

陽性患者数の予測

ARIMA モデルと局所線形構造モデル

N. Murata

2021年1月8日

1 はじめに

本稿では、時系列の基本的なモデルを用いて今後の患者数の予測を試みる。

図示やモデル化に際しては以下の package を用いる。

```
1  ## パッケージの読み込み
2  library(forecast)
3  library(tidyverse)
4  library(scales) # 年月日表示
5  library(plotly)
6  library(zoo)    # 時系列表示
7  library(ggfortify)
```

2 データの取得

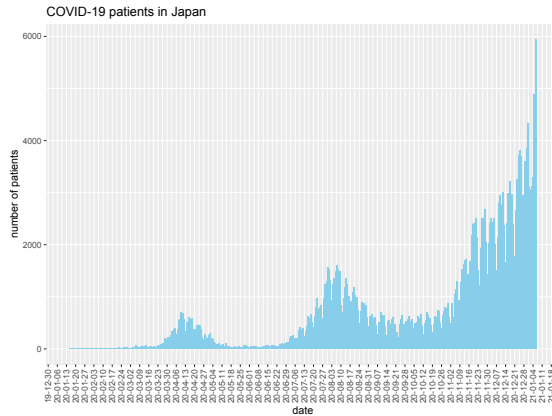
以下では厚生労働省が公開している COVID-19 の全国の感染者数データを利用する。¹

¹ 厚生労働省の患者数データ

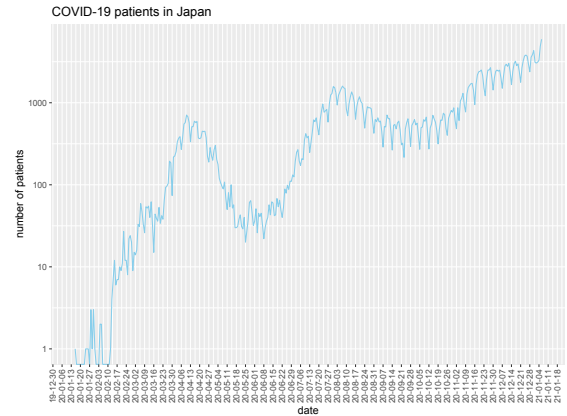
```
1  ## データの取得と整理
2  patients <-
3    read.csv("https://www.mhlw.go.jp/content/pcr_positive_d_
4      ↪ aily.csv")
5      ↪ %>%
6      dplyr::rename(date=1, patients=2) %>%
7      dplyr::mutate(date=as.Date(date))
8  ## 時系列データ (zooクラス) への変更
9  patients <- with(patients,
10    zoo(x=patients, order.by=date))
```

図1は全国の陽性患者数の推移を図示したものである。左図は横軸に日付を、縦軸に観測された患者数を表示したものである。SIRなどの感染症の基本的なモデルによれば、感染が拡大しはじめる際の人数の増加は指数関数的であるので、右図では縦軸を患者数の常用対数として表示している。週日・休日での検査数に波はあるが、9月以降はほぼ単調に増加していることがわかる。

```
1  ## データの視覚化
2  p <-
3    ggplot(data = fortify(patients, melt = TRUE),
4      mapping = aes(x = Index,
5        y = Value)) +
6    scale_x_date(labels = date_format("%y-%m-%d"), # 年月日表示
7      breaks = date_breaks("1 week")) + # 週毎
8    theme(axis.text.x = element_text(angle = 90,
9      vjust = 0.5, hjust=1)) +
```



(a) 日毎の患者数の変遷



(b) 患者数の対数表示

図 1: 全国の陽性患者数の推移.

```

10     labs(title = "COVID-19 patients in Japan",
11           x = "date",
12           y = "number of patients")
13   ## 棒グラフ
14   print(p + geom_col(fill="skyblue")) # グラフ出力
15   ggplotly() # plotly表示 (browser)
16   ## 折れ線グラフ+常用対数表示
17   print(p + geom_line(colour="skyblue") + scale_y_log10())
18   ggplotly()

```

3 基礎分析

次節以降で確認するモデルの精度を確認するために、推定と予測のために系列を分けて考える。² 9月中旬から11月のデータをモデルの推定に用い、12月以降のデータで予測の精度を検証することとする。

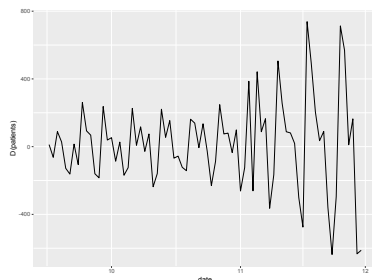
² 東京都の分析結果から9月中旬あたりが第3波の開始と考えらるので、ここでは9月15日以降のデータを用いることとする。

```

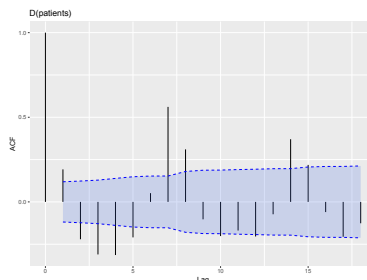
1   ## 9月以降の第3波を対象とする
2   train <- window(patients,
3                   start="2020-09-15",
4                   end="2020-11-30")
5   test <- window(patients,
6                  start="2020-12-01")

