

心電図データにおける特徴抽出

田中 直哉 藤井 章博[†] 清水 宏泰

[†] 法政大学理工学部 〒184-8584 東京都小金井市梶野町 3-7-2

E-mail: [†] fujii@hosei.ac.jp

あらまし 現在, 心疾患の特定や突然死リスクの高低度合いについて, 心電図データから目視で分類が行われている.

本研究では, 心電図データから, 特徴波形を抽出, それをデータセットとして機械学習によるリスク評価を自動的に行う手法を提案する. また, リスク評価について, 複数の機械学習アルゴリズムを比較することにより, 提案手法において最適と思われるアルゴリズムの検討も行った.

キーワード 心電図, 機械学習, 自動診断

Extract feature of electro-cardiogram

Naoya Tanaka Akihiro Fujii[†] and Hiroyasu Shimizu

[†] Faculty of Science and Engineering, Hosei University 3-7-2 Kajino-cho, Koganei-shi, Tokyo 184-8584, Japan

E-mail: [†] fujii@hosei.ac.jp

Abstract Heart diseases have risk of sudden death. The diagnosis of the diseases is usually performed by visual identification from the electrocardiogram data. In this study, we extract waveforms from electrocardiogram data, then the characteristic of the form is classified by machine learning scheme. Risk evaluation is done automatically based on this classification.

We have proposed several machine learning algorithms in terms of risk assessments and compared them to find out optimal methodology for the data set.

Keywords ECG, deep learning, automated diagnosis

1. はじめに

現在, 医療系研究機関および医療系企業では心電図のデータから, 特徴波形の位置と数値を記録するために印刷紙にあるマス目や定規を用いた人の手によって行われている[1]. 近年, 機械学習を用いた診断が医師と同程度に発展し医療に対する機械学習の利用が期待されている[2]. 日本人間ドッグ協会によると心電図の特徴波形の位置データおよび数値データを用いて心疾患の特定や突然死のリスクの高低度合いを分類するための項目が用意されている[3]. そこで, 本研究ではプログラムで特徴波形を捉え, 日本人間ドッグ協会の提唱する項目でデータセットを作成し, 機械学習によってリスクを分類する手法を提案し, その評価を行った.

2. 心電図の自動診断

本研究では時系列データとして電圧値および時間データのある心電図データを用い, 特徴波形の検出, 診断項目に沿ったデータセットの作成, 及びリスク評価までの自動化手法を提案する. そして, 各診断項目における特徴量を分析し, 自動診断の精度を向上させる. また, 精度向上の要因を見つけ, 人の手による診断で誤診を減らす診断項目を発見する手法を提案する.

3. 提案手法

3.1. ホルター心電図の特徴波形

特徴波形には, 肺へ血液を送る時に見られる P 波および全身へ血液を送る時に見られる Q 波・R 波・S 波・T 波からなる[5]. 本研究では, データセット作成のため Q 波から T 波までを検出する. モニター心電図の波形および特徴波形を図 1 に示す.

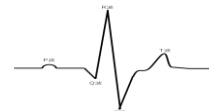


図 1 心電図波形および特徴波形の位置

3.2. 教師あり機械学習のアルゴリズム

既知の情報を学習して未知の情報を分類または回帰する[6]. 本研究で使用するアルゴリズムは「ナイーブベイズ」, 「線形 SVM」, 「SVM(カーネル法)」, 「ランダムフォレスト」, 「kNN」, 「ニューラルネットワーク」である[7].

3.3. 主成分分析

情報量をなるべく失わないように次元を圧縮させる学習アルゴリズムで, 本研究では多次元データの可視化および因子分析に用いる[7].

