

機械学習における 評価指標 ~AUC&C-index~

長崎大学病院 初期研修医 中島誉也



参考文献

- イベント予測モデルの評価指標
- 岡田遥平先生のnote
 - 臨床における良い予測モデルとは？
 - 一番わかりやすいROC曲線とAUCとC統計量

Table of contents

1. 評価指標の紹介
2. Discriminationについて
3. C-indexについて

Table of contents

1. 評価指標の紹介

2. Discriminationについて

3. C-indexについて

評価指標の紹介

- 識別能 (Discrimination) :

「予測モデルが正しく予測対象のイベント発生のリスクが高い患者とリスクが低い患者を区別できる性能」

- 較正能 (Calibration) :

「予測がどれくらい実際に当たるか」

- 臨床的有用性 (Clinical utility)

予測モデルの究極的な目標は

「予測に基づき治療方針が変化し重要なアウトカムが改善する」

Table of contents

1. 評価指標の紹介

2. AUC(C統計量)について

3. C-indexについて

混同行列(Confusion Matrix)

		予測	
		負(Negative)	正(Positive)
実際	負(Negative)	真陰性 (True Negative/TN)	偽陽性 (False Positive/FP)
	正(Positive)	偽陰性 (False Negative/FN)	真陽性 (True Positive/TP)

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall \text{ (True Positive Rate, TPR)} = \frac{TP}{TP + FN}$$

$$F\text{-measure} = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

$$Specificity \text{ (True Negative Rate, TNR)} = \frac{TN}{TN + FP}$$

$$False \text{ Positive Rate (FPR)} = \frac{FP}{TN + FP}$$

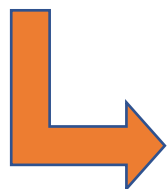
混同行列の問題点

1. データに偏りがある場合

元のデータでPositiveが99%, Negativeが1%のような場合, 評価指標の数値は偏ったものになる.

2. カットオフ値の設定

カットオフ値をどこにするかで当然PositiveとNegativeの割合が変化し, そのカットオフ値が適切なのかが曖昧になる.