```

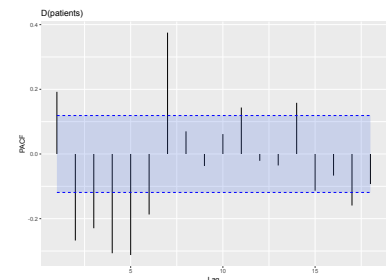
まず、9-11月のデータを用いて時系列の基礎的な性質を確認する。データは明らかに定常ではないので、まず階差系列の性質を確認する。図3に階差系列の推移、および自己相関と偏自己相関を示す。系列の分散は時間とともに増大しており、また相関関係から7日毎の関係が強いことがわかる。これは検査機関の稼働状況の影響だと考えられる。



(a) 階差系列



(b) 自己相関



(c) 偏自己相関

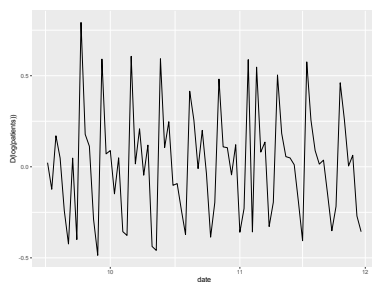
図 3: 階差系列の基礎的な性質.

```

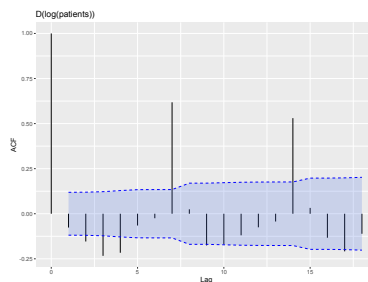
1  ## 階差系列の性質
2  autoplot(diff(train)) +
3    labs(x = "date",
4         y = "D(patients)")
5  autoplot(acf(diff(train), plot = FALSE), # 自己相関
6           conf.int.fill = "royalblue",
7           conf.int.alpha = 0.2,
8           conf.int.value = 0.7,
9           conf.int.type = "ma") +
10   labs(title = "D(patients)")
11  autoplot(pacf(diff(train), plot = FALSE), # 偏自己相関
12           conf.int.fill = "royalblue",
13           conf.int.alpha = 0.2,
14           conf.int.value = 0.7) +
15   labs(title = "D(patients)",
16        y = "PACF")

```

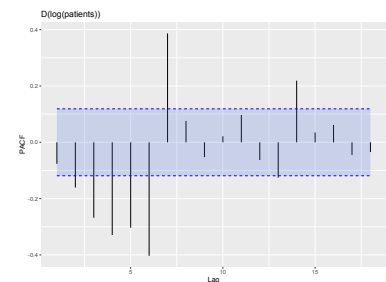
SIR のような感染症のモデルの解析から、感染者数は指数型の増加を示すことが知られているので、次に感染者数の対数を取ったものの階差系列の性質を調べる。この結果を図 5 に示す。系列の推移より、分散の変動が定常化されていることがわかる。自己相関からは 7 日毎の関係が同様に確認される。



(a) 階差系列



(b) 自己相関



(c) 偏自己相関

図 5: 対数変換後の階差系列の性質。

```

1  ## 対数変換を確認する
2  ltrain <- log(train)
3  autoplot(diff(ltrain)) +
4    labs(x = "date",
5         y = "D(log(patients))")
6  autoplot(acf(diff(ltrain), plot = FALSE), # 自己相関
7           conf.int.fill = "royalblue",
8           conf.int.alpha = 0.2,
9           conf.int.value = 0.7,
10          conf.int.type = "ma") +
11   labs(title = "D(log(patients))")
12  autoplot(pacf(diff(ltrain), plot = FALSE), # 偏自己相関
13           conf.int.fill = "royalblue",
14           conf.int.alpha = 0.2,
15           conf.int.value = 0.7) +
16   labs(title = "D(log(patients))",
17        y = "PACF")