4. 検出方法

本研究では特徴波形を検出後に検出された特徴波形の位置および数値を用いてデータセットを作成し、機械学習で突然死リスクの高低度合いを分類する。

4.1. 特徴波形の検出手法

- ① 10秒間の心電図データから極大値を検出
- ② 極大値における前後の勾配の大きさ上位30個を記録
- ③ 記録された極大値の周囲-0.45秒から+0.4秒で超える値がない場合にR波として記録
- ④ R波の+0.15秒から+0.4秒間までの最大値をT波として記録
- ⑤ R波の-0.2秒までの極小値の最小値をQ波として記録
- ⑥ R波とT波の間にある最小値をS波として記録

5. データセットの作成方法

検出された特徴波形のピーク値から以下の手法で各診断項目の算出基準を用いて機械学習用のデータセットを作成した。

I. 高いT波(尖度)

$$Tk = \frac{1}{2n} \sum_{i=0}^n \left\{ \left(\frac{Ty_i - Ty_{i-1}}{Tx_i - Tx_{i-1}} \right) + \left(\frac{Ty_i - Ty_{i+1}}{Tx_i - Tx_{i+1}} \right) \right\} \dots ①$$

II. 高いT波(数値)

$$Tv = \frac{1}{n} \sum_{i=0}^n \frac{Ty_i}{Ry_i} \dots ②$$

III. QT 間隔

$$QT = \frac{1}{n} \sum_{i=0}^n (Tx_i - Qx_i) \dots ③$$

IV. R-R' 間隔

$$aveR = \frac{1}{n-1} \sum_{i=0}^{n-1} (Rx_{i+1} - Rx_i)$$

$$RR = \sqrt{\frac{1}{n} \sum_{i=0}^n (Rx_i - aveR)^2} \dots ④$$

V. ST 上昇

$$ST = \frac{1}{n} \sum_{i=0}^n (Ty_i - Sy_i) \dots ⑤$$

VI. T波オルタナンス

$$TS = \frac{Ty_i - \min(Ty)}{\max(Ty) - \min(Ty)} - 0.5 \dots ⑥$$

$$sumT = \sum_{i=0}^{n-2} \{ |(TS_i + TS_{i+2}) - (TS_i + TS_{i+1})| \} \dots ⑦$$

$$TWA = \frac{|sumT|}{(n-2) \cdot \text{std}(TS)} \dots ⑧$$

VII. QT・R-R' 間隔

$$QTRR = \text{corr}(QT, RR) \dots ⑨$$

6. 突然死リスクにおける高低度合の分類方法

作成したデータセットを用いて機械学習による分類を行う。使用した機械学習のアルゴリズムは「ナイーブベイズ」「線形SVM」「SVM(カーネル法)」「ランダムフォレスト」「kNN」「ニューラルネットワーク」の6種類である。診断において最も適切なアルゴリズムを検討するために各種機械学習によって訓練データとテストデータを毎回変えて150回分類テストを行い正解率を記録する。

7. 突然死リスクの分類結果

7.1. 特徴波形の検出結果

4.の手法で検出した心電図データの成功例を図2に、失敗した例を図3に示す。ここで、特徴波形の検出は180個のデータセットの中で126個のデータで成功した。失敗したデータの内訳としては正常波形が14個で低リスク波形は19個となっており高リスク波形では21個であった。

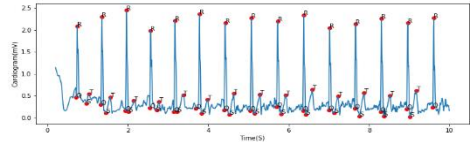


図2 検出に成功した心電図の例(正常波形)

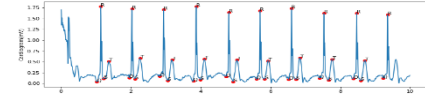


図3 検出に成功した心電図の例(低リスク波形)

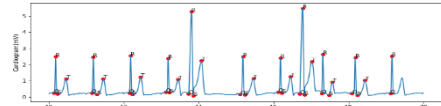


図4 検出に成功した心電図の例(高リスク波形)

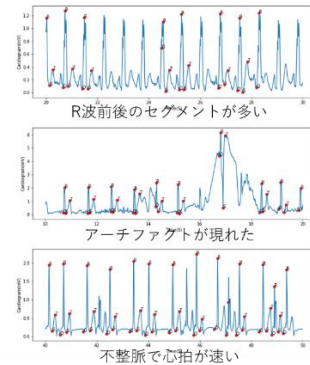


図5 検出に失敗した心電図データの例

7.2. 機械学習による分類

6.の手法で分類した正解率の分布を図4に示す。また、各学習アルゴリズムにおける分類精度の詳細を表1に記す。図4から中央値および平均値とともに最も高いランダムフォレストおよびニューラルネットワークの正解率上位25%から75%まで精度は約80%台となった。

