

به نام خدا

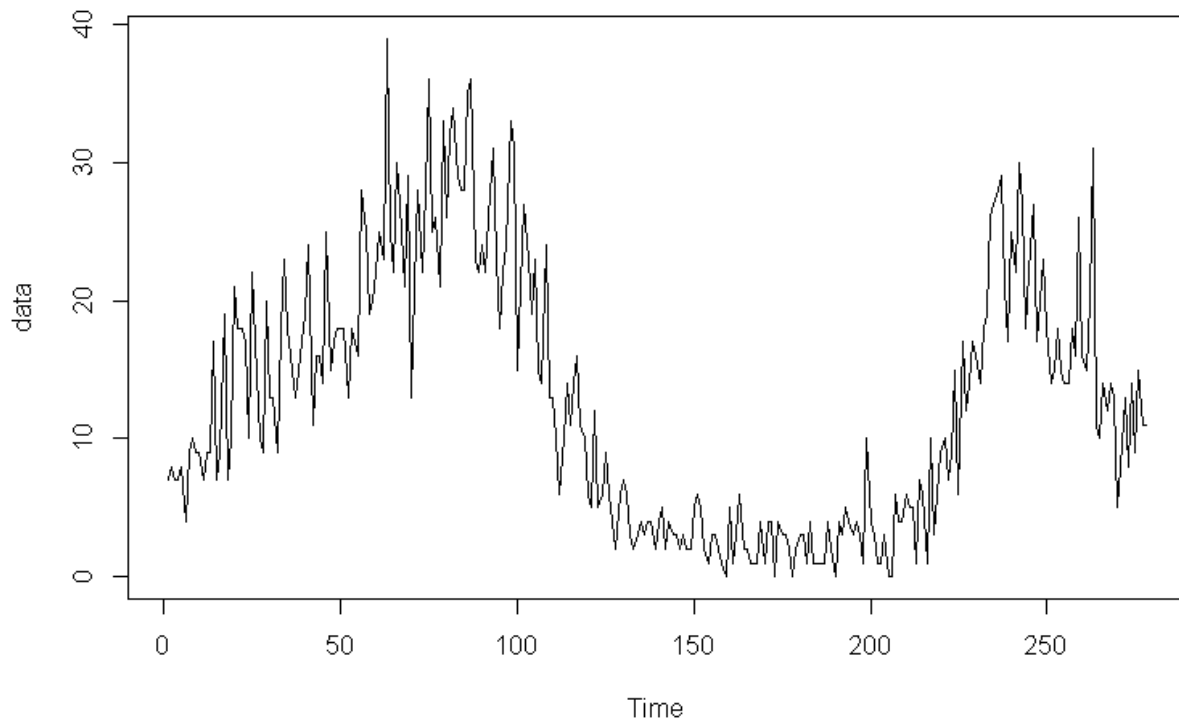
گزارش آمار فوتی در استان اصفهان با نرم افزار R