```

また、この階差系列の基本統計量 (平均と中央値) を見ると、若干正に偏っており、元の系列が増大する傾向があることが確認される。

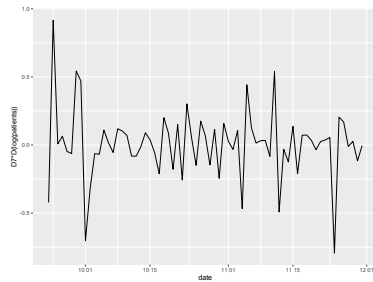
```

1  ## 基本統計量の確認
2  summary(as.numeric(diff(ltrain)))

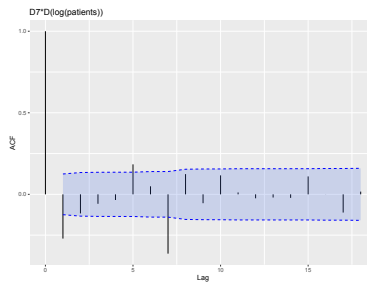
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-0.48603	-0.23086	0.02945	0.01303	0.17210	0.79224

次に7日毎の階差を取り、7日毎の相関の性質がどのように変化するか確認する。これを図7に示す。相関関係の正負が逆転し、強い相関が残ることから、単純な周期成分ではないことが示唆される。³

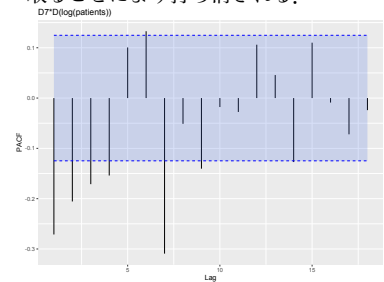


(a) 階差系列



(b) 自己相関

³ 7日毎の周期成分であれば7日階差を取るにより打ち消される。



(c) 偏自己相関

図7: 1日および7日の階差を取った系列の性質。

```
1  ## 7日の周期性を確認する
2  autoplot(diff(diff(ltrain), lag=7)) +
3      labs(x = "date",
4           y = "D7*D(log(patients))")
5  autoplot(acf(diff(diff(ltrain), lag=7), plot = FALSE), # 自己相関
6           conf.int.fill = "royalblue",
7           conf.int.alpha = 0.2,
8           conf.int.value = 0.7,
9           conf.int.type = "ma") +
10     labs(title = "D7*D(log(patients))")
11 autoplot(pacf(diff(diff(ltrain), lag=7), plot = FALSE), # 偏自己相関
12           conf.int.fill = "royalblue",
13           conf.int.alpha = 0.2,
14           conf.int.value = 0.7) +
15     labs(title = "D7*D(log(patients))",
16          y = "PACF")
```

4 ARIMA モデル

まず ARIMA モデルによる推定を行う。ARIMA モデルは階差系列に対して ARMA モデルを用いたものである。直感的には ARMA モデルに従う時系列のランダムウォークであり、非定常な時系列の基本的なモデルの一つである。

時系列を X_t , $t = 1, 2, \dots$ とし、ラグ作用素 L (lag operator, backshift operator) を以下で定義する。

$$LX_t = X_{t-1}$$

これを用いると、階差系列は

$$Y_t = X_t - X_{t-1} = (1 - L)X_t$$

と書け、高階の階差は作用素 $(1 - L)$ の冪で書くことができる。また次数 (p, d, q) の ARIMA モデルは以下で定義される。

$$(1 - a_1L - a_2L^2 - \dots - a_pL^p)(1 - L)^dX_t = (1 + b_1L + b_2L^2 + \dots + b_qL^q)\epsilon_t, \epsilon_t \sim \mathcal{N}(0, \sigma^2)$$

今回対象とするデータは1階の階差で定常となっていると考えられるため、 $d = 1$ のモデルを考えることになる。また、増大する傾向を一種のトレンドとしてモデル化するためにドリフト (drift) 項を加えて以下のモデルを考える。

$$(1 - a_1L - a_2L^2 - \dots - a_pL^p)(1 - L)(X_t + \beta t) = (1 + b_1L + b_2L^2 + \dots + b_qL^q)\epsilon_t, \epsilon_t \sim \mathcal{N}(0, \sigma^2)$$

次数 (p, d, q) と対応する係数とドリフト, およびホワイトノイズの分散 σ^2 を推定した結果は以下のようになる。⁴

