

海洋生態学のための関数データ解析

江口 真透 数理推論研究系

紹介： 生態学の主目的の一つはある地域の生態システムの健全さの評価を行うことにある。その地域と対象となる生物たちに着目して定期的に観測が継続された結果、関数 $\{X(t): t \in T\}$ が得られる。その生態系が健全であるかないかで2値の確率変数 Y が0または1を取るとしよう。このようにして、 n 組のデータ $(\{x_i(t): t \in T_i\}, y_i)$ $i = 1, \dots, n$ が得られる。このデータ解析のために関数一般化線形モデルなどについて考察したい。特に線形判別解析とロジスティック回帰モデルと混合正規モデルの関係を関数データの文脈で拡張を試みたい。水産資源の内容でおおよそ200系群を超える約55年間のデータが蓄えられている公開データベース ‘RAM Legacy Stock Assessment’の全域解析を試みる。

関数ロジスティックモデル： 一般化線形モデルの内容で線形予測量が

$$\eta = \alpha + \int_T \beta(t)X(t) dt$$

と表されるとし、リンク関数 g と分散関数 σ^2 によって次のように表されると仮定する：

$$\mathbb{E}(Y | X(t), t \in T) = g(\eta), \quad \text{Var}(Y | X(t), t \in T) = \sigma^2(\eta)$$

さらに基底関数 $\{\phi_j(t): j = 1, \dots, p\}$ に対してパラメトリック形： $\beta(t) = \beta_1\phi_1(t) + \dots + \beta_J\phi_J(t) = \beta^T \varphi(t)$ を仮定する。ロジスティック回帰モデルは

$$\mathbb{E}(Y = y | X(t), t \in T) = \frac{\exp\{y(\alpha + \beta^T \int_T \varphi(t)X(t) dt)\}}{1 + \exp(\alpha + \beta^T \int_T \varphi(t)X(t) dt)}$$

となる。ここでは共変量が単一の関数の場合を扱っているが、 J 本の関数 $\{X_j(t): j = 1, \dots, J, t \in T\}$ の場合や通常の共変量ベクトル x が追加されても線形予測量は同様に記述される：

$$\eta = \alpha + \beta_0^T x + \beta_1^T \int_T \varphi_1(t)X_1(t) dt + \dots + \beta_J^T \int_T \varphi_J(t)X_J(t) dt$$

関数線形判別分析： 線形判別分析の内容で $Y = y$ が与えられた $\{X(t)\}_{t \in T}$ の条件付分布を正規分布

$$\{X(t)\}_{t \in T} | Y = y \sim N(\{\mu_y^T \varphi(t)\}_{t \in T}, \Omega + \sigma^2 \mathbf{I})$$

を仮定すると線形判別関数は

$$\eta = \alpha + \beta^T \int_T \varphi(t)X(t) dt$$

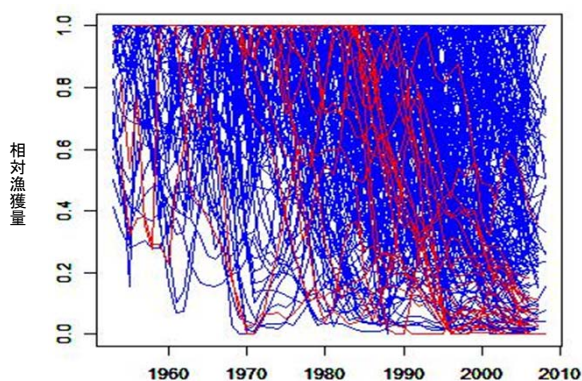
関数クラスタリング： モデルによるクラスタリングは正規混合モデル

$$\{X(t)\}_{t \in T} \sim \pi_1 N(\{\mu_1^T \varphi(t)\}_{t \in T}, \Omega + \sigma^2 \mathbf{I}) + \pi_0 N(\{\mu_0^T \varphi(t)\}_{t \in T}, \Omega + \sigma^2 \mathbf{I})$$

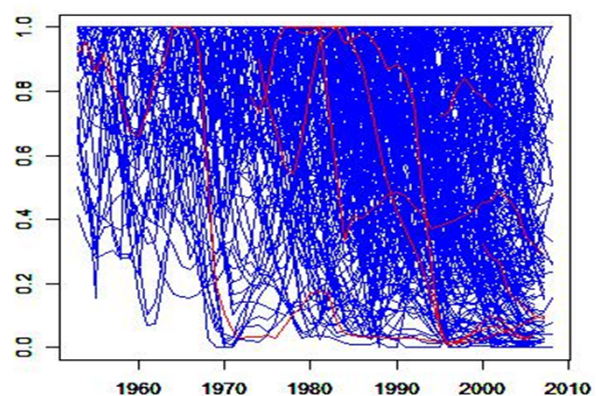
を仮定し、事後分布の推定によって実行される。

3つの方法の関連： 上記の3つのモデルの尤度解析は対数尤度関数の意味で線形判別解析はロジスティック回帰と正規混合モデルに分解される。正確には $L_J(\theta) = L_C(\beta) + L_M(\theta)$ の関係で特徴付けられる。

RAM Legacy Stock Assessment： RAMの232ストックの1953-2008年の相対漁獲データ (t 年の相対漁獲量 $R(t) = \frac{X(t)}{\max\{X(s): s \leq t\}}$)



真のラベル
(赤: 枯渇 青: 非枯渇)



FGLMの予測ラベル, AUC = 0.81