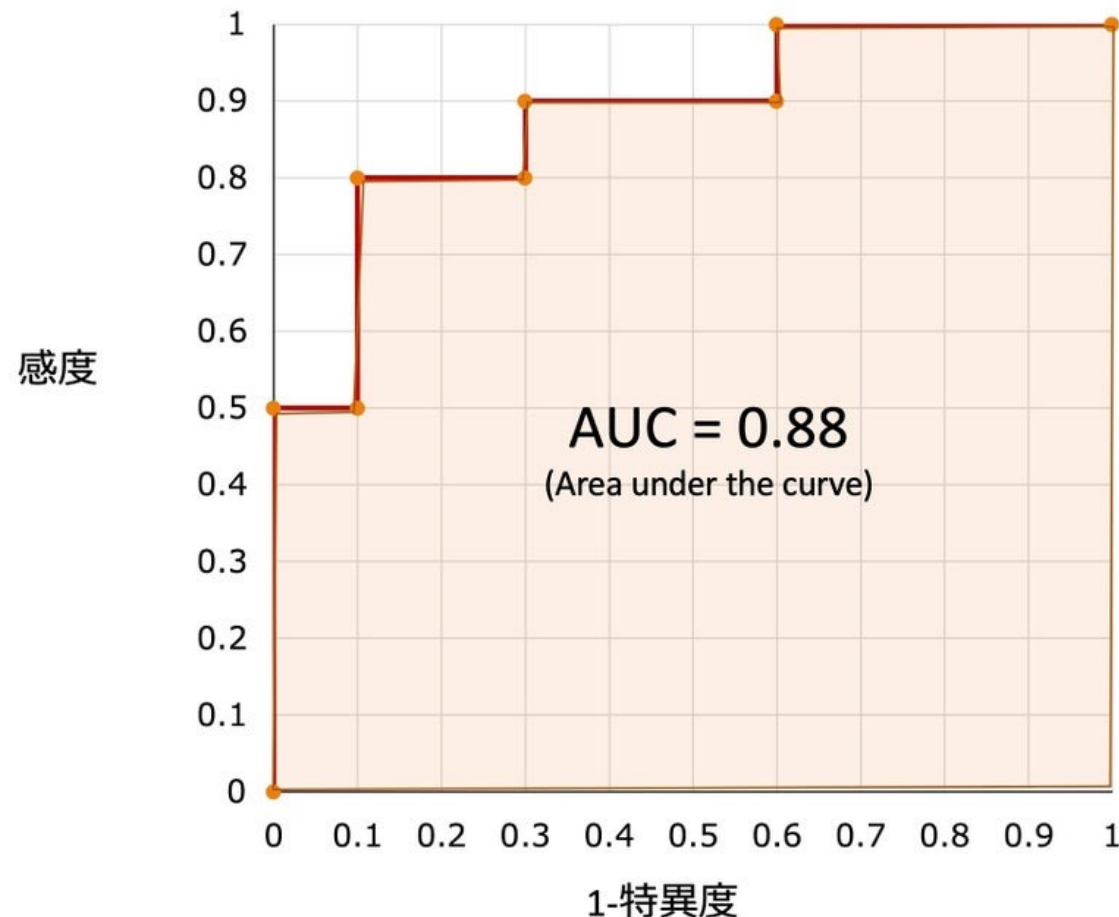
AUC(C統計量)による評価

AUCについて

[岡田遥平先生のnote](#)を引用しています

「対象としている連続変数（例：検査の値、スコア、または予測確率）が、二値のアウトカム（例えば生存/死亡）をどの程度正確に識別できるか」

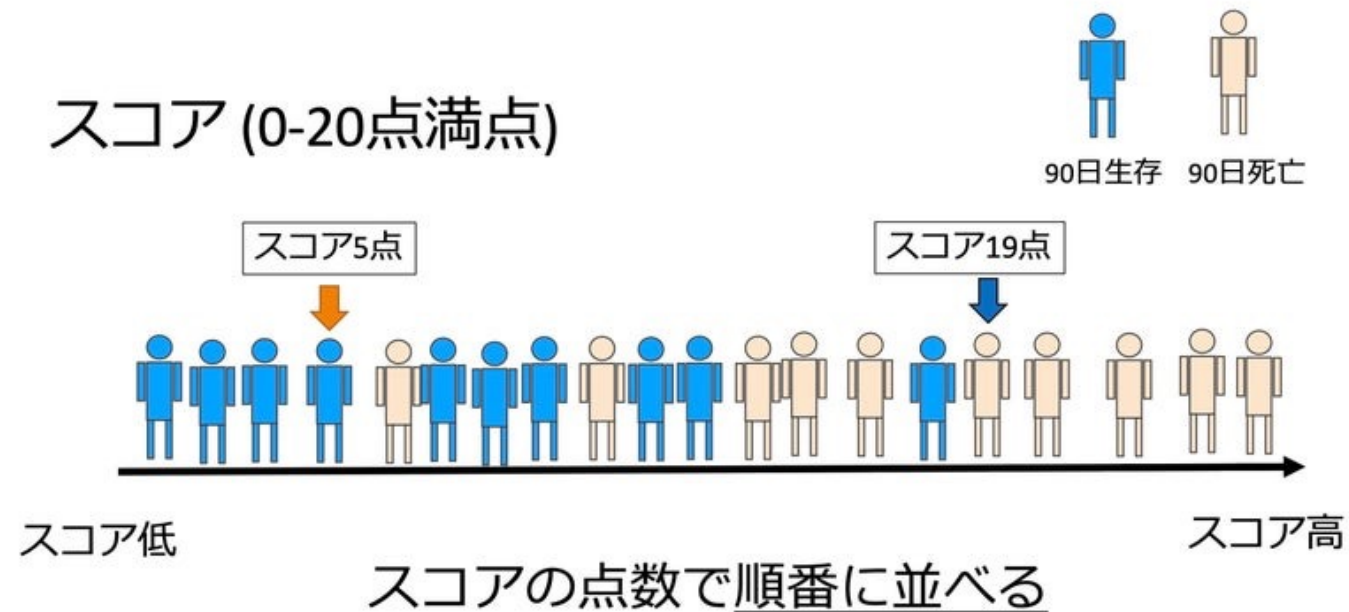
ROC曲線



スコア	アウトカム
0	生存
1	生存
2	生存
5	生存
6	死亡
7	生存
8	生存
10	生存
11	死亡
12	生存
15	生存
16	死亡
16	死亡
17	死亡
18	生存
19	死亡
19	死亡
19	死亡
20	死亡
20	死亡

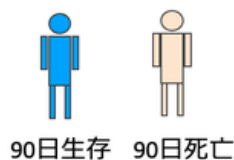
SOFAスコア：敗血症患者の重症度を表すスコア

スコア (0-20点満点)

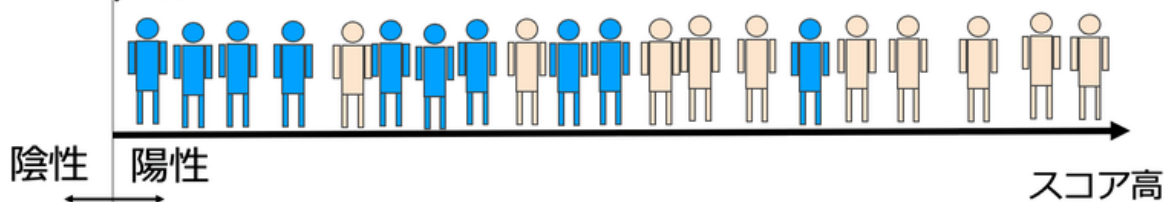


- 1 と 0 の間をカットオフにする

スコア (0-20点満点)