```
1  ## drift 付きの ARIMA モデルの次数を自動推定
2  est.arima <- forecast::auto.arima(ltrain)
3  ## 推定されたモデルを表示
4  print(est.arima)
5  ## SARIMA モデルを当て嵌める場合は周期を指定する.
6  ## frequency(ltrain) <- 7 # 7日周期の成分を仮定
7  ## est.arima <- auto.arima(ltrain)
8  ## このデータではモデルの推定はうまくいかない
```

Series: ltrain
ARIMA(5,1,1) with drift

Coefficients:

	ar1	ar2	ar3	ar4	ar5	ma1	drift
	-0.0394	-0.3175	-0.3379	-0.3265	-0.2757	-0.4083	0.0191
s.e.	0.1722	0.1102	0.1083	0.1120	0.1263	0.1532	0.0075

sigma² estimated as 0.06513: log likelihood=-1.28
AIC=18.55 AICc=20.7 BIC=37.2

図9は推定されたモデルの良さを診断するためのプロットで, 上から順に残差の推移, 残差の自己相関, 自己相関のLjung-Box 統計量(自己相関が0であるか検定)を表している. 残差は無作為になっているように見えるが, 7日毎の自己相関は若干残っており, 7日以降のラグの無相関性は棄却されない.

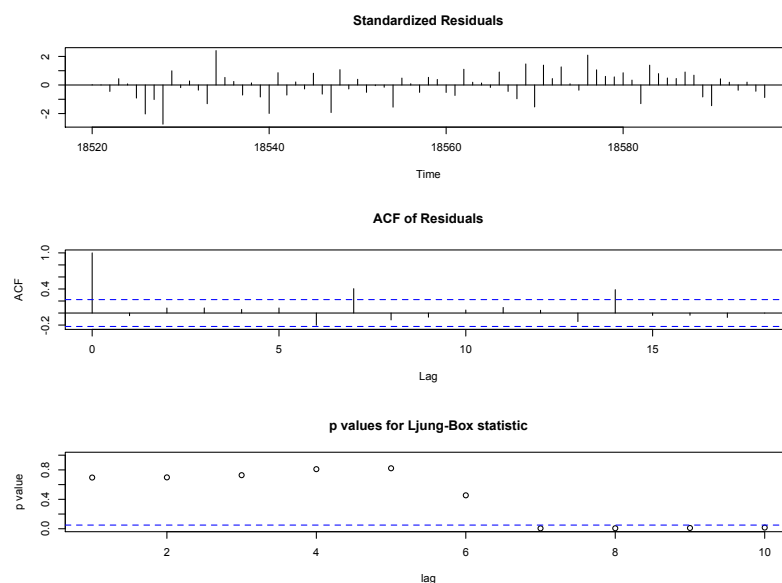


図9: ARIMA モデルの診断.

```
1  ## 診断プロット
2  tsdiag(est.arima)
3  ## 残差に相関が残っているので, 優れたモデルという訳ではない
```

推定されたモデルによる当て嵌め結果(モデルによる1期先の予測)を図10に示す. 概周期成分による誤差(遅れ)はあるものの, それなりに良く追従していることがわかる.

⁴ 季節成分 (seasonal) を加えた SARIMA モデルを考えることもできるが, このデータでは周期性が曖昧なため, モデルの推定はうまくいかなかった. コードのコメント部分を参照.

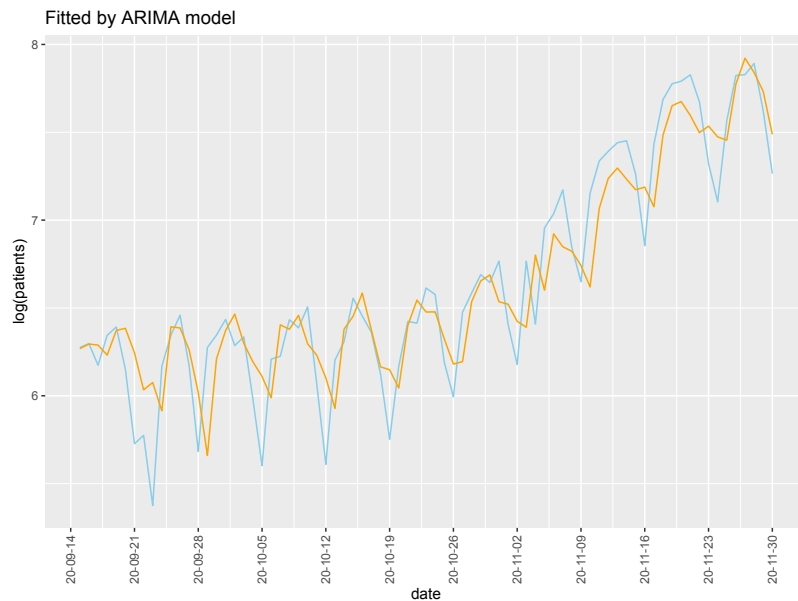


図 10: ARIMA モデルによる当て嵌め結果.