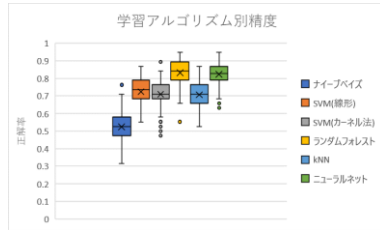


図 6 学習アルゴリズム別における正解率の分布

表 1 学習アルゴリズム別精度の詳細

	ナイーブベイズ	SVM(線形)	SVM(カーネル法)	ランダムフォレスト	kNN	ニューラルネットワーク
最大値	0.763158	0.868421	0.894737	0.947368	0.868421	0.947368
上位25%	0.578947	0.782895	0.763158	0.888158	0.763158	0.868421
中央値	0.526316	0.736842	0.710526	0.842105	0.710526	0.828947
平均値	0.523509	0.724211	0.709649	0.831404	0.707368	0.823158
下位25%	0.473684	0.684211	0.684211	0.789474	0.657895	0.789474
最小値	0.315789	0.552632	0.473684	0.552632	0.526316	0.631579
標準偏差	0.085485	0.070634	0.074271	0.070197	0.071302	0.059286

7.3. データの分布と決定境界

特徴波形の検出に成功した 126 個のデータを主成分分析で二次元に圧縮し、機械学習の決定境界を可視化したグラフを図 7 に示す。

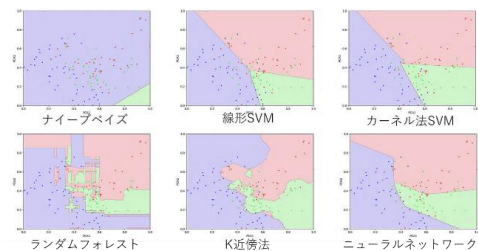


図 7 圧縮したデータの分布および各種機械学習の決定境界

8. データの特徴分析

分類テストで用いたニューラルネットワークで精度が中央値に近かった時のテストデータと予測のメトリクス表を表 2 に記す。また、この時の適合率および再現率ならびに F 値を表 3 に記す。

表 2 分類テストにおけるメトリクス表

	予 測			
テ ス ト		正 常 波 形	低 リ ス ク	高 リ ス ク
	正 常 波 形	15	0	2
	低 リ ス ク	0	17	2
	高 リ ス ク	0	4	11

表 3 分類テストにおける適合率および再現率ならびに F 値

	適合率	再現率	F 値	総数
正常波形	1.00	0.88	0.94	17
低リスク	0.81	0.89	0.85	19
高リスク	0.73	0.73	0.73	15

表 3 で適合率が低く、また表 8-1 から分類結果の誤りにおいて高リスク波形が最も原因に寄与していることが分かる。

9. 分類に寄与する因子の分析

分類テストにおいて、7 個の説明変数を用いて分類

を行っているが、この中で分類に重要な因子および重要でない因子を見つけることで分類の精度および速度の向上を行う。

9.1. 重要度に基づく因子分析

因子の重要度の指標として、教師あり学習を用いて分類経過から因子分析を行う。使用するアルゴリズムはランダムフォレストとする。分類経過において 150 回テストしてサンプル数およびジニ係数から算出した重要度数の代表値を表 4 に記す。また、リスク別各項目の分布を図 8 に示す。

表 4 ランダムフォレストを用いた各因子における重要度数の代表値

因子	平均値	中央値
高い T 波(尖度)	0.12	0.12
高い T 波(数値)	0.12	0.12
QT 間隔	0.37	0.38
R-R' 間隔	0.15	0.15
ST 上昇	0.12	0.12
T 波オルタナンス	0.05	0.05
QT・R-R' 間隔	0.06	0.06

