

レポート「信用調査データを用いた雇用傾向の把握」で 用いられているモデリング

中谷 太洋*

2023年9月20日

本資料は『信用調査データを用いた雇用傾向の把握 (2022年12月データ)』(<https://www.tdb-di.com/posts/2023/08/e2023083001.php>) で公開されたレポートのガウス過程回帰による欠測値の補完方法について数理的な背景を詳細にまとめたものである。

一般に統計や機械学習モデルでは、観測値の組 $(\mathbf{x}_i, \mathbf{y}_i)$ ($i = 1, \dots, n$) について、 \mathbf{x}_i を入力として \mathbf{y}_i を出力する関数 $f(\mathbf{x})$ を推定する。現実的には観測値は誤差を伴うため、 $\mathbf{y} = f(\mathbf{x}) + \varepsilon$ ($\varepsilon \sim N(0, \sigma^2)$) として推定を行う。ガウス過程回帰はこの $f(\mathbf{x})$ を非線形かつベイズ的に推定する。

まずガウス過程の1つの仮定として、 $f(\mathbf{X}) \sim N(m(\mathbf{X}), V(\mathbf{X}, \mathbf{X}))$ とする。これは、モデルが次元の入力に対して次元の出力を行い、誤差がない場合の出力が平均 $m(\mathbf{X}_i)$ 、分散 $V(\mathbf{X}_i, \mathbf{X}_i)$ の正規分布にしたがうことを表している。これに伴って観測値全体の同時分布が多変量正規分布で表されることになる。この仮定は便利で、モデルの入力の次元が増えても平均行列と2点間の関係の記述だけで観測値の同時分布が求まる。よって、連続的な無限次元の入力でも機能する。ガウス過程回帰の非線形かつベイズ的に関数を推定するという特徴は、以上の仮定とその推定方法によるものである。先の記述から、 $f(\mathbf{X}) \sim N(m(\mathbf{X}_i), V(\mathbf{X}_i, \mathbf{X}_i))$ として事前分布からランダムに乱数を生成する。ここでは乱数を関数値とみなすことで、関数 f そのものを抽出していることになる。乱数の組をいくつも生成すれば、異なる関数が得られる。ここで $V(\mathbf{X}_i, \mathbf{X}_i)$ の要素について、ガウス過程回帰ではカーネル関数 (共分散関数) $k(\mathbf{x}_i, \mathbf{x}_i)$ を用いる。カーネル関数には、動径基底関数 (RBF カーネル、ガウスカーネル)、線形カーネル、Matern カーネルなど様々な種類があり、これらの選択やその組み合わせによって直線や曲線の表現が可能になる。ここからはカーネル行列を $K(\mathbf{x}_i, \mathbf{x}_i)$ と表す。

次に事後分布について考える。ガウス過程回帰では、関数を乱数列としているため、明示的なパラメータがない。つまり、通常の統計モデルや機械学習モデルのようにパラメータと入力値の和積でその入力に対する出力が得られない。しかし、ガウス過程の仮定から新しい入力点に対する出力を求めることができる。新しい入力ベクトルを \mathbf{x}^* 、対応する出力を $f(\mathbf{x}^*)$ とすると、出力の同時分布は、

$$\begin{pmatrix} \mathbf{y} \\ f(\mathbf{x}^*) \end{pmatrix} \sim N \left[\begin{pmatrix} m(\mathbf{X}_i) \\ m(\mathbf{X}^*) \end{pmatrix}, \begin{pmatrix} K(\mathbf{X}_i, \mathbf{X}_i) + \sigma^2 \mathbf{I} & K(\mathbf{X}_i, \mathbf{X}^*) \\ K(\mathbf{X}_i, \mathbf{X}^*)^T & K(\mathbf{X}^*, \mathbf{X}^*) \end{pmatrix} \right]$$

として表すことができる。分散共分散行列のカーネル関数では、既存の入力と新しい入力と関心がある $f(\mathbf{x}^*)$ の分布は、 \mathbf{y} が与えられたもとの条件付き分布となるため、

$$f(\mathbf{X}^*) | \mathbf{y} \sim N(m(\mathbf{X}^*) + K(\mathbf{X}_i, \mathbf{X}^*)^T (K(\mathbf{X}_i, \mathbf{X}_i) + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - m(\mathbf{X}_i)), \\ K(\mathbf{X}^*, \mathbf{X}^*) - K(\mathbf{X}_i, \mathbf{X}^*)^T (K(\mathbf{X}_i, \mathbf{X}_i) + \sigma^2 \mathbf{I})^{-1} K(\mathbf{X}_i, \mathbf{X}^*))$$