-1/0点

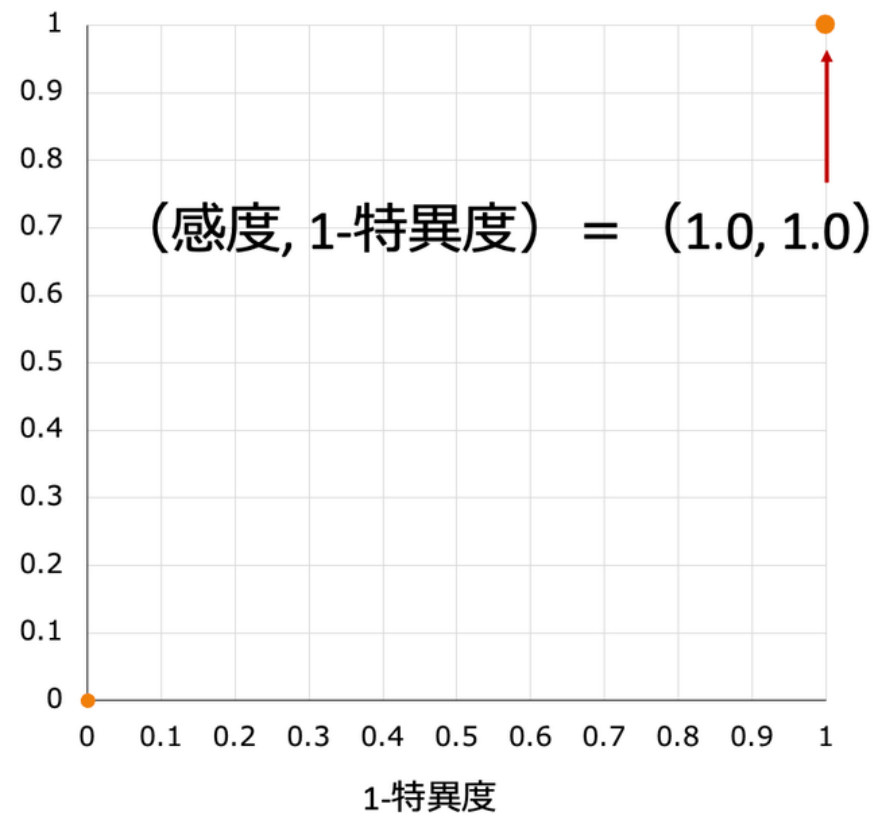


-1/0点をカットオフにする

感度100%
特異度0%

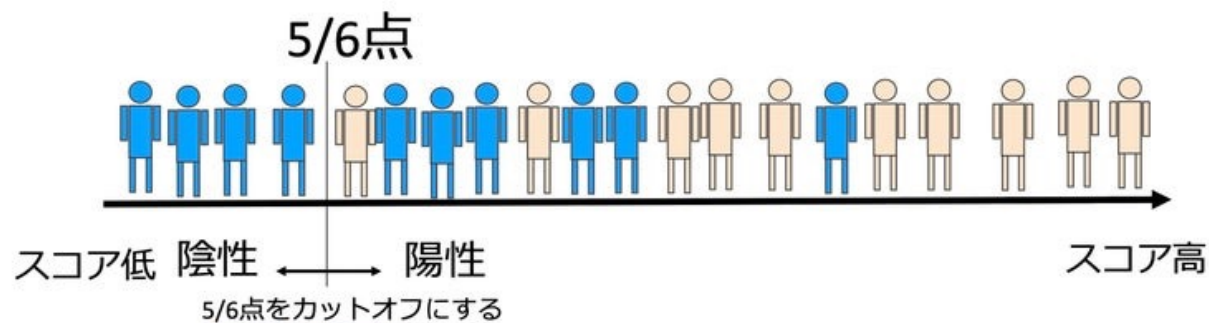
	死亡	生存
陽性	10	10
陰性	0	0

感度



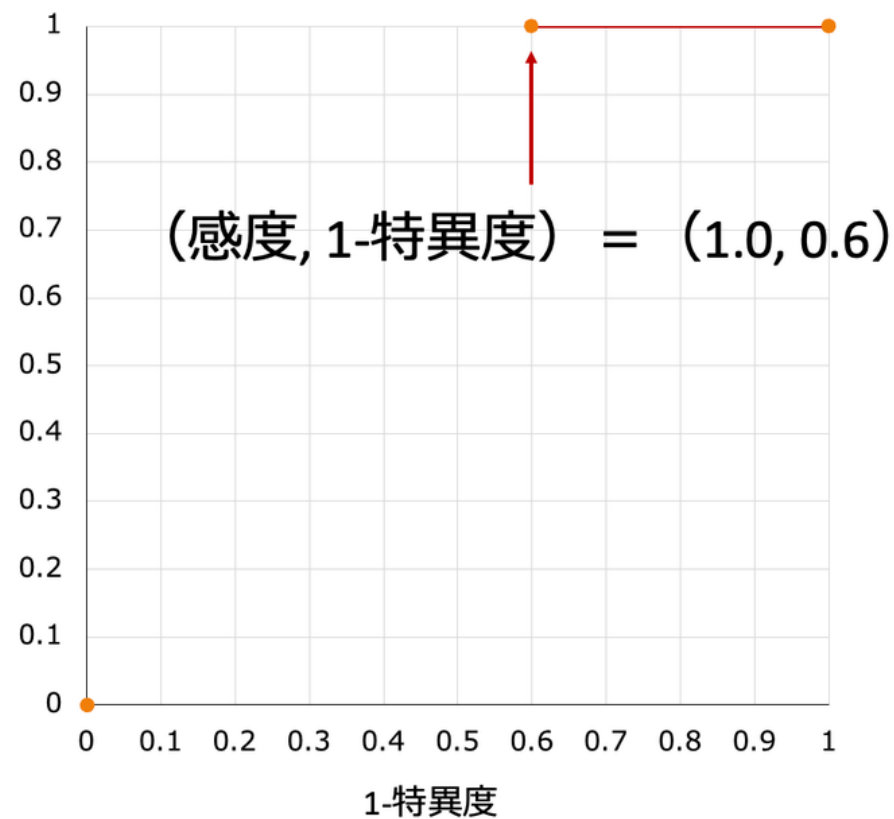
5と6の間をカットオフにする

スコア (0-20点満点)