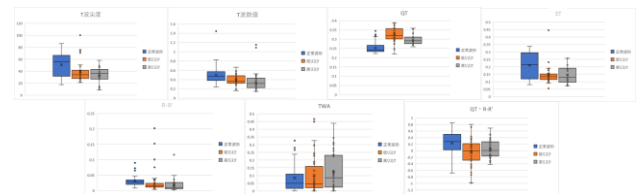


図 8 リスク別各項目の分布

9.2. シミュレーションに基づく因子分析

データセットとして記録してあるデータすべてを教師データとして最大最小正規化して学習し、テストデータには各項目に 0 から 1 の乱数を用いて分類を行う。分類結果における説明変数の分布を記録して有意差が見られる診断項目を選定する。シミュレーション結果から得られた分類時の分布を図 9 に示す。

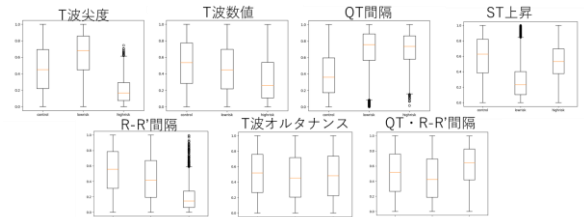


図 9 シミュレーション結果から得られた分類時の分布

10. 因子を選定した場合における分類結果

9 項において得られた重要度およびシミュレーション結果から表 5 に記す因子の選定を行い、ニューラルネットワークを用いて 70%を訓練データとして 30%をテストデータにして毎回訓練データとテストデータを変え 150 回精度を検証した結果得られた精度の分布を図 10 に、統計的代表的値を表 6 に示す。

表 5 データセットに用いる項目

診断項目	重要度	シミュレーション
高い T 波(尖度)	○	○
高い T 波(電圧)	○	×
QT 間隔	○	○
ST 上昇	○	○
R-R' 間隔	○	○
T 波オルタナンス	×	×
QT・R-R' 間隔	×	×

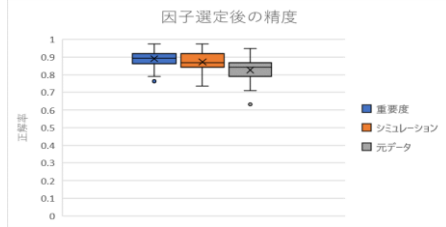


図 10 新しいデータセットにおける学習アルゴリズム別における正解率の分布

表 6 新しいデータセットにおける学習アルゴリズム別精度の詳細

	重要度	シミュレーション	元データ
最大値	0.973684211	0.973684211	0.947368421
上位25%	0.921052632	0.921052632	0.868421053
中央値	0.894736842	0.868421053	0.842105263
平均値	0.890526316	0.872105263	0.825789474
下位25%	0.868421053	0.842105263	0.789473684
最小値	0.763157895	0.736842105	0.631578947
標準偏差	0.048118014	0.052986265	0.054681894

11. 時間と波形の特徴からデータセットの作成

11.1. 波形データの特徴からデータセットの作成

図 8 から電圧値に着目した場合に T 波の電圧に患者のリスクごとの特徴が現れた。また、時間データは計測機器・計測方法に依存しないため時間データとして波形の間隔の時間と標準偏差を用いて表 7 のようにデータセットを作成した。

表 7 T 波の電圧と時間に着目したデータセット

診断項目	算出方法
T 波(電圧)	$\text{mean}(T_y)$
QT 間隔(時間)	$\text{mean}(T_x - Q_x)$
QT 間隔(標準偏差)	$\text{std}(T_x - Q_x)$
R-R' 間隔(標準偏差)	$\text{std}(R_{x+1} - R_x)$
R-R' 間隔(時間)	$\text{mean}(R_{x+1} - R_x)$
Q-Q' 間隔(標準偏差)	$\text{std}(Q_{x+1} - Q_x)$
Q-Q' 間隔(時間)	$\text{mean}(Q_{x+1} - Q_x)$

作成したデータセットを用いて 7.2 項と同様に分類を行い精度を計測した。得られた機械学習別精度の分布を図 11 に示す。また、得られた機械学習別精度の統計的代表値を表 8 に記す。