پرهام حسن

در مرحله اول داده را به نرم افزار پیوست میکنیم و آن را تمییز میکنیم.

```
> #####read data#####  
> dead = read.csv(file.choose())  
> dead[۱۵۳:۴۳۲,۴]  
[۱] ۷ ۸ ۷ ۷ ۸ ۴ ۹ ۱۰ ۹ ۹ ۷ ۹ ۹ ۱۷ ۷ ۹ ۱۹ ۷ ۱۰ ۲۱ ۱۸ ۱۸ NA ۱۷ ۱۰ ۲۲ ۱۷ ۱۰ ۹ ۲۰  
[۳۱] ۱۳ ۱۳ ۹ ۱۸ ۲۳ ۱۸ ۱۶ ۱۳ ۱۵ ۱۷ ۱۹ ۲۴ ۱۱ ۱۶ ۱۶ ۱۴ ۲۵ ۱۵ ۱۷ ۱۸ ۱۸ ۱۸ ۱۳ ۱۸ ۱۷ ۱۶ ۲۸ ۲۵ ۱۹  
۲۰  
[۶۱] ۲۲ ۲۵ ۲۳ ۳۹ ۲۵ ۲۲ ۳۰ ۲۷ NA ۲۱ ۲۹ ۱۳ ۲۱ ۲۸ ۲۲ ۲۸ ۳۶ ۲۵ ۲۶ ۲۱ ۳۳ ۲۶ ۳۲ ۳۴ ۲۹ ۲۸ ۲۸ ۳۵  
۳۶ ۲۳  
[۹۱] ۲۲ ۲۴ ۲۲ ۲۶ ۳۱ ۲۳ ۱۸ ۲۲ ۲۴ ۳۳ ۳۱ ۱۵ ۲۱ ۲۷ ۲۳ ۱۹ ۲۳ ۱۵ ۱۴ ۲۴ ۱۳ ۱۳ ۱۰ ۶ ۱۰ ۱۴ ۱۱ ۱۴ ۱۶  
۱۱  
[۱۲۱] ۱۰ ۶ ۵ ۱۲ ۵ ۶ ۹ ۶ ۴ ۲ ۶ ۷ ۶ ۳ ۲ ۳ ۴ ۳ ۴ ۴ ۲ ۴ ۵ ۲ ۴ ۳ ۳ ۲ ۳ ۲  
[۱۵۱] ۲ ۵ ۶ ۵ ۲ ۱ ۳ ۳ ۲ ۱ ۰ ۵ ۱ ۳ ۶ ۲ ۲ ۱ ۱ ۱ ۴ ۱ ۴ ۴ ۰ ۴ ۳ ۳ ۲ ۰  
[۱۸۱] ۲ ۳ ۳ ۱ ۴ ۱ ۱ ۱ ۱ ۴ ۲ ۰ ۴ ۳ ۵ ۴ ۳ ۴ ۳ ۱ ۱۰ ۴ ۳ ۱ ۱ ۳ ۰ ۰ ۶ ۴  
[۲۱۱] ۴ ۶ ۵ ۵ ۱ ۷ ۶ ۱ ۱۰ ۳ ۶ ۹ ۱۰ ۷ ۹ ۱۵ ۶ ۱۷ ۱۲ ۱۴ ۱۷ ۱۶ ۱۴ ۱۸ ۱۹ ۲۶ ۲۷ ۲۸ ۲۹ ۲۱  
[۲۴۱] ۱۷ ۲۵ ۲۲ ۳۰ ۲۷ ۱۸ ۲۲ ۲۷ ۱۷ ۲۰ ۲۳ ۱۸ ۱۴ ۱۵ ۱۸ ۱۵ ۱۴ ۱۴ ۱۸ ۱۶ ۲۶ ۱۶ ۱۵ ۲۰ ۳۱ ۱۱ ۱۰ ۱۴  
۱۲ ۱۴  
[۲۷۱] ۱۳ ۵ ۸ ۱۳ ۸ ۱۴ ۹ ۱۵ ۱۱ ۱۱  
> #dead[۱۵۳:۴۳۲,۲]  
> d = data.frame(dead[۱۵۳:۴۳۲,۲],dead[۱۵۳:۴۳۲,۴])  
> View(d)  
> data = d[,۲]  
> data = data[-c(۲۳,۶۹)]
```