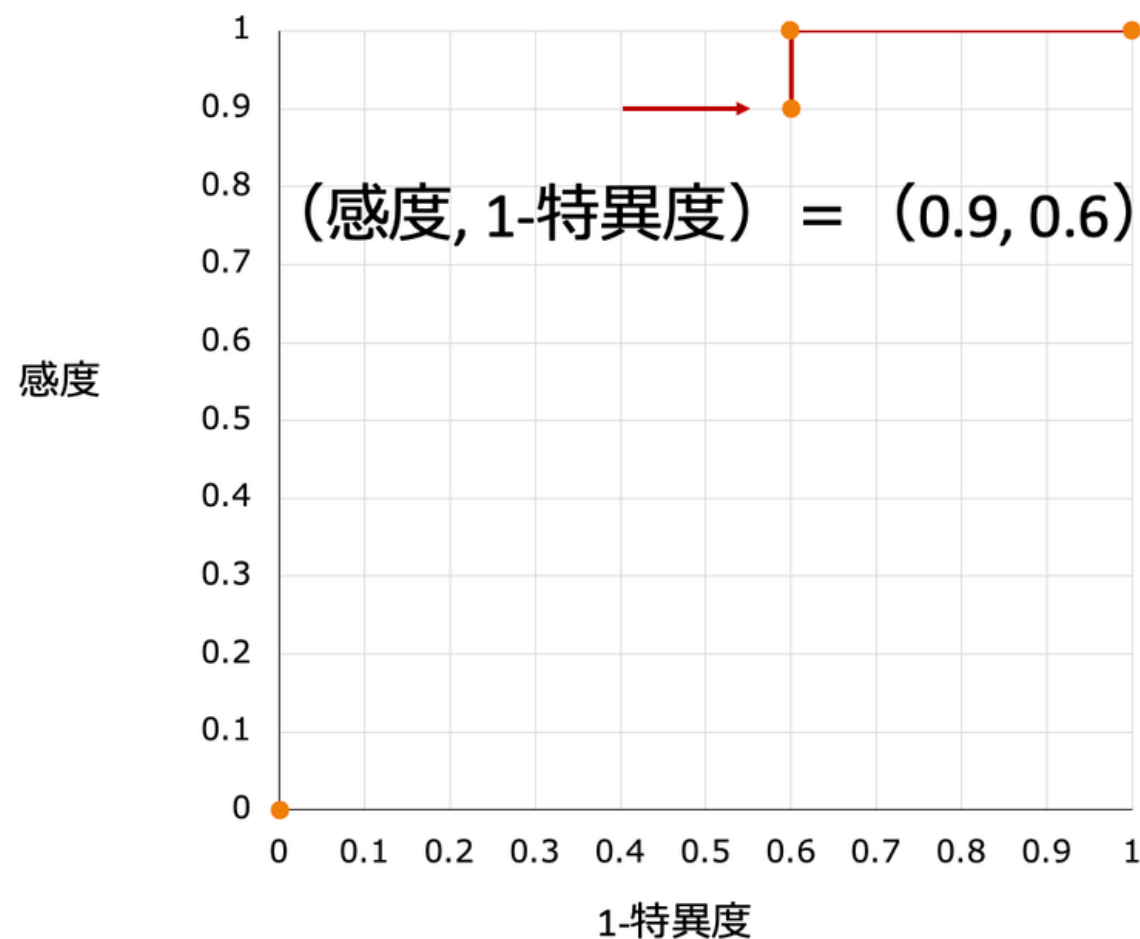
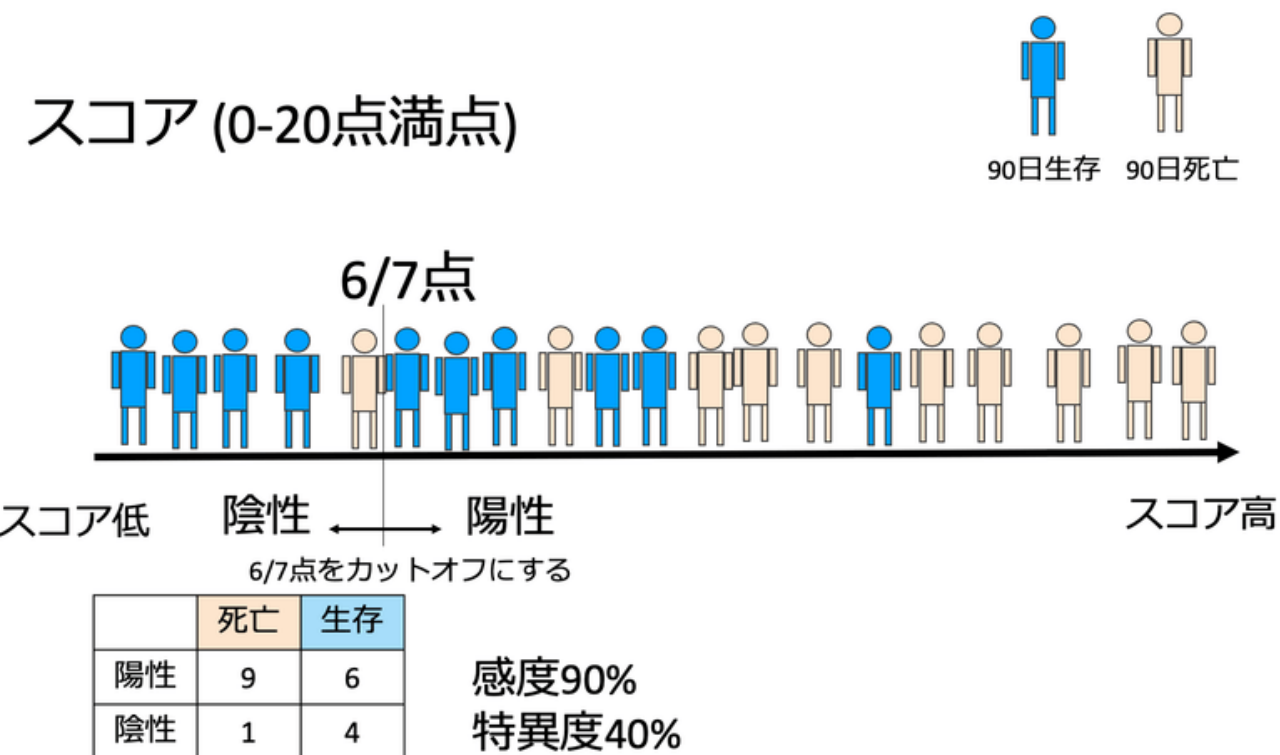