* 博士前期課程2年

と表せる。計算を簡便化させるために予め観測から平均を引いて中心化することが多い。実際のモデルは関数の出力に誤差を加えたものであるため、予測分布は、

$$\mathbf{y}^* | \mathbf{y} \sim N(m(\mathbf{X}^*) + K(\mathbf{X}_i, \mathbf{X}^*)^T (K(\mathbf{X}_i, \mathbf{X}_i) + \sigma^2 \mathbf{I})^{-1} (\mathbf{y} - m(\mathbf{X}_i)), \\ K(\mathbf{X}^*, \mathbf{X}^*) - K(\mathbf{X}_i, \mathbf{X}^*)^T (K(\mathbf{X}_i, \mathbf{X}_i) + \sigma^2 \mathbf{I})^{-1} K(\mathbf{X}_i, \mathbf{X}^*) + \sigma^2)$$

となる。本レポートの分析では、各企業に対してガウス過程回帰を適用し、各四半期の中心となる月の値を取得している。つまり、1年前までの直近10年で毎年調査が入っている企業を対象にしていることから最低でも10点の入出力の組をもってガウス過程回帰のモデルを構築し、毎年2、5、8、11月の値を取得している。値の取得には不確実性を考慮して、平均曲線上の値とその±1標準偏差の値を用いている。最終的に、平均、+1標準偏差、-1標準偏差のそれぞれで指標を算出する。全体的に LOCF よりもなだらかな結果となった。

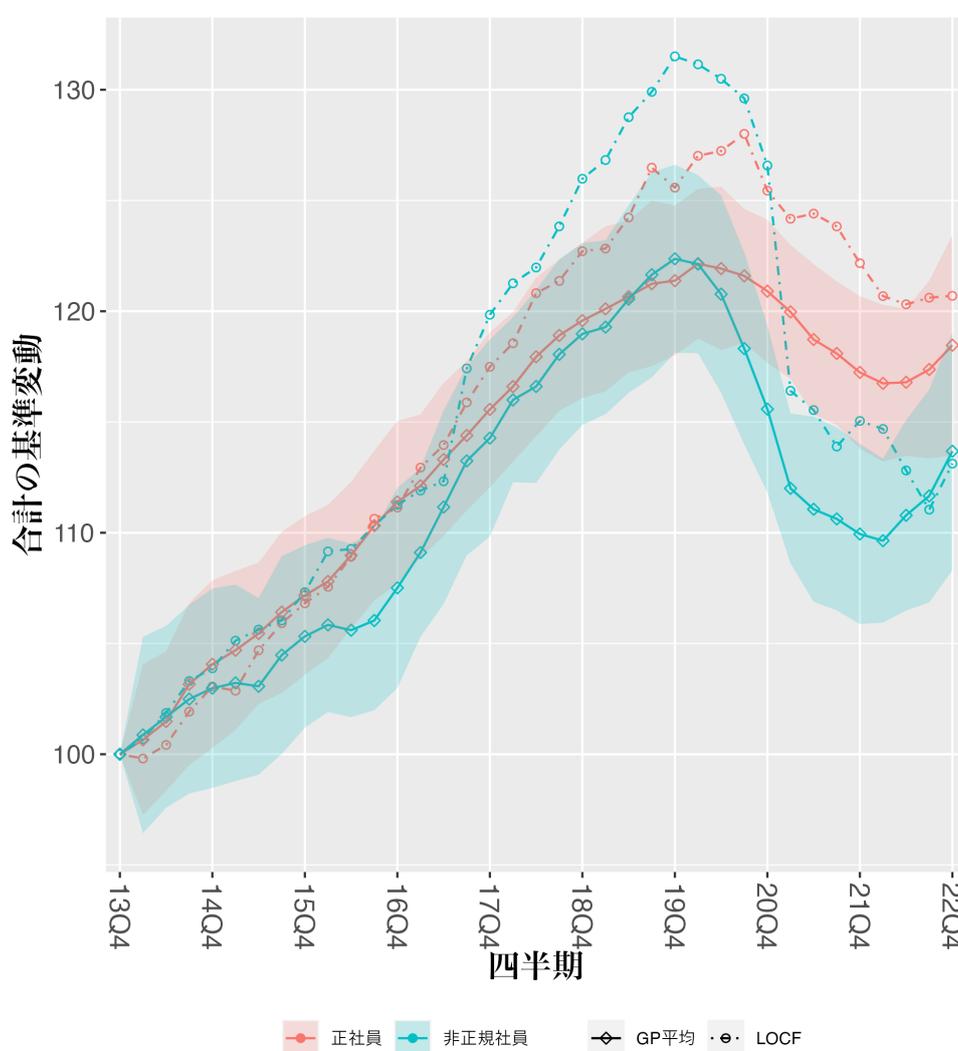


図1 2022年第4四半期の合計の基準変動