در مرحله بعدی نمودار سری زمانی داده را رسم میکنیم. با استفاده از `ts.plot(data)`



برای بررسی اینکه داده ها رگرسیون هستند یا خیر مدل رگرسیونی زیر را به آن برازش می دهیم.

```

> #####model#####
> t = 1:length(data)
> tY = t ^ 2
> model = lm(data ~ t + tY)
> summary(model)

Call:
lm(formula = data ~ t + tY)

Residuals:
    Min       1Q   Median       3Q      Max
-16.535  -7.054  -1.568   6.011  24.410

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  2.13e+01  1.70e+00  13.280 < 2e-16 ***
t           -1.298e-01  2.700e-02 -4.809 1.73e-06 ***
tY           3.798e-03  9.210e-05  41.137 1.49e-08 ***
---
Signif. codes:  . '***' . . . '***' . . . '*' . . . '.' . . . ' '

Residual standard error: 8.80 on 270 degrees of freedom
Multiple R-squared:  0.108, Adjusted R-squared:  0.0989
F-statistic: 17.2 on 2 and 270 DF, p-value: 2.23e-07

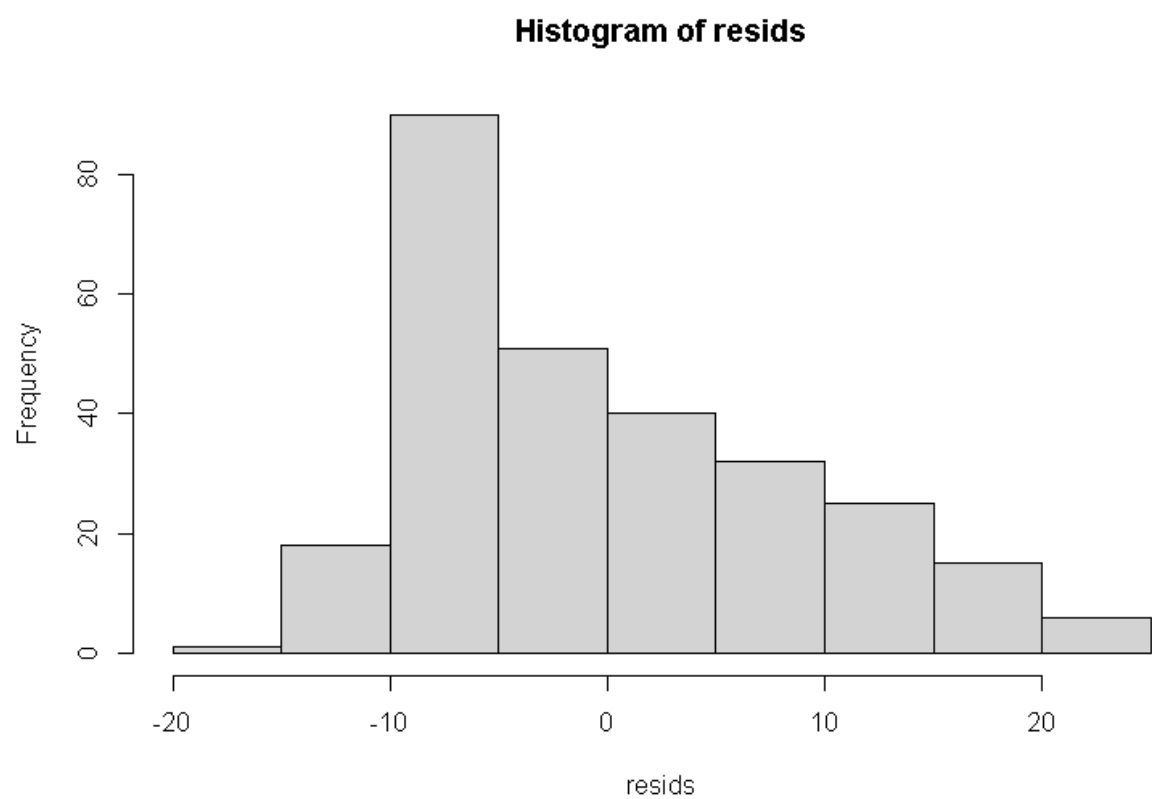
> resids = residuals(model)
> fits = fitted(model)

```

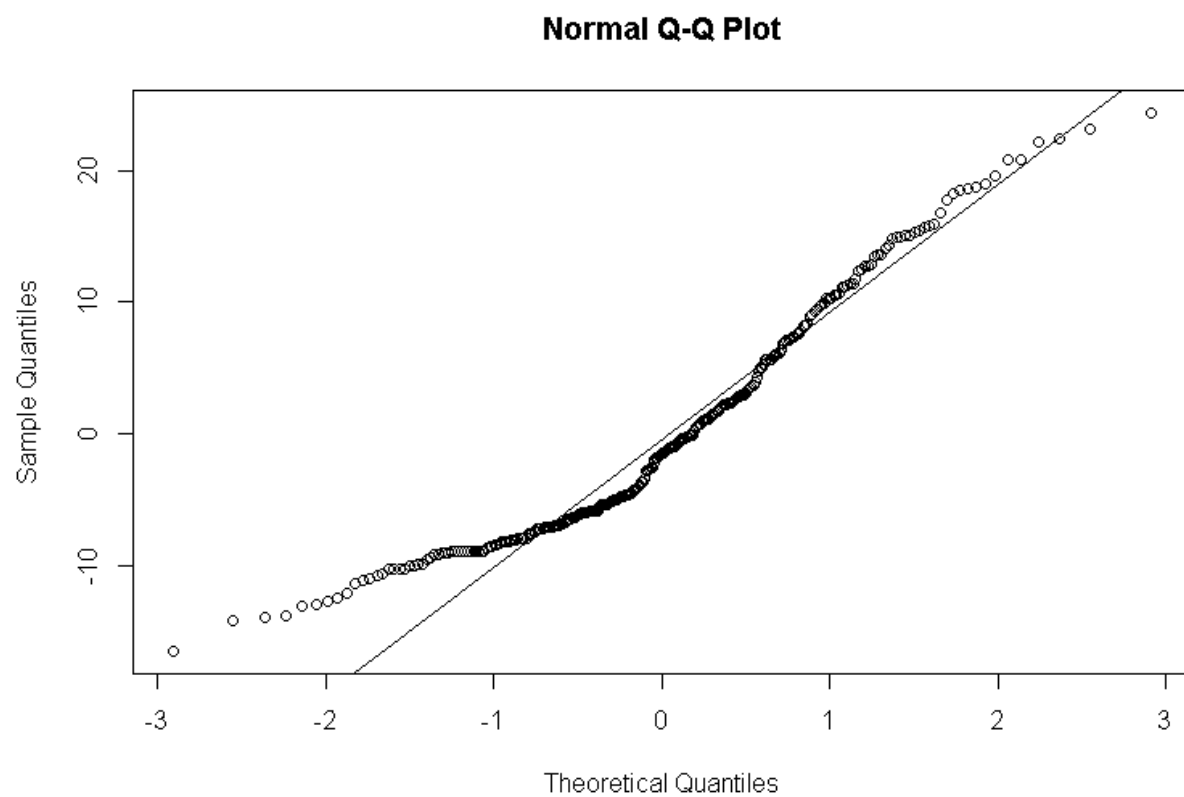
مقادیر p -value نشان می‌دهند که ضرایب رگرسیونی معنی دار هستند و R -squared مدل ۰.۰۹ می‌باشد و بسیار کم است که به معنای بزرگ بودن SSE یعنی فاصله مقادیر برازش داده شده با داده اصلی می‌باشد.