	死亡	生存	
陽性	10	6	感度100%
陰性	0	4	特異度40%

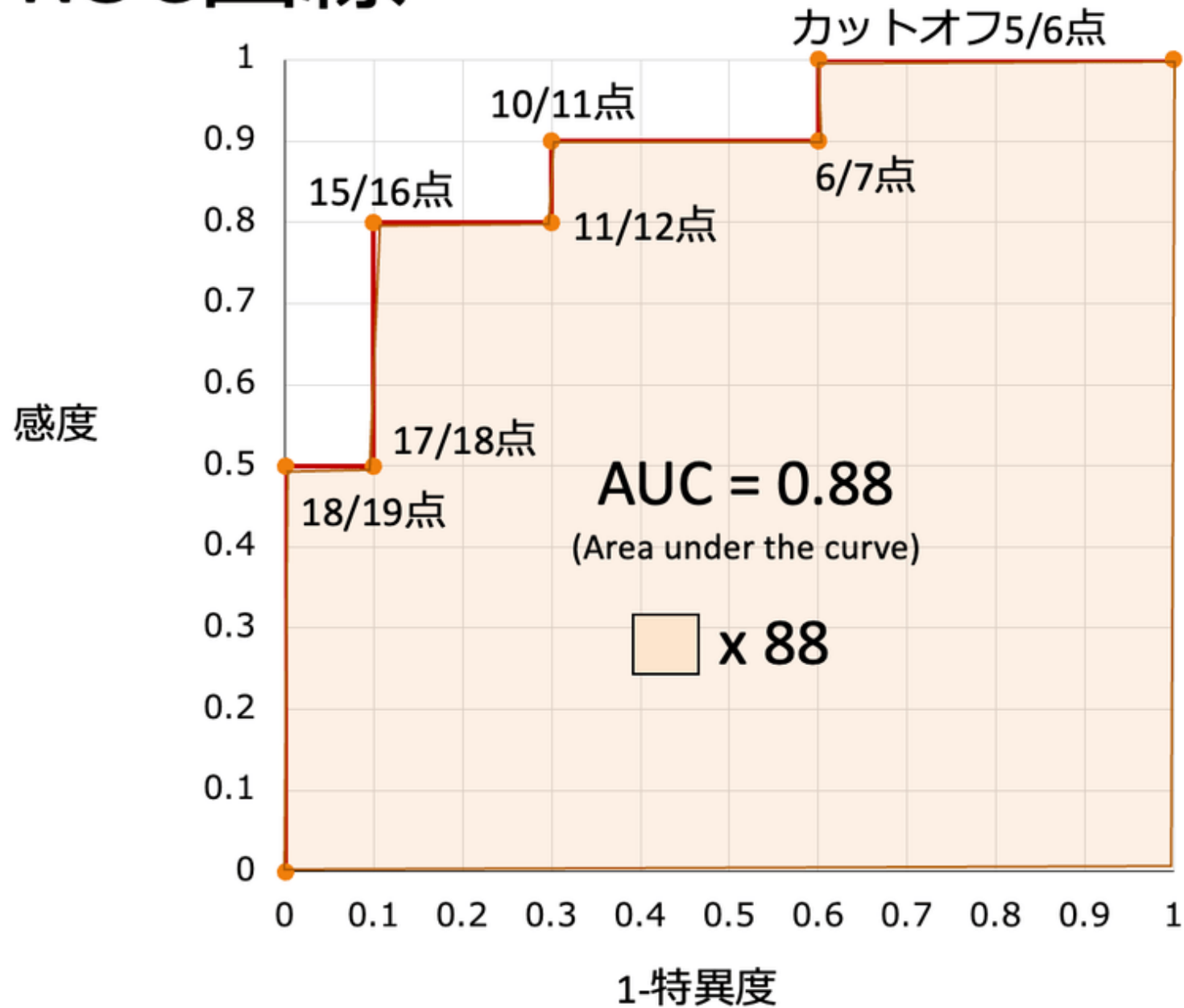
感度



6と7の間をカットオフにする



ROC曲線



C統計量について

[岡田遥平先生のnote](#)を引用

「アウトカムのイベントなし（生存）の患者とアウトカムのイベントあり（死亡）の患者を、ランダムに1人ずつ抽出した時に、それぞれの患者のスコアの大小関係が、
（抽出した生存の患者のスコア） < （抽出した死亡の患者のスコア） となる確率」



