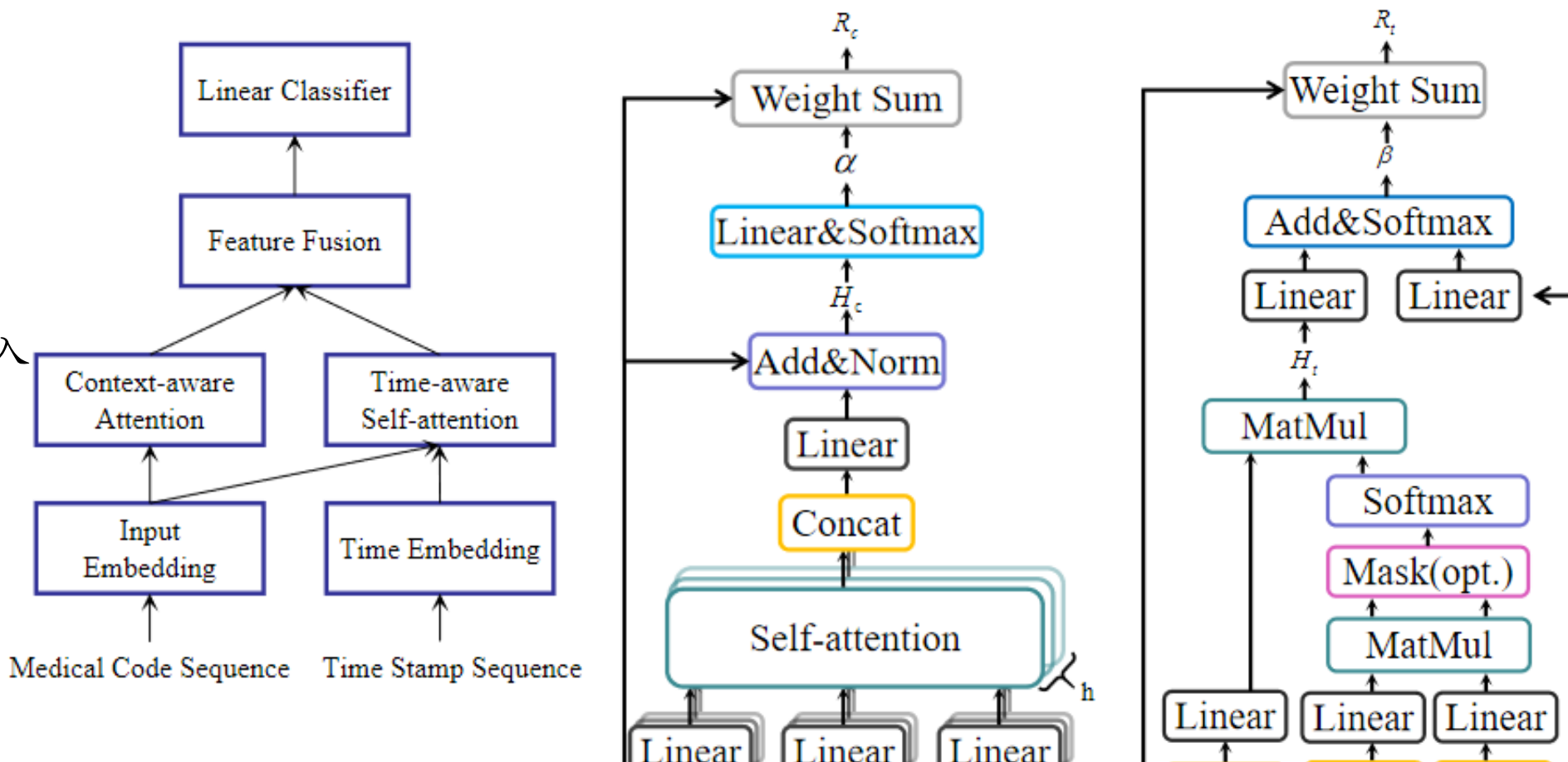
$$E_t \in \mathbb{R}^{S \times d_m}$$

$$R_c = \alpha E_c$$

$$R_t = \beta E_t$$

$$R_f = W_f^T \text{concat}(R_c, R_t) + b_f$$

$$\hat{y} = \text{softmax}(W_p^T R_f + b_p)$$



where $T \in \mathbb{R}^{S \times 1}$ is the input time sequence, $\max(T)$ is the date that we perform prediction, and $T' \in \mathbb{R}^{S \times 1}$ denotes the converted time sequence. Then, we use a 1×1 convolutional