در مرحله بعدی فرض های اصلی رگرسیون را بررسی میکنیم.

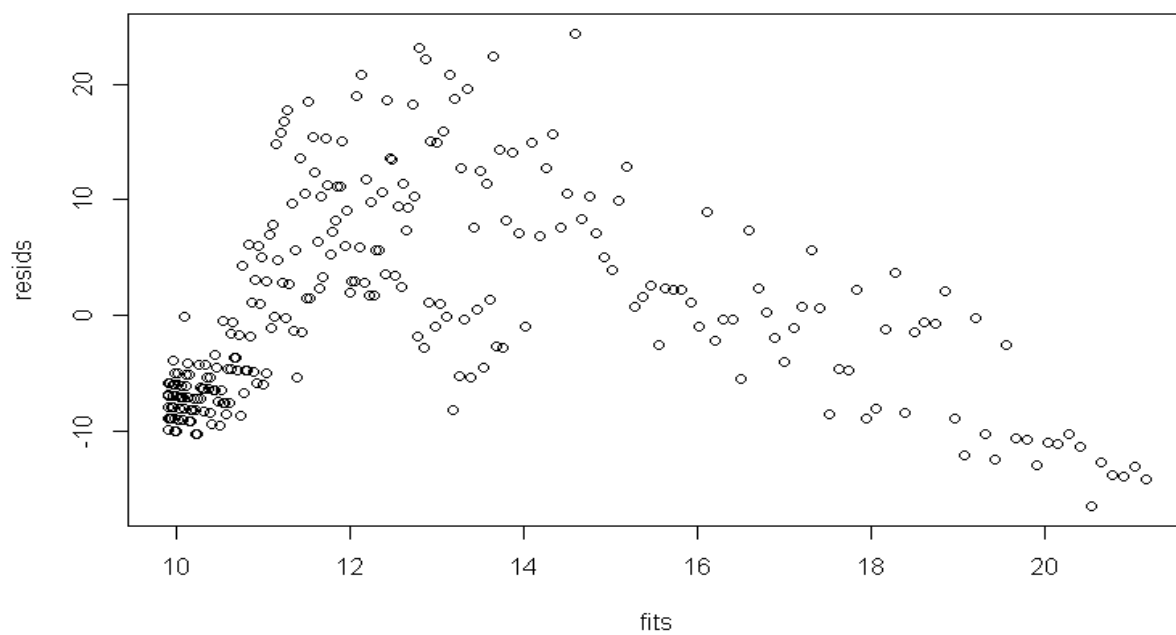
نمودارهای هیستوگرام ، qq و مانده ها در برابر مقادیر برازش داد شده و مانده ها در برابر زمان را رسم میکنیم.