構築したスコアやモデルが適切に識別できている確率

計算してみよう！！

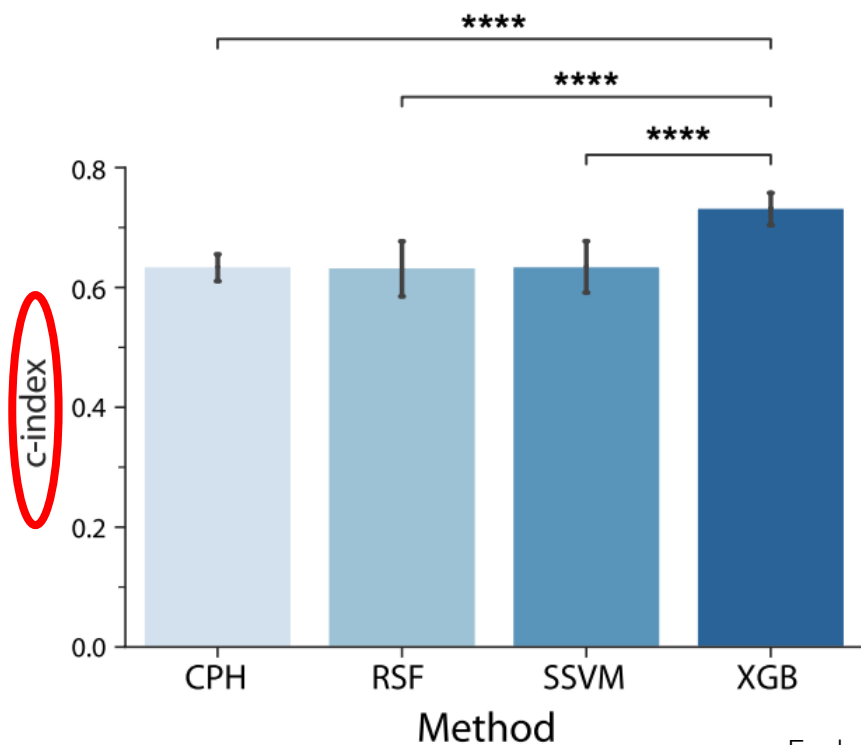
[岡田遥平先生のnote](#)を参照

Table of contents

1. 評価指標の紹介
2. AUCについて
- 3. C-indexについて**

C-indexについて

- 生存時間解析の分野においても、機械学習を用いた予測モデルの開発の研究が盛んに行われている。
- 打ち切りの概念を取り入れなければならないため、AUCで評価することは困難。
- 打ち切りの概念を取り入れつつ、AUCのように評価したい！！ → **C-index**の登場



様々な生存時間解析 × 機械学習モデル

Task	Model	sksurv	lifelines	statsmodels	pycox
Survival Function Estimation	Non-parametric	✓	✓	✓	✓
	Parametric	✗	✓	✗	✗
CHF Estimation	Non-parametric	✓	✓	✗	✗
	Parametric	✗	✓	✗	✗
Linear Regression	Cox	✓	✓	✓	✓
	Cox + Elastic-Net	✓	✓	✓	✗
	Log-normal AFT	✓	✓	✗	✗
	Log-logistic AFT	✗	✓	✗	✗
	Weibull AFT	✗	✓	✗	✗
	Piecewise exponential	✗	✓	✗	✓
	Aalen	✗	✓	✗	✗
	Gradient Boosted AFT	✓	✗	✗	✗
	Gradient Boosted Cox	✓	✗	✗	✗
Non-linear Regression	Heterogenous Ensemble	✓	✗	✗	✗
	NN (Grouped survival times)	✗	✗	✗	✓
	NN (Proportional hazards)	✗	✗	✗	✓
	NN (Piecewise exponential)	✗	✗	✗	✓
	Random Survival Forest	✓	✗	✗	✗
	Survival SVM	✓	✗	✗	✗
	Survival Tree	✓	✗	✗	✗
	Brier Score	✓	✗	✗	✓
	Concordance Index	✓	✓	✗	✓
Evaluation	Time-dependent ROC	✓	✗	✗	✗

Table 1: Availability of methods. AFT: Accelerated Failure Time. CHF: Cumulative Hazard Function. NN: Neural Network. SVM: Support Vector Machine.

C-indexについて

C統計量(Harrell's C-Statistic.) :

異なるアウトカムじゃないと×

イベント発症者($Y_i=1$)の予測確率 p_i の分布と非発症者($Y_j=0$)の予測確率 p_j

の分布それぞれからランダムサンプリングした時、発症者での予測確率の方が高くなる確率。



同じアウトカムでも良い！！

C-index(C-statistic by Uno et al.) :

標本からランダムに2つの異なるデータ i, j を取り出してペアにした時、観測打ち切り時点 t までの生存時間 T の短長と予測確率 $P(t)$ の大小が一致する確率。

標本からランダムに2つの異なるデータ i, j を取り出してペアにした時、
観測打ち切り時点 c までの生存時間 T の短長と予測確率 $P(t)$ の大小が一致する確率。



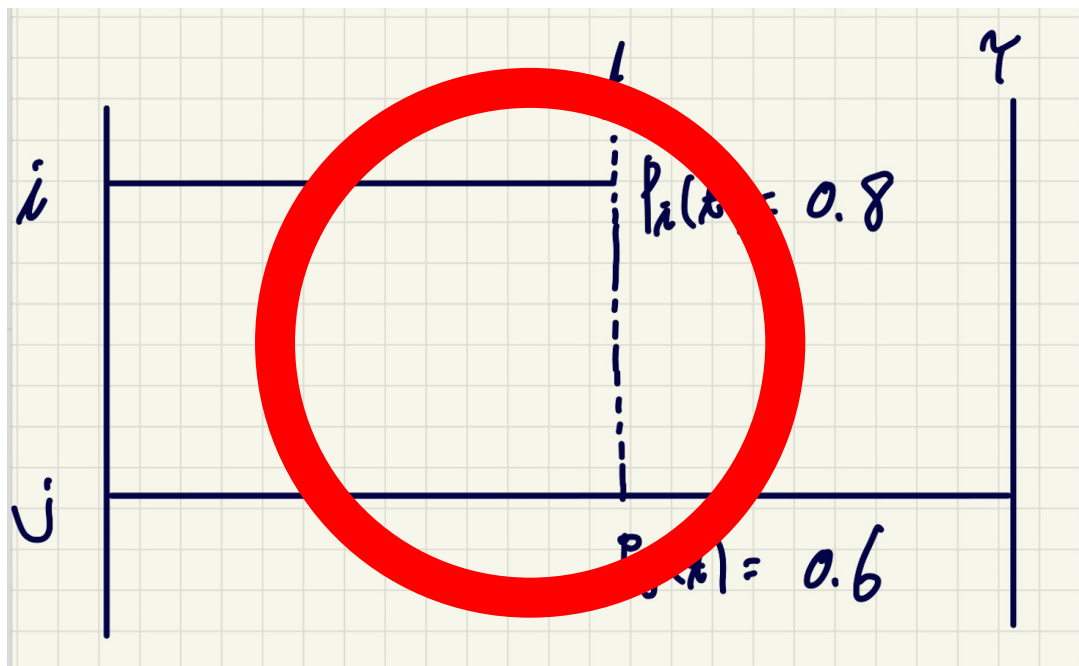
ランダムに2人(i と j)を抽出した時に、観測期間 c までの生存期間の長さの短長と
そのモデルによって予測される時点 t の死亡確率の大小が一致する確率。