二 研究方法

可解释性

验证本文方法中，文本注意力和时间注意力机制分别对结果产生的贡献。

$$R_f = W_f^\top \text{concat}(R_c, R_t) + b_f, \quad \hat{y} = \text{softmax}(W_p^\top R_f + b_p).$$

$$\hat{y} = \text{softmax}(W_p(W_f \begin{bmatrix} R_c \\ R_t \end{bmatrix} + b_f) + b_p)$$

$$\hat{y} = \text{softmax}(W_p(W_f \begin{bmatrix} \alpha E_c \\ \beta E_c \end{bmatrix} + b_f) + b_p)$$

$$= \text{softmax}(W_p(W_f \begin{bmatrix} \sum_i^S \alpha[i] E_c[i] \\ \sum_i^S \beta[i] E_c[i] \end{bmatrix} + b_f) + b_p)$$

$$= \text{softmax}((\sum_i^S W_p W_f[:, d_m] \alpha[i] E_c[i]$$

$$+ \sum_i^S W_p W_f[d_m : 2d_m] \beta[i] E_c[i] + b_f) + b_p),$$

i denotes the index of the medical code in the sequence.

$$W_f \mathbb{R}^{2d_m \times d_m}$$

- *Rely on context-aware attention only.*

$$\mathcal{C}_i^{(c)} = W_p W_f[:, d_m] \alpha[i] E_c[i]$$

- *Rely on time-aware attention only.*

$$\mathcal{C}_i^{(t)} = W_p W_f[d_m : 2d_m] \beta[i] E_c[i].$$

- *Rely on both context-aware and time-aware attention.*

$$\mathcal{C}_i^{(a)} = \mathcal{C}_i^{(c)} + \mathcal{C}_i^{(t)}$$

三 实验结果

验证本文方法在精度和速度方面的优势。

NLL：损失数值

AUROC：预测类别的效果

Train & Test time/epoch:一个epoch的运算时间

TABLE II
PREDICTION PERFORMANCE OF INPLIM AND BASELINES ON THREE COHORTS

	Model	HF		DIA		CKD	
		NLL	AUROC	NLL	AUROC	NLL	AUROC
Baselines	LR	0.3509	0.6172	0.5550	0.6248	0.4994	0.7239
	MLP	0.3463	0.6632	0.5782	0.6373	0.4968	0.7310
	CNN	0.3187	0.7064	0.5763	0.6193	0.4889	0.7429
	RNN	0.3553	0.7218	0.5935	0.6348	0.4915	0.7626
	RETAIN	0.3317	0.7291	0.5481	0.6457	0.4849	0.7460
	Dipole	0.3367	0.7489	0.5636	0.6509	0.4868	0.7606
	SETF-ATTN	0.3123	0.7483	0.5470	0.6498	0.4794	0.7639
	AdaCare	0.3300	0.7421	0.5536	0.6477	0.4860	0.7648
	SAnD	0.3320	0.7432	0.5560	0.6553	0.4883	0.7646
Ours	INPLIM	0.3083	0.7541	0.5373	0.6658	0.4799	0.7650
	INPLIM _{c-}	0.3101	0.7491	0.5385	0.6634	0.4810	0.7634
	INPLIM _{t-}	0.3102	0.7486	0.5385	0.6631	0.4800	0.7644
	INPLIM _p	0.3284	0.7302	0.5677	0.6336	0.4953	0.7435

Model	HF		DIA		CKD	
	Train time / epoch	Test time / epoch	Train time / epoch	Test time / epoch	Train time / epoch	Test time / epoch
RNN	1.27s	0.11s	0.75s	0.07s	1.57s	0.14s
RETAIN	2.31s	0.17s	2.31s	0.17s	2.84s	0.25s
Dipole	1.88s	0.18s	1.19s	0.09s	2.35s	0.19s
INPLIM	1.65s	0.05s	0.50s	0.03s	1.97s	0.07s

