

二値データのロジスティック回帰分析にもとづく予測と予測性評価指標。

薬物治療の際には、予後日数や副作用発生の有無などのアウトカムを患者ごとに予測することが有用である。予測のためには何らかの経験的な規則や数式、あるいは新たに作成した数式を用いて新規患者でのアウトカムを求める。予測式を作成するために用いるデータを「トレーニングデータ」と呼び、予測性の評価に用いるデータを「テストデータ（バリデーションデータ）」と呼ぶ。

例えば 150 人の患者から得たデータを 100 人分のトレーニングデータと 50 人分のバリデーションデータに分け、まずトレーニングデータについて予測因子（X：複数個あってもよい）とアウトカム（Y）との関係を数式で導く。アウトカムが副作用の有無といった二値の場合にはロジスティック（重）回帰分析が適している。次にその予測式のうちの X の値にバリデーションデータの X 値を代入し、その患者での Y の予測値を得てバリデーションデータの実測値と比較することで予測性を評価できる。二値データの場合、ロジスティック回帰にもとづく Y の予測値は「副作用発生確率」として得られるので、例えばその確率が 0.5 以上であれば「副作用あり」、0.5 未満であれば「副作用なし」と換算する。

- ・（トレーニングデータの X、Y 値）→ 回帰式（予測式）
- ・（バリデーションデータの X 値）+ 予測式 → （バリデーションデータの Y 予測値（理論値））
- ・（バリデーションデータの Y 実測値） VS. （バリデーションデータの Y 理論値）（比較評価）

例えば、ある患者群の予後日数について 3 週間未満（<3w）か 3 週間以上（≧3w）かを予測することを考える。予後日数に影響する因子として「呼吸困難の有無（0：なし、1：あり）」、「せん妄の有無（0：なし、1：あり）」、「倦怠感の有無（0：なし、1：あり）」を考えて、トレーニングデータをロジスティック重回帰分析して次の式が得られたとする（架空データ）：

$$\begin{aligned} \text{Logit} = & -2.2 + 1.0 \times \text{呼吸困難の有無 (0 または 1)} \\ & + 1.8 \times \text{せん妄の有無 (0 または 1)} \\ & + 1.5 \times \text{倦怠感の有無 (0 または 1)} \end{aligned}$$

このとき、係数 -2.2 は切片（定数項）、1.0、1.8、1.5 が予測式の中での各症状の影響の程度を表す係数となる。また、オッズ比（OR）とこれらの係数（ β ）の間には $OR = \exp(\beta)$ の関係があるので、それぞれの症状についてのオッズ比は $\exp(1.0) = 2.8$ 、 $\exp(1.8) = 5.9$ 、 $\exp(1.5) = 4.5$ などと得られる。

次にバリデーションデータに含まれるあるひとりの患者で、呼吸困難はなかったがせん妄と倦怠感がみられたとすると、予測式から得られるこの患者での logit は $-2.2 + 1.0 \times 0 + 1.8 \times 1 + 1.5 \times 1 = 1.1$ となり $\text{logit} = \ln(p / (1 - p))$ の関係から $p = 0.750$ となる。別の患者で呼吸困難のみがみられた場合には $\text{logit} = -1.2$ 、 $p = 0.231$ となる。 $p = 0.5$ を境界値として予後日数の予測（<3w か ≧3w か）を二値データで行うとしたら、前者では $p \geq 0.5$ なので <3w、後者では $p < 0.5$ なので ≧3w と予測することになる。