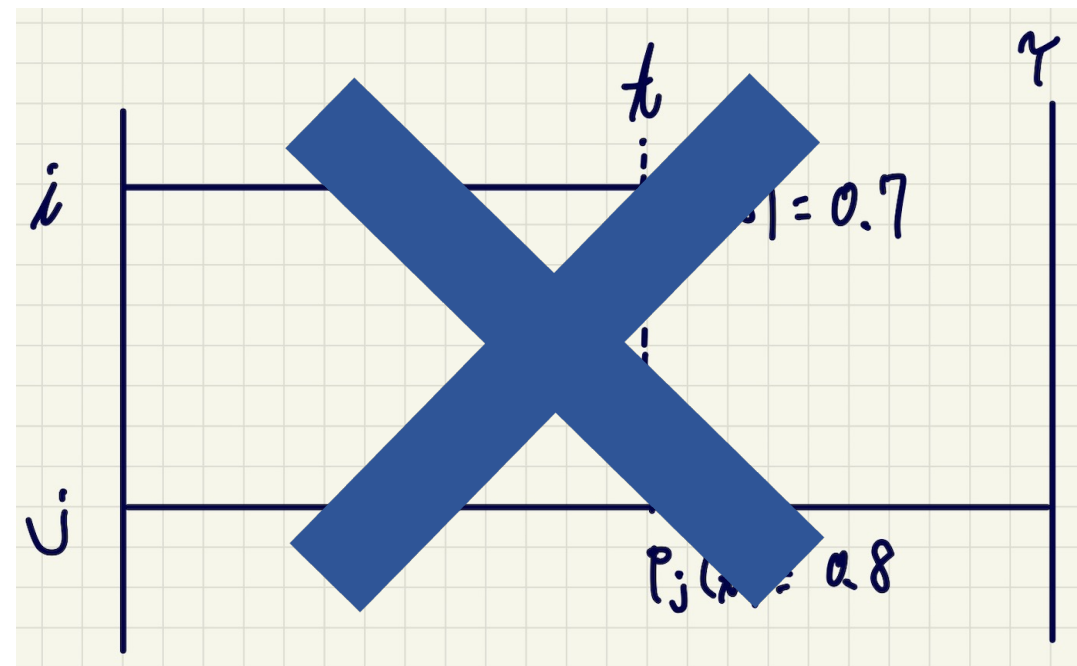
適切に識別できてる！！

i と j を比較した時に、 i の方が、実際のデータでも j より早く死亡(=時点 t)しており、
かつ

モデルによって予測される時点 t における死亡確率が j よりも大きければcountできる！



i の方が、早く死亡(=時点 t)しており、
予測される時点 t での死亡確率が j よりも大きい
(=正しい予測)



i の方が、早く死亡(=時点 t)しているが、
予測される時点 t での死亡確率が j よりも小さい
(=誤った予測)

C-indexの式(Uno, et al., 2011)

$$\hat{C}_\tau = \frac{\sum_{i,j} D_i I(T_i^* < T_j^*) I\{P_i(t) > P_j(t)\}}{\sum_{i,j} D_i I(T_i^* < T_j^*)}$$

- 分母が比較可能な（comparable）ペアの総数
- 分子はそのうち T と P(t) の短長と大小が揃っている（concordant）ペアの総数を表す。
- ペアが比較可能である必要十分条件は、ペア内で観察時間 T * の短い方のイベントが観測されていること（= i にアウトカムが発生していること）である。（D_i = 1）

計算してみよう！！

No.	Outcome	Time(year)
1	生存	10
2	死亡	3
3	死亡	7
4	生存	10
5	生存	10
6	死亡	2
7	死亡	8
8	生存	10
9	死亡	1
10	生存	10

心筋梗塞による死亡をアウトカムとした生存時間解析
観測期間は10年とする.

Random Survival Forestを用いて行った死亡確率の予測結果

観測時点 (時点t)	No.1	No.2 (死亡)	No.3 (死亡)	No.4	No.5	No.6 (死亡)	No.7 (死亡)	No.8	No.9 (死亡)	No.10
1	0.03	0.12	0.1	0.03	0	0.2	0.11	0.06	0.31	0.08
2	0.06	0.19	0.12	0.04	0.02	0.3	0.17	0.06	0.34	0.11
3	0.1	0.22	0.12	0.05	0.06	0.31	0.25	0.07	0.41	0.14
7	0.1	0.3	0.13	0.09	0.07	0.4	0.23	0.08	0.42	0.17
8	0.12	0.32	0.2	0.13	0.17	0.44	0.27	0.09	0.43	0.17
10	0.16	0.36	0.26	0.18	0.18	0.47	0.3	0.13	0.46	0.19