```

1  ## モデルによる当て嵌めの視覚化
2  p <-
3      ggplot(data = fortify(est.arima) %>%
4          dplyr::mutate(Index=as.Date(Index)),
5          mapping = aes(x = Index,
6                        y = Data)) +
7      geom_line(colour = "skyblue") +
8      geom_line(mapping = aes(y = Fitted),
9                  colour = "orange") +
10     scale_x_date(labels = date_format("%y-%m-%d"),
11                 breaks = date_breaks("1 week")) +
12     theme(axis.text.x = element_text(angle = 90,
13                                       vjust = 0.5, hjust=1)) +
14     labs(title = "Fitted by ARIMA model",
15          x = "date",
16          y = "log(patients)")
17 print(p)
18 ggplotly()

```

このモデルを用いて、12月以降の患者数を予測した結果を 80%信頼区間とともに示したのが図 11 である。

```

1  ## 12月以降 (最大 60 日) を予測してみる
2  p <-
3      ggplot(data = fortify(forecast(est.arima,
4                                    h=min(length(test),60))) %>%
5          dplyr::mutate(Index=as.Date(Index)) %>%
6          left_join(fortify(test), by = "Index"),
7          mapping = aes(x = Index,
8                        y = exp(Data)),
9          na.rm = TRUE) +
10     geom_line(colour = "skyblue",
11               na.rm = TRUE) +
12     geom_line(mapping = aes(y = test),
13               colour = "red",
14               na.rm = TRUE) +
15     geom_line(mapping = aes(y = exp(`Point Forecast`)),
16               colour = "royalblue",

```

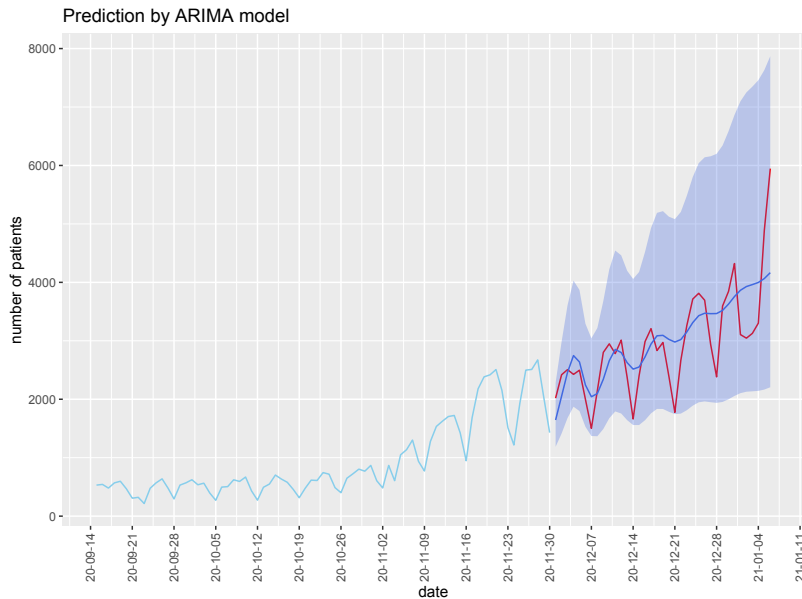


図 11: ARIMA モデルによる予測。

```

17     na.rm = TRUE) +
18     geom_ribbon(mapping = aes(ymin = exp(`Lo 80`),
19                               ymax = exp(`Hi 80`)),
20               fill = "royalblue", alpha = 0.3,
21               na.rm = TRUE) +
22     ## geom_ribbon(mapping = aes(ymin = exp(`Lo 95`),
23                               ##   ymax = exp(`Hi 95`)),
24     ##           fill = "royalblue", alpha = 0.1,
25     ##           na.rm = TRUE) +
26     scale_x_date(labels = date_format("%y-%m-%d"),
27                 breaks = date_breaks("1 week")) +
28     theme(axis.text.x = element_text(angle = 90,
29                                       vjust = 0.5, hjust=1)) +
30     labs(title = "Prediction by ARIMA model",
31          x = "date",
32          y = "number of patients")
33 print(p)
34 ggplotly()

```

5 局所線形構造モデル

もう一つは局所線形構造モデルと呼ばれるもので、時系列に以下の構造を仮定したものである。

観測された時系列 X_t は、トレンド μ_t とホワイトノイズ ϵ の2つの成分からなると仮定する。

$$X_t = \mu_t + \epsilon_t, \epsilon_t \sim \mathcal{N}(0, \sigma_\epsilon^2)$$

さらにトレンドは以下のような力学系に従うと仮定する。

$$\mu_{t+1} = \mu_t + \nu_t + \xi_t, \xi_t \sim \mathcal{N}(0, \sigma_\xi^2)$$

$$\nu_{t+1} = \nu_t + \zeta_t, \zeta_t \sim \mathcal{N}(0, \sigma_\zeta^2)$$

直感的には ν_t はトレンドの傾きに相当する。傾きがランダムウォークにより増減することにより、トレンドの変化の速度が変わるため、多項式より複雑な形状のトレンドを表すことができる。

このモデルには3つのホワイトノイズが含まれているが、その分散の最尤推定を行うことによって、モデルの推定が行われる。

```

1  ## StructTS による方法
2  est.sts <- StructTS(ltrain)
3  ## 推定されたモデルを表示
4  print(est.sts)

```

Call:

StructTS(x = ltrain)

Variances:

```

      level      slope  epsilon
0.07039  0.00000  0.01102

```

ARIMA モデルと同様に、推定されたモデルの診断プロットを図 12 に示す。ARIMA より残差に相関が残っていることが確認できる。

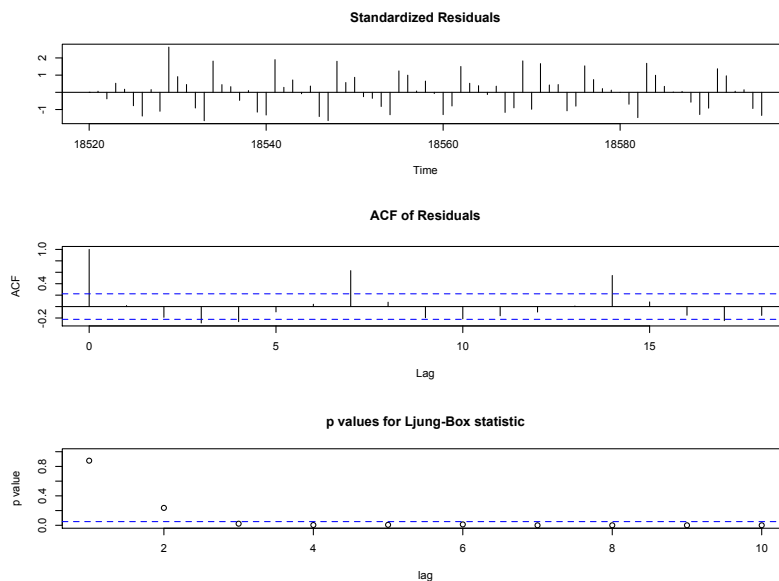


図 12: 局所線形構造モデルの診断。

```

1  ## 診断プロット
2  tsdiag(est.sts)
3  ## こちらも残差に相関が残っているので、優れたモデルという訳ではない

```

推定されたモデルによる元の時系列の分解結果は図 13 のようになる。

```

1  ## StructTS による時系列の分解結果
2  autoplot(est.sts)

```

モデルによる当て嵌め結果を図 14 に示す。ARIMA モデルより推定のラグが大きいことが見て取れる。

```

1  ## モデルによる当て嵌めの視覚化
2  p <-
3    ggplot(data = fortify(forecast(est.sts)) %>%
4      dplyr::mutate(Index=as.Date(Index)),
5      mapping = aes(x = Index,
6        y = Data),
7      na.rm = TRUE) +
8    geom_line(colour = "skyblue",
9      na.rm = TRUE) +
10    geom_line(mapping = aes(y = Fitted),
11      colour = "orange",
12      na.rm = TRUE) +

```

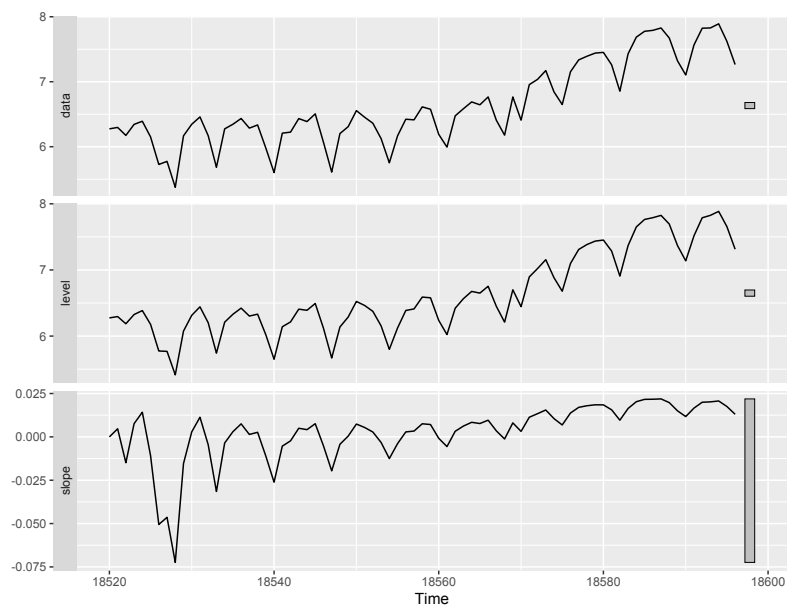



図 13: 局所線形構造モデルによる時系列の分解. level が μ , slope が ν に対応する.

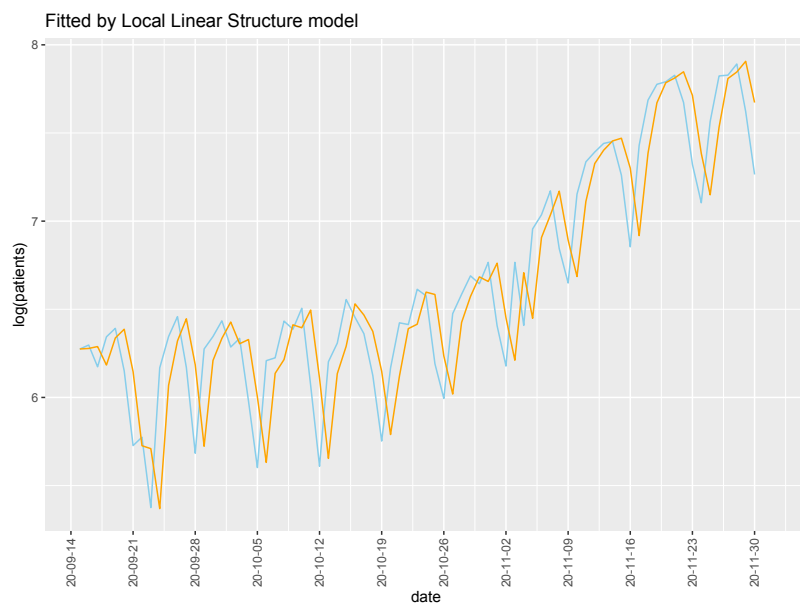


図 14: 局所線形構造モデルによる当て嵌め結果.