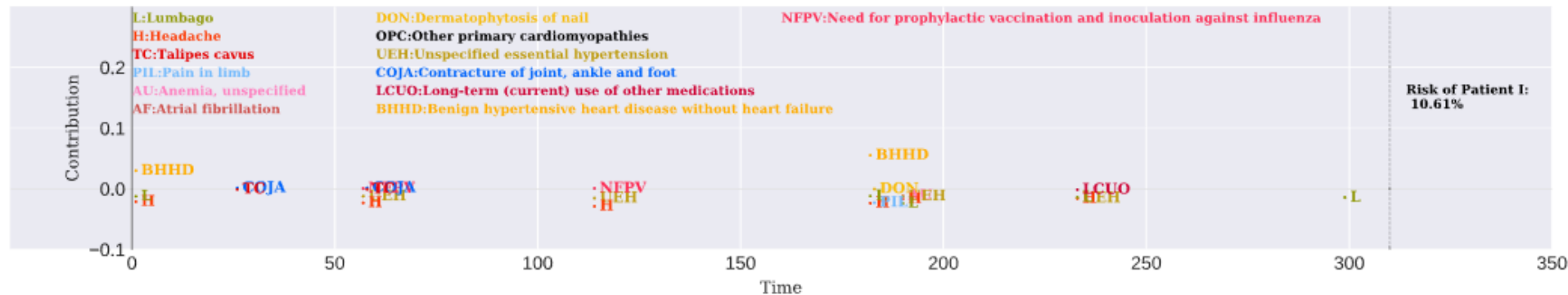
三 实验结果

测试集中对心力衰竭影响最大的前 10 项医疗事件

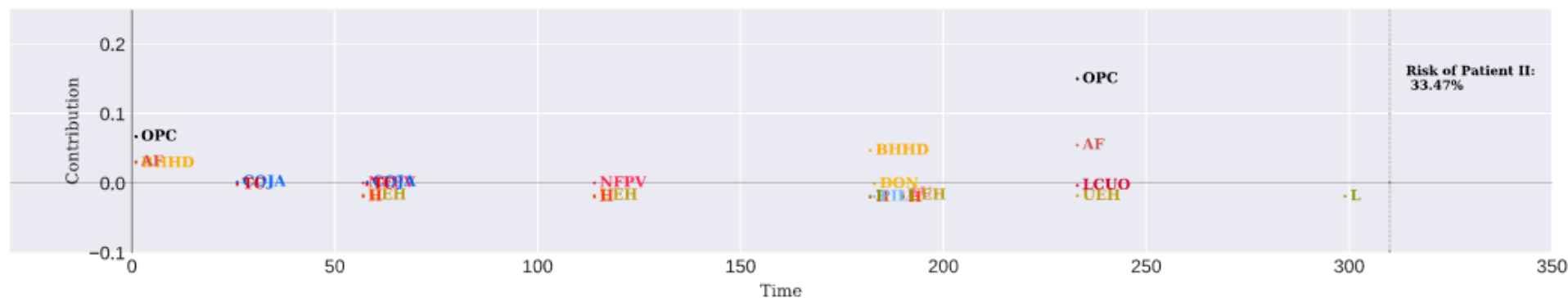
	ICD-9	Code description	Contribution rate	Reasonable
$\mathcal{C}^{(c)}$	585.3	Chronic kidney disease, stage iii (moderate)	6.31	✓
	425.4	Other primary cardiomyopathies	4.40	✓
	402.10	Benign hypertensive heart disease without heart failure	3.37	✓
	585.4	Chronic kidney disease, stage iv (Severe)	1.56	✓
	496	Chronic airway obstruction, not elsewhere classified	1.05	
	461.9	Acute sinusitis, unspecified	0.91	
	278.01	Morbid obesity	0.77	✓
	585.1	Chronic kidney disease, stage i	0.69	✓
	553.21	Incisional hernia without mention Of obstruction or gangrene	0.69	
	V45.02	Automatic implantable cardiac defibrillator In Situ	0.58	✓
$\mathcal{C}^{(t)}$	250.00	Diabetes mellitus without mention of complication, type ii or unspecified type, not stated as uncontrolled	17.34	✓
	427.31	Atrial fibrillation	7.99	✓
	250.02	Diabetes mellitus without mention of complication, type ii or unspecified type, uncontrolled	4.11	✓
	402.10	Benign hypertensive heart disease without heart failure	1.96	✓
	729.1	Myalgia and myositis, unspecified	1.78	
	496	Chronic airway obstruction, not elsewhere classified	1.74	
	414.01	Coronary atherosclerosis of native coronary artery	1.31	✓
	414.00	Coronary atherosclerosis of unspecified type of vessel, native or graft	1.27	✓
	443.9	Peripheral vascular Disease, Unspecified	1.05	✓
	424.0	Mitral valve disorders	0.82	✓
$\mathcal{C}^{(a)}$	250.00	Diabetes mellitus without mention of complication, type ii or unspecified type, not stated as uncontrolled	17.66	✓
	427.31	Atrial fibrillation	8.54	✓
	585.3	Chronic kidney disease, stage iii (moderate)	6.35	✓
	402.10	Benign hypertensive heart disease without heart failure	5.33	✓
	425.4	Other primary cardiomyopathies	5.08	✓
	250.02	Diabetes mellitus without mention of complication, type ii or unspecified type, uncontrolled	4.25	✓
	496	Chronic airway obstruction, not elsewhere classified	2.79	
	729.1	Myalgia and myositis, unspecified	1.93	
	278.01	Morbid obesity	1.84	✓
	414.00	Coronary atherosclerosis of uUnspecified type of vessel, native Or graft	1.81	✓