観測時点 (時点t)	No.1	No.2 (死亡)	No.3 (死亡)	No.4	No.5	No.6 (死亡)	No.7 (死亡)	No.8	No.9 (死亡)	No.10
1	0.03	0.12	0.1	0.03	0	0.2	0.11	0.06	0.31	0.08
2	0.06	0.19	0.12	0.04	0.02	0.3	0.17	0.06	0.34	0.11
3	0.1	0.22	0.12	0.05	0.06	0.31	0.25	0.07	0.41	0.14
7	0.1	0.3	0.13	0.09	0.07	0.4	0.23	0.08	0.42	0.17
8	0.12	0.32	0.2	0.13	0.17	0.44	0.27	0.09	0.43	0.17
10	0.16	0.36	0.26	0.18	0.18	0.47	0.3	0.13	0.46	0.19

$$\hat{C}_\tau = \frac{\sum_{i,j} D_i I(T_i^* < T_j^*) I\{P_i(t) > P_j(t)\}}{\sum_{i,j} D_i I(T_i^* < T_j^*)}$$

e.g.)

i をNo.9とすると、時点t = 1, τ = 10であるから、
比較可能となる候補 j は、No.1, 2, 3, 4, 5, 6, 7, 8, 10 (9組)

例えば、i = No.9, j = No.1とすると、

$$P_i(t) = P_{No.9}(t = 1) = 0.31$$

$$P_j(t) = P_{No.1}(t = 1) = 0.03$$

$$P_i(t) > P_j(t) \text{であるから、}$$

i の方が、早く死亡(=時点t)しており、

予測される時点tでの死亡確率がjよりも大きい → countされる

観測時点 (時点t)	No.1	No.2 (死亡)	No.3 (死亡)	No.4	No.5	No.6 (死亡)	No.7 (死亡)	No.8	No.9 (死亡)	No.10
1	0.03	0.12	0.1	0.03	0	0.2	0.11	0.06	0.31	0.08
2	0.06	0.19	0.12	0.04	0.02	0.3	0.17	0.06	0.34	0.11
3	0.1	0.22	0.12	0.05	0.06	0.31	0.25	0.07	0.41	0.14
7	0.1	0.3	0.13	0.09	0.07	0.4	0.23	0.08	0.42	0.17
8	0.12	0.32	0.2	0.13	0.17	0.44	0.27	0.09	0.43	0.17
10	0.16	0.36	0.26	0.18	0.18	0.47	0.3	0.13	0.46	0.19

$$\hat{C}_\tau = \frac{\sum_{i,j} D_i I(T_i^* < T_j^*) I\{P_i(t) > P_j(t)\}}{\sum_{i,j} D_i I(T_i^* < T_j^*)}$$

J = No.2の場合もcountされる. (::0.31>0.12)
 J = No.3の場合もcountされる. (::0.31>0.1)
 J = No.4の場合もcountされる. (::0.31>0.03)
 J = No.5の場合もcountされる. (::0.31>0)
 J = No.6の場合もcountされる. (::0.31>0.2)
 J = No.7の場合もcountされる. (::0.31>0.11)
 J = No.8の場合もcountされる. (::0.31>0.06)
 J = No.10の場合もcountされる. (::0.31>0.08)



$$C_{i=No.9} = \frac{9}{9} = 1$$

観測時点 (時点t)	No.1	No.2 (死亡)	No.3 (死亡)	No.4	No.5	No.6 (死亡)	No.7 (死亡)	No.8	No.9 (死亡)	No.10
1	0.03	0.12	0.1	0.03	0	0.2	0.11	0.06	0.31	0.08
2	0.06	0.19	0.12	0.04	0.02	0.3	0.17	0.06	0.34	0.11
3	0.1	0.22	0.12	0.05	0.06	0.31	0.25	0.07	0.41	0.14
7	0.1	0.3	0.13	0.09	0.07	0.4	0.23	0.08	0.42	0.17
8	0.12	0.32	0.2	0.13	0.17	0.44	0.27	0.09	0.43	0.17
10	0.16	0.36	0.26	0.18	0.18	0.47	0.3	0.13	0.46	0.19

$$\hat{C}_\tau = \frac{\sum_{i,j} D_i I(T_i^* < T_j^*) I\{P_i(t) > P_j(t)\}}{\sum_{i,j} D_i I(T_i^* < T_j^*)}$$

同様に、

i = No.6の時、 比較可能群(= j)はNo. 1, 2, 3, 4, 5, 7, 8, 10 (8組) , $C_{i=No.6} = \frac{8}{8} = 1$

i = No.2の時、 比較可能群(= j)はNo. 1, 3, 4, 5, 7 8, 10 (7組) , $C_{i=No.2} = \frac{6}{7} \doteq 0.86$

i = No.3の時、 比較可能群(= j)はNo. 1, 4, 5, 7, 8, 10 (6組) , $C_{i=No.3} = \frac{4}{6} \doteq 0.67$

i = No.7の時、 比較可能群(= j)はNo. 1, 4, 5, 8, 10 (5組) , $C_{i=No.7} = \frac{5}{5} = 1$

平均 $\hat{C} \doteq 0.91$

AUC(C統計量), C-indexの弱点

予測確率がランダムに選ばれたもう1つのサンプルが大きいかどうかしか見ていない

→ どれくらい正確な予測ができているかどうかまでは分からない

→ 正しい予測ができているのかをCalibrationしてあげる必要がある

おまけ

1. 予測確率が同じ p のデータを集めると確率 p でイベントを生じるように
「同じ予測値を与える状況ではその予測値は**較正(calibration)**されるべき」
2. 不確実な個々の観測には確率 1 か 0 を個々に割り当てる予測より, $0 < p_i < 1$
という確率 p_i で予測する方が良い
3. このような確率的な予測 p_i においては「なるべく p_i の分布が極端になるように(0か1に近くなるように)予測確率を割り当てるべき」

Brier Score

$$BS = \frac{1}{N} \sum_{t=1}^N (f_t - o_t)^2$$