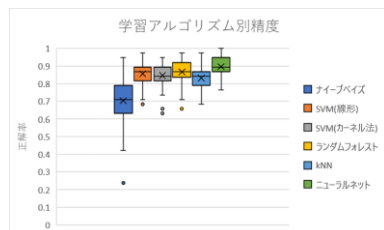


図 11 新しく作成した診断項目における精度分布

表 8 新しく作成した診断項目における精度の分布の詳細

	ナイーブベイズ	SVM(線形)	SVM(カーネル法)	ランダムフォレスト	kNN	ニューラルネットワーク
最大値	0.947368421	0.973684211	0.947368421	0.973684211	0.973684	1
上位25%	0.789473684	0.894736842	0.894736842	0.921052632	0.868421	0.947368421
中央値	0.710526316	0.868421053	0.842105263	0.868421053	0.842105	0.894736842
平均値	0.703859649	0.855964912	0.847017544	0.866491228	0.831404	0.896491228
下位25%	0.631578947	0.815789474	0.815789474	0.842105263	0.789474	0.868421053
最小値	0.236842105	0.684210526	0.631578947	0.657894737	0.684211	0.763157895
標準偏差	0.116862926	0.055107488	0.056661889	0.056391376	0.055141	0.048750079

11.1.1. データの分布と決定境界

特徴波形の検出に成功した 126 個のデータを主成分分析で二次元に圧縮し、機械学習の決定境界を可視化したグラフを図 12 に示す。

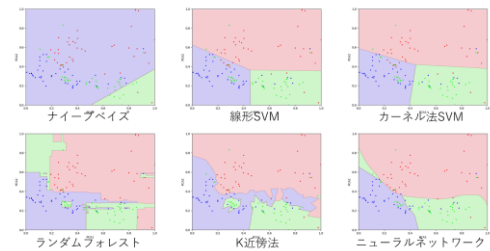


図 12 圧縮したデータの分布および各種機械学習の決定境界

11.2. 分類に寄与する因子の分析

分類テストにおいて、7 個の説明変数を用いて分類を行っているが、この中で分類に重要な因子および重要でない因子を見つけることで分類の精度および速度の向上を行う。

11.2.1. 重要度に基づく因子分析

9.1 項と同様の手法を用いて各項目の重要度を測定した。サンプル数およびジニ係数から算出した重要度数の代表値を表 9 に記す。また、リスク別各項目の分布を図 13 に示す。

表 9 ランダムフォレストを用いた各因子における重要度数の代表値

	中央値	平均値
T 波(電圧)	0.099629	0.099026
QT 間隔(時間)	0.338388	0.341234
QT 間隔(標準偏差)	0.16785	0.167144
R-R 間隔(標準偏差)	0.117919	0.118178
R-R 間隔(時間)	0.095974	0.096084
Q-Q 間隔(標準偏差)	0.079247	0.081329
Q-Q 間隔(時間)	0.096855	0.097004

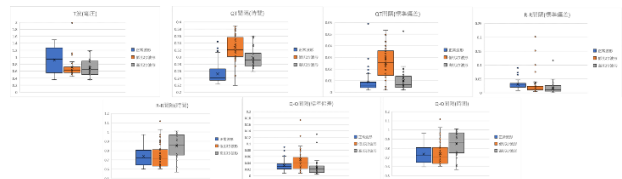


図 13 リスク別各項目の分布

11.2.2. シミュレーションに基づく因子分析

データセットとして記録してあるデータすべてを教師データとして最大最小正規化して学習し、テストデータには各項目に0から1の乱数を用いて分類を行う。分類結果における説明変数の分布を記録して有意差が見られる診断項目を選定する。シミュレーション結果から得られた分類時の分布を図14に示す。

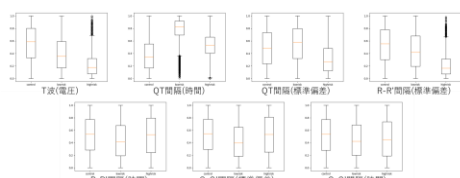


図14 シミュレーション結果から得られた分類時の分布

11.3. 因子を選定した場合における精度

11.1項において得られた重要度およびシミュレーション結果から表10に記す因子の選定を行い、ニューラルネットワークを用いて70%を訓練データとして30%をテストデータにして毎回訓練データとテストデータを変え150回精度を検証した結果得られた精度の分布を図15に、統計的代表的値を表11に示す。