```

13     scale_x_date(labels = date_format("%y-%m-%d"),
14                 breaks = date_breaks("1 week")) +
15     theme(axis.text.x = element_text(angle = 90,
16                                       vjust = 0.5, hjust=1)) +
17     labs(title = "Fitted by Local Linear Structure model",
18          x = "date",
19          y = "log(patients)")
20     print(p)
21     ggplotly()

```

予測の結果は図 15 に示す。増加傾向は予測できており、信頼区間の中に実際の値は含まれているが、信頼区間は非常に広く、予測の精度は ARIMA に劣ると考えられる。

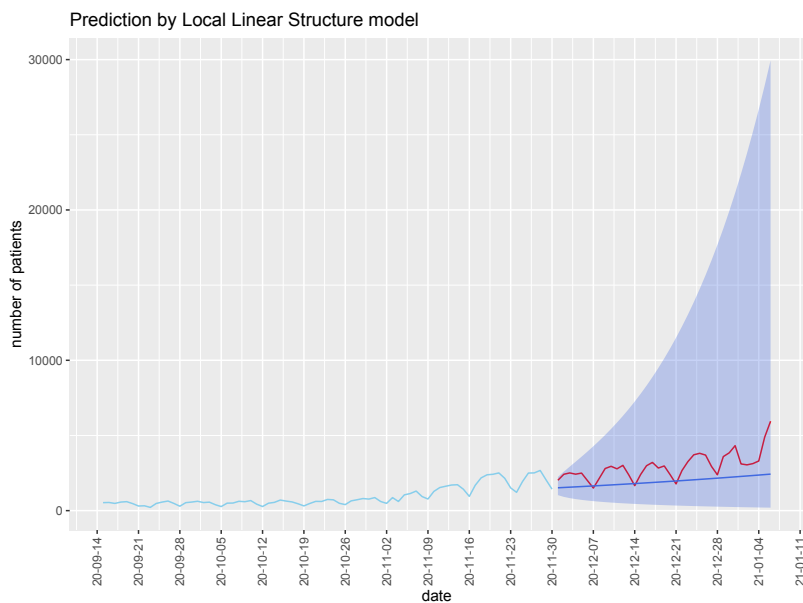


図 15: 局所線形構造モデルによる予測。

```

1  ## 12月以降 (最大 60 日) を予測してみる
2  p <-
3      ggplot(data = fortify(forecast(est.sts,
4                                  h=min(length(test),60))) %>%
5          dplyr::mutate(Index=as.Date(Index)) %>%
6          left_join(fortify(test), by = "Index"),
7          mapping = aes(x = Index,
8                        y = exp(Data)),
9          na.rm = TRUE) +
10     geom_line(colour = "skyblue",
11              na.rm = TRUE) +
12     geom_line(mapping = aes(y = test),
13              colour = "red",
14              na.rm = TRUE) +
15     geom_line(mapping = aes(y = exp(`Point Forecast`)),
16              colour = "royalblue",
17              na.rm = TRUE) +
18     geom_ribbon(mapping = aes(ymin = exp(`Lo 80`),
19                             ymax = exp(`Hi 80`)),
20               fill = "royalblue", alpha = 0.3,
21               na.rm = TRUE) +
22     ## geom_ribbon(mapping = aes(ymin = exp(`Lo 95`),
23     ##                               ymax = exp(`Hi 95`)),
24     ##                               fill = "royalblue", alpha = 0.1,

```