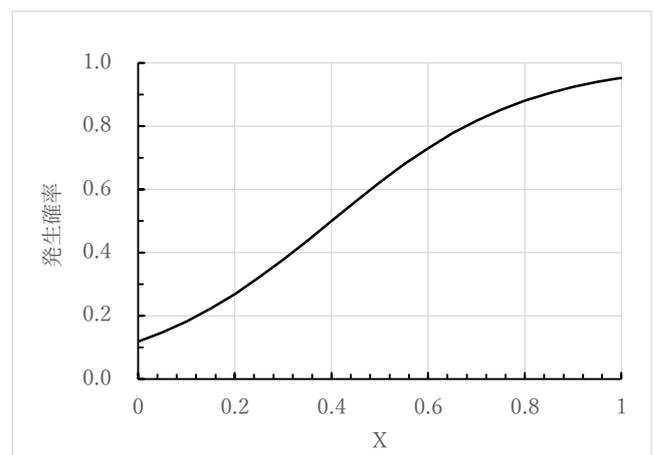
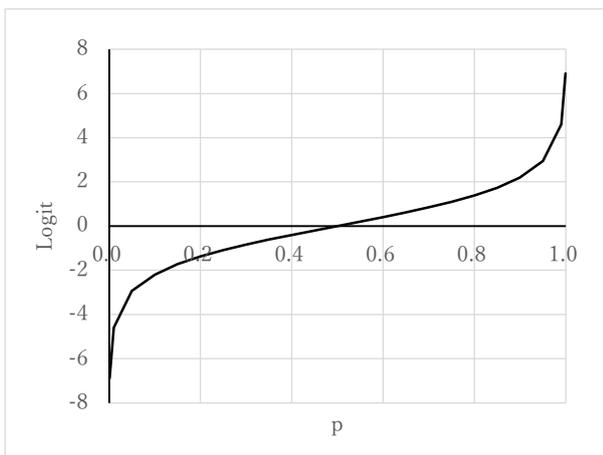
バリデーションデータの理論値と実測値を比較するとき、連続的な変数の場合には理論値と実測値との相関性を評価する（末尾【補足】参照）が、二値データの予測性評価には分割表を利用できる。

あるアウトプット（副作用の有無など）が「ある,Z」あるいは「ない,non-Z」と判定される場合、実測値（事実）との組み合わせで表のように4通りが考えられる。ここで「Zであるときに正しくZと予測する割合」を感度（Sensitivity）といい $A / (A + C)$ で計算できる。また「Zでないときに正しく non-Z と予測する割合」特異度（Specificity）といい $D / (B + D)$ で計算できる。その他に、「Zであると判断された被験者が実施にZであった割合を陽性的中率（Positive Predictive Value, PPV）といい $A / (A + B)$ で、「Zでないと判断された被験者が実際にZでなかった割合を陰性的中率（Negative Predictive Value, NPV）といい $D / (C + D)$ 」で得られる。表の値を例にすると、

予測値（理論値）	真実（実測値）		
	副作用あり	副作用なし	合計
副作用あり	A	B	A + B
副作用なし	C	D	C + D
合計	A + C	B + D	A+B+C+D

予測値（理論値）	真実（実測値）		
	副作用あり	副作用なし	合計
副作用あり	34	9	43
副作用なし	6	46	52
合計	40	55	95

それぞれ順に 0.85、0.84、0.79、0.88 となる。Logit と p の間には左図のような関係が、また、ある X と確率 p には右図のような関係がある（グラフと表や本文の数値とは一致させていない）。



【補足】連続変数の予測性評価指標

- ME（Mean Error、平均誤差）～偏り（bias）の指標
各ポイントの予測値と実測値の差を合計し、それを個数で割った値
- MAE（Mean Absolute Error、平均絶対誤差）～全体のばらつき（予測精度）の指標
各ポイントの予測値と実測値の差『の絶対値』を合計し、それを個数で割った値
- RMSE（Root Mean Squared Error、平均平方和の平方根）～全体のばらつき（予測精度）の指標
各ポイントの予測値と実測値の差を『二乗して』合計し、それを個数で割った値について『さらに平方根をとった』値

（この資料を作成するにあたっては T. Morita et al., *Jpn. J. Clin. Oncol.*, **29**(3), 156-159 (1999) を参考にした）