表10 データセットに用いる項目

Q-Q間隔(時間)	×	×
Q-Q間隔(標準偏差)	×	×
R-R'間隔(時間)	×	×
R-R'間隔(標準偏差)	○	○
QT間隔(標準偏差)	×	○
QT間隔(時間)	○	○
T波(電圧)	○	○

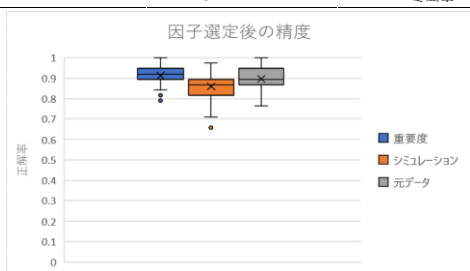


図15 新しいデータセットにおける学習アルゴリズム別における正解率の分布

表11 新しいデータセットにおける学習アルゴリズム別精度の詳細

	重要度	シミュレーション	元データ
最大値	1	0.973684211	1
上位25%	0.947368421	0.894736842	0.947368421
中央値	0.921052632	0.868421053	0.894736842
平均値	0.912105263	0.859122807	0.896491228
下位25%	0.894736842	0.815789474	0.868421053
最小値	0.789473684	0.657894737	0.763157895
標準偏差	0.044269485	0.054466143	0.048750079

12. データの要約

12.1. 因子負荷量

主成分分析を行い7項および11項におけるデータセットに対して7次元データを2次元に圧縮を行い分布を可視化を行った。7項のデータセットにおける因子負荷量を表12に、11項のデータセットにおける因子負荷量を表13に記す。

表12 7項のデータセットにおける因子負荷量

	第1主成分	第2主成分
T波尖度	-0.30834	0.184338
T波数値	-0.15531	0.097305
QT	0.300984	0.319573
ST	0.088041	0.140123
R-R'	-0.26973	0.209883
TWA	0.212245	-0.02959
QT_RR	-0.2004	-0.13134

表13 11項のデータセットにおける因子負荷量

	第1主成分	第2主成分
T波(電圧)	0.015086	-0.07611
QT間隔(時間)	0.358954	0.16709
QT間隔(標準偏差)	0.116533	0.291213
R-R'間隔(標準偏差)	0.108564	0.136536
R-R'間隔(時間)	0.340861	-0.20143
Q-Q間隔(標準偏差)	0.143794	0.208734
Q-Q間隔(時間)	0.340602	-0.20241

表12の結果から二次元に圧縮したデータにおいて表14のように要約できる。

表14 7項のデータを二次元に圧縮したデータの分布

第一主成分			
第二主成分		負	正
	正	TWA(小) T波尖度(大) R-R'間隔(大)	QT間隔(大) TWA(大) QT・RR(小)
	負	QT間隔(小) TWA(小) QT・RR(大)	TWA(大) T波尖度(小) R-R'間隔(小)

また、表13の結果から二次元に圧縮したデータにおいて表15のように要約できる。

表15 11項のデータを二次元に圧縮したデータの分布

第一主成分			
第二主成分		負	正
	正	RRt(小) QQt(小) QTs(大) QQs(大) Tv(小)	QTt(大) QTs(大) QQs(大) Tv(小)
	負	QTt(小) QTs(小) QQs(小) Tv(大)	RRt(大) QQt(大) QTs(小) QQs(小) Tv(大)

12.2. データの分布

7 項におけるデータセットを主成分分析で 2 次元に圧縮した結果得られた分布を図 16 に示す. また, 11 項におけるデータセットを主成分分析で 2 次元に圧縮した結果得られた分布を図 17 に示す.