```

25     ##           na.rm = TRUE) +
26     scale_x_date(labels = date_format("%y-%m-%d"),
27                 breaks = date_breaks("1 week")) +
28     theme(axis.text.x = element_text(angle = 90,
29                                       vjust = 0.5, hjust=1)) +
30     labs(title = "Prediction by Local Linear Structure model",
31          x = "date",
32          y = "number of patients")
33     print(p)
34     ggplotly()

```

6 ARIMA モデルによる予測

前2節の結果から、このデータに対しては ARIMA モデルによる予測の方が精度が良いと考えられる。階差系列が ARMA モデルで良く近似されるということは、9月中旬以降の感染拡大の動特性があまり変化していないということであり、集団の行動がこの期間でほとんど変様していないことが示唆される。

現在までのデータを用いて ARIMA モデルを再度構築し、これを用いて 60 日間の予測を行った結果を図 16 に示す。内側は 80%信頼区間、外側は 95%信頼区間である。

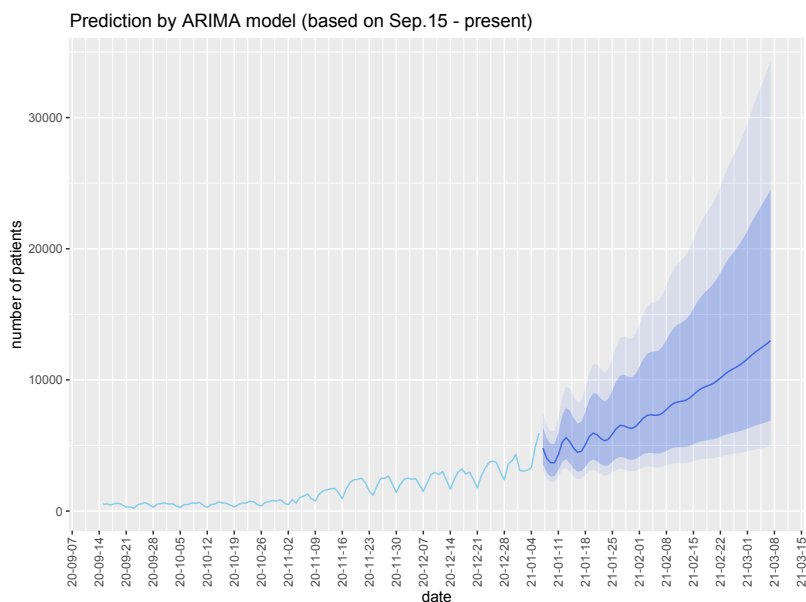


図 16: ARIMA モデルによる今後 60 日間の予測。

```

1  ## モデルの推定
2  est <- auto.arima(log(window(patients,start="2020-09-15")))
3  ## 推定されたモデルの表示
4  print(est)
5  ## 現在から 60 日先まで予測してみる
6  p <-
7      ggplot(data = fortify(forecast(est,h=60)) %>%
8              dplyr::mutate(Index=as.Date(Index)),
9              mapping = aes(x = Index,
10                           y = exp(Data)),
11              na.rm = TRUE) +
12      geom_line(colour = "skyblue",
13                na.rm = TRUE) +
14      geom_line(mapping = aes(y = exp(`Point Forecast`)),

```

```

15         colour = "royalblue",
16         na.rm = TRUE) +
17     geom_ribbon(mapping = aes(ymin = exp(`Lo 80`),
18                               ymax = exp(`Hi 80`)),
19               fill = "royalblue", alpha = 0.3,
20               na.rm = TRUE) +
21     geom_ribbon(mapping = aes(ymin = exp(`Lo 95`),
22                               ymax = exp(`Hi 95`)),
23               fill = "royalblue", alpha = 0.1,
24               na.rm = TRUE) +
25     scale_x_date(labels = date_format("%y-%m-%d"),
26                 breaks = date_breaks("1 week")) +
27     theme(axis.text.x = element_text(angle = 90,
28                                       vjust = 0.5, hjust=1)) +
29     labs(title = "Prediction by ARIMA model (based on Sep.15 - present)",
30          x = "date",
31          y = "number of patients")
32 print(p)
33 ggplotly()

```

```

Series: log(window(patients, start = "2020-09-15"))
ARIMA(5,1,1) with drift

```

Coefficients:

	ar1	ar2	ar3	ar4	ar5	ma1	drift
	-0.0594	-0.3813	-0.3579	-0.3409	-0.3318	-0.4022	0.0191
s.e.	0.1324	0.0869	0.0891	0.0886	0.1006	0.1193	0.0051

```

sigma^2 estimated as 0.04994: log likelihood=11.54
AIC=-7.07 AICc=-5.69 BIC=14.74

```

7 おわりに

対象とした患者数の推移は極めて非定常なデータであるが、時系列を適切に変換し、生成モデルがある程度定常な区間を捉えることができれば、基本的な ARMA モデルでも良い推定を行うことができる。