

3章冒頭

3.1節、3.2節

前処理

データマイニングの分析目的


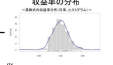
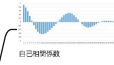

- 学習型分析 → 既知のデータから学習モデルを作り、未知のデータの特徴を明らかにする
- 記述型分析 → データに対して、特徴や相関関係の相関性を抽出する

3章では予測的解析における後処理を扱う

1. データクレンジング（整理、標準化）【3.1節】

- ① 目的変数データの設定
  - やることは4つ
    - (1) データが反映している情報 → 分析目的に合致したデータを予測対象として選択する
    - (2) データの入事可能性を確認 → 時間間隔が短いものや、詳細な情報は必要の場合が多い
    - (3) データの範囲・計算方法を事前に確認する
      - 定額市況データの例 → 株式指数：組み入れ銘柄の増減、計算時の重み
      - 定額市況データの例 → 定額指数：株式分割や合併などによる調整の範囲
      - 定額市況データの例 → 定額指数：株式分割や合併などによる調整の範囲
    - (4) 学習期間の長さを決める
      - 決めるときに主に考慮する点
        - データの期間
        - 学習期間
- ② サンプリング間隔の決定
  - 目的は、不明な時間系列を等間隔時間系列にすること → 例えば、ティックデータ等は等間隔ではない
  - やり方は主に2つ
    - ・ダウンサンプリング（データの間引き） → 十分に高頻度の場合はこっち
    - ・内挿（補間）
  - 注意点
    - もともと持っているデータの特性が失われないようにする → よくあるのが、周知性
    - グラフや基本統計量を用いてから調整しよう
- ③ 欠損値と外れ値の処理
  - 欠損値
    - 時系列の場合、欠損は補間する
    - ただし、補間が明らかに概念は理解しないこともある → 例えば、週末であるとき
  - 外れ値
    - 理論上ありえない値の中には欠損値と同じ扱い
    - 分析目的による

2. 特徴量抽出【3.2節】

- ① データの値を決めている要因について仮説を立てる
  - 【Point】金融・経済の時系列データの値を決める要因
    - 時系列データの値 = 自己相関 + 長期トレンド + 短期変動 + 外生変数 + ノイズ
  - 自己相関 → 過去の自分の値が将来の値に与える影響
  - 長期トレンド → 長い期間に存在する緩急した変動の傾向
  - 短期変動 → ある決まった期間で繰り返して発生する変動の傾向
  - 外生変数 → 予測対象とは別のメカニズムで生じた事象（外生変数）が原因となって引き起こされる変動
  - ノイズ → 特定の傾向や特徴を含まないランダムな変動
- ② 可視化 → 最終結果でよくグラフ
  - 時系列プロット
    - 例 
  - 変動性のヒストグラム
    - 例 
  - コロログラム
    - 異なるラグに対して自己相関係数もしくは交差相関係数を算出し、縦軸にラグ、横軸に自己相関係数をとったグラフ
    - 自己相関係数 
    - 交差相関係数 
- ③ 定常性の確認
  - 定常性の仮説 → この資料では省略
  - 金融時系列データが定常性をもつのは、通常定常な仮説であるため確認する。
  - 確認方法 → 検定法では扱っていないが
- ④ 見せかけの相関の確認
  - 単回帰分析とは → 単回帰分析であり、部分相関係数であるか確認 → 例：ランダムウォーク
  - 見せかけの相関とは → 単回帰分析である場合、関係のない単回帰分析結果を導く → 例：ランダムウォーク
  - 単回帰分析 → 単回帰分析だが、回帰係数にのみ注目し、関係性を確認する
- ⑤ 変数変換
  - 変換後のデータでも定常性がないか確認する
    - 定常性があるから、自己相関による変動要因を排除した残差に含まれる外生変数
    - 市には長期トレンド要因、短期変動も
  - 差分変換 → 階層分だけでなく、2階層分も差分を取るとなる
  - 対数変換 → 大きく変化した金額が、データの特徴を減らすに均一な変換に変換できる
  - 対数変換の変換 → 対数値の差分
  - 対数変換の変換 → 金融データでは、収益率の代わりに、対数リターン（対数値の差分）を用いることがある
  - ベンチマーク調整 → 市場の平均利回りからの超過リターンなど
  - 正の変数変換のやり方は以下の通り
    - ロジット変換 →  $Z = \log\left(\frac{Y}{1-Y}\right)$  → 0以上の範囲のデータが、(0, 1)の範囲に変換される
    - 乗数調整 → 例えば、月別平均値
    - 乗数調整 → 例えば、パーセント法
    - 乗数平均