三 实验结果

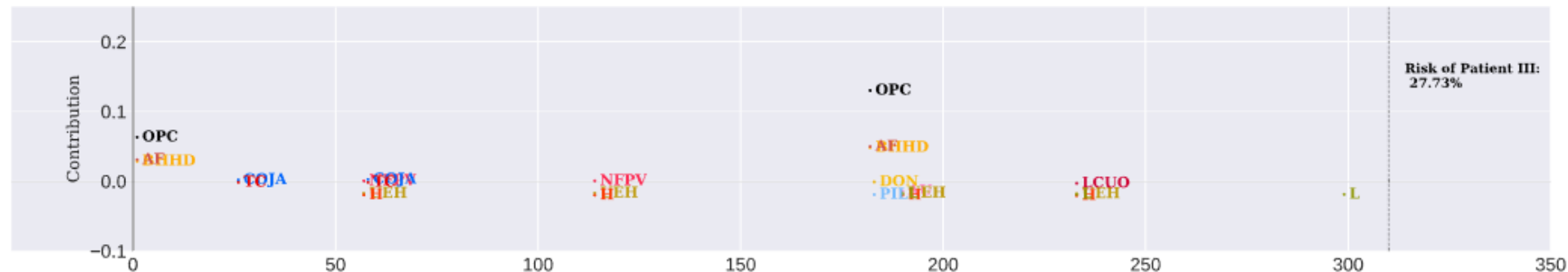
Patient I: A real-world patient's visit records in a test set of Heart Failure cohort.