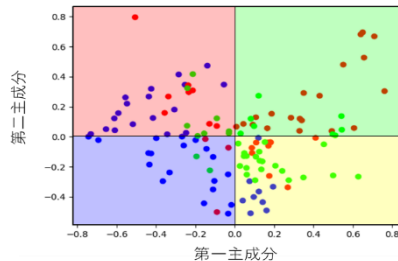


図 16 7 項のデータを 2 次元に圧縮して得られた分布

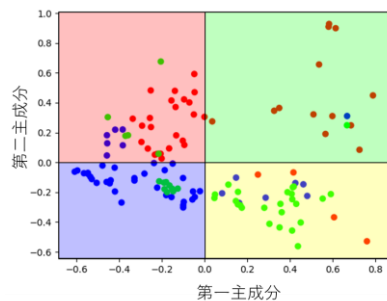


図 17 11 項のデータを 2 次元に圧縮して得られた分布

図 16 および図 17 において青点は正常波形. 赤点は正常波形, 緑点は高リスク波形を意味している.

図 16 から 7 項のデータセットにおいて正常波形は第二・第三象限に集中していることため T 波尖度・R-R' 間隔・QTRR 間隔の数値が高く QT 間隔が小さいことが分かる. 次に低リスク波形は第一・第二象限に集中しているため QT 間隔・R-R' 間隔・T 波尖度が高く QTRR 間隔が小さい事が分かる. また, 高リスク波形は第四象限に集中しているため T 波オルタナンスが高く T 波尖度・R-R' 間隔が小さいことが分かる.

図 17 から 11 項のデータセットにおいて正常波形は第三象限に集中しているため QT 間隔の時間・標準偏差・Q-Q' 間隔の標準偏差が小さく T 波の電圧が高いことが分かる. 次に低リスク波形は第一・第二象限に集中しているため QT 間隔の時間・標準偏差・Q-Q' 間隔の標準偏差が大きく T 波の電圧・Q-Q' 間隔の時間・R-R 間隔の時間が小さいことが分かる. また, 高リスク波形は第 4 象限に集中しているため R-R' 間隔の時間・Q-Q' 間隔の時間・T 波の電圧が高く Q-Q' 間隔の標準偏差と QT 間隔の標準偏差が低いことが分かる.

13. 結論

13.1. 現時点における自動診断

現時点で 84% の確率で心電図データから特徴波形を検出できた. また米 NCBI の発表する誤診率は 31.7% となっており [8], 今回行ったテストでニューラルネットワークを用いた場合に最小値を誤診率としたときでは表 11 から 21.1% となりリスク判定そのものでは下回った.

13.2. システムを用いない場合の診断項目

表 11 より, QT 間隔(標準偏差)の追加の有無で精度が大きく変わることが分かり, 図 13 および図 17 から低リスク波形となる生活習慣病の判別に利用できる.

また, T 波オルタナンス図 8 および図 16 から高リスク波形となる重大な心疾患の判別に使うことができるが, 表 4 より重要度は高くないため T 波オルタナンスが出た波形を危険な波形として扱うことができる.

14. 今後の展望

現時点で自動診断をする際, 障害となることが特徴波形の検出である. 特徴波形の検出はグラフを用いた画像データであるため図 2 から図 3 のように検出に成功した画像および失敗した画像を用いて機械学習による画像認識を行い自動で成功および失敗を選別する.

また, 今回使用したデータセットは 30 人と少ないデータを用いて検証したため今後はより多くのデータを用いて検証を行う必要がある.

文 献

- [1] やさしい心電図の見方 - ecg1-4.pdf
<https://med.toaeiyo.co.jp/contents/ecg/pdf/ecg1-4.pdf>
- [2] 子どもの病気, AI が医師と同精度で診断
<https://www.afpbb.com/articles/-/3210672>
- [3] 心電図 - 日本人間ドッグ協会
<https://www.ningen-dock.jp/public/inspection/electrocardiogram>
- [4] 法政大学 理工学部 応用情報工学科 八名研究室
<http://bsi.ws.hosei.ac.jp/index.html>
- [5] 基本心電図波形
http://www.cardiac.jp/view.php?lang=ja&target=normal_ecg_pattern.xml
- [6] 人工知能と機械学習 - ict_skill_3_5.pdf
http://www.soumu.go.jp/ict_skill/pdf/ict_skill_3_5.pdf
- [7] Python Data Science Handbook | Python Data Science Handbook
<https://jakevdp.github.io/PythonDataScienceHandbook/>
- [8] Clinical and autopsy diagnoses in the intensive care unit: a prospective study.
<https://www.ncbi.nlm.nih.gov/pubmed/14980989>