Patient II: 替换I的L和H为OPC和AF——风险提高的同时，两个检查的共献也提高



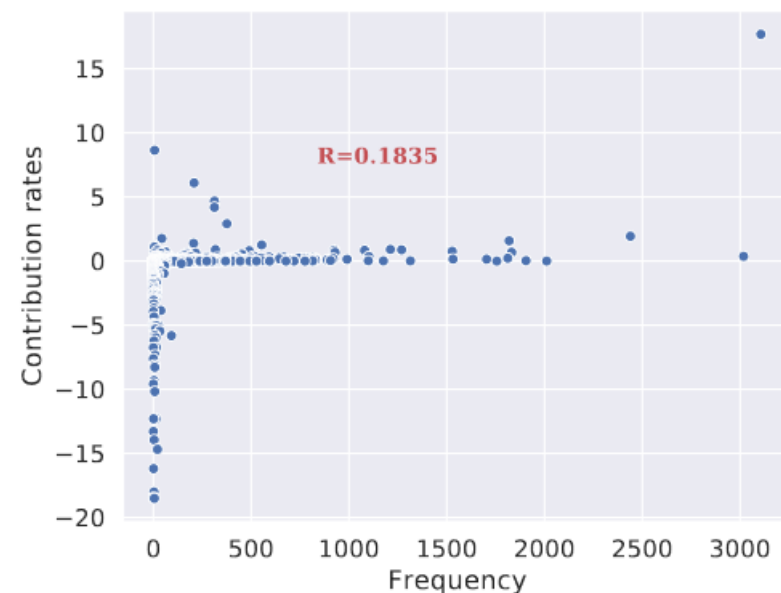
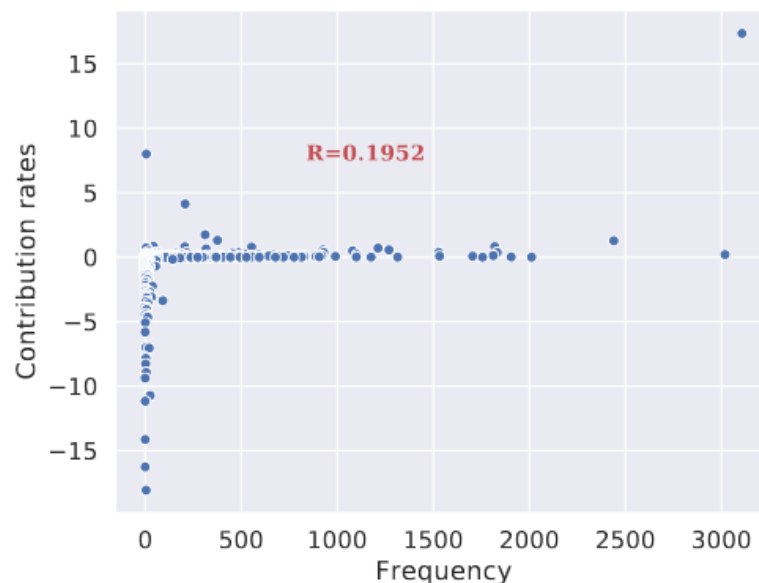
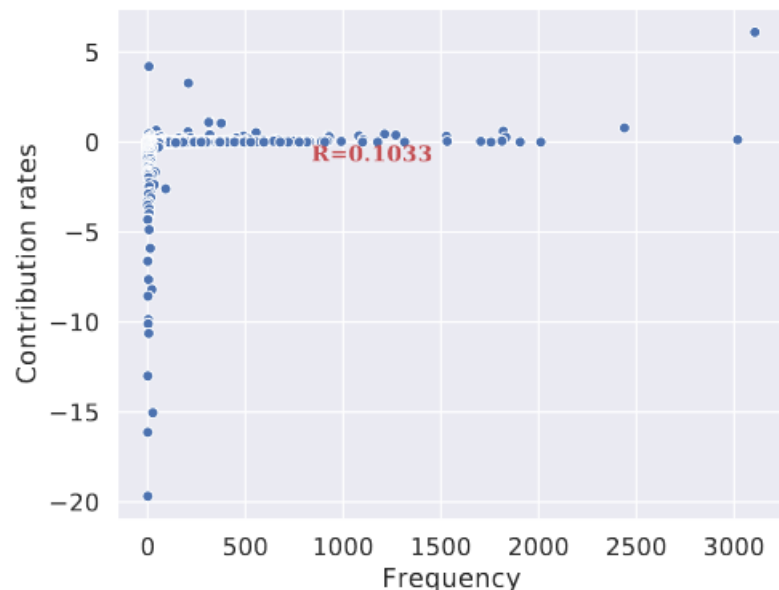
Patient III: 针对于II而言，将OPC和AF提前，发现不同时间段的共献程度不一样。



心力衰竭数据集 L:腰痛 H:头疼 OPC:其他原发性心肌病 AF:心房颤动 BHHD:无心力衰竭的良性高血压心脏病

三 实验结果

可视化训练集中每个医学代码的频率与测试集中相应贡献率之间的相关性。



(a) Correlation between the frequency of each medical code in the training set and the corresponding contribution rate $\mathcal{C}^{(c)}$ in the test set.

(b) Correlation between the frequency of each medical code in the training set and the corresponding contribution rate $\mathcal{C}^{(t)}$ in the test set.

(c) Correlation between the frequency of each medical code in the training set and the corresponding contribution rate $\mathcal{C}^{(a)}$ in the test set.

R是皮尔逊相关系数

Q&A



合肥工業大學