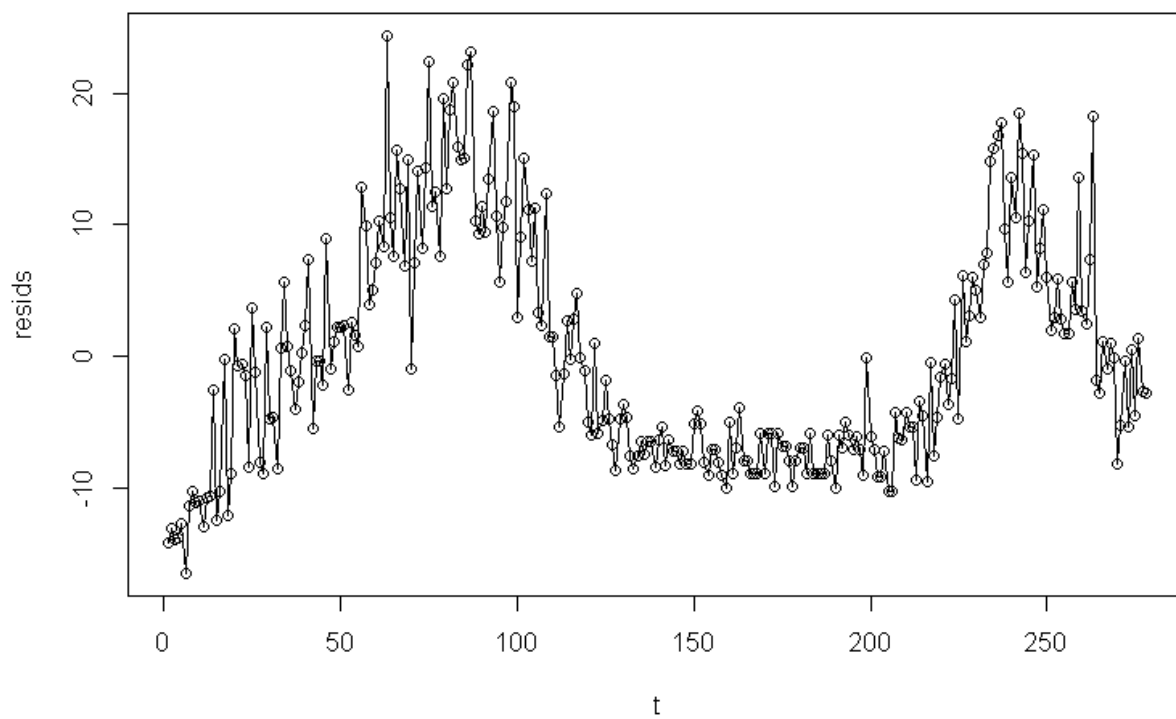
نمودار هیستوگرام شبیه به توزیع نرمال نمی‌باشد.



برای مقادیر منفی فاصله نقاط از خط زیاد می‌شود.



از نمودار مانده ها در برابر مقادیر برازش یافته همگن نبودن واریانس و روند داشتن مانده ها قابل مشاهده است.



تست نرمال بودن مانده های رگرسیونی را با استفاده از شاپیرو و اندرسون دارلینگ انجام میدهیم.

نرمال بودن مانده های رگرسیونی را با استفاده از شاپیرو و اندرسون دارلینگ انجام میدهیم.

```
> #####norm test#####  
> shapiro.test(resids)  
  
      Shapiro-Wilk normality test  
  
data:  resids  
W = 0.94087, p-value = 4.024e-09  
  
> library(nortest)  
> ad.test(resids)  
  
      Anderson-Darling normality test  
  
data:  resids  
A = 5.6747, p-value = 5.19e-14
```

مقادیر **p-value** در هر دو آزمون کمتر از ۰.۰۵ هستند که به معنای رد شدن فرض صفر این آزمون می باشد و فرض نرمال بودن مانده ها رد می شود.

از آزمون ناپارامتری **runs** برای بررسی تصادفی بودن یا روند داشتن مانده ها استفاده میکنیم.


```
> #####runs test####
```

```
> library(randtests)
```

```
> runs.test(resids)
```

Runs Test

```
data: resids
```

```
statistic = -۱۱,۸۹۷, runs = ۴۱, n۱ = ۱۳۹, n۲ = ۱۳۹, n = ۲۷۸, p-value < ۲,۲e-۱۶
```

```
alternative hypothesis: nonrandomness
```

چون p-value این آزمون کمتر از ۰.۰۵ است پس داده ها روند دارند و تصادفی نیستند.

و در نهایت آزمون همگنی واریانس را با استفاده از دستور آزمون `leveneTest` انجام می‌دهیم و برای اینکار داده ها را به ۲۳ گروه ۱۲ تایی دسته بندی میکنیم.

```

> #####variance test####
> fac = factor(rep(۱:۲۳, each = ۱۲))
> fac = c(fac , c(۲۳,۲۳))
> #fac = c(fac , c(۳۶,۳۶,۳۶))
> library(car)
Loading required package: carData
> leveneTest(resids, group = fac)
Levene's Test for Homogeneity of Variance (center = median)

      Df F value    Pr(>F)
group  ۲۲  ۳,۳۹۰.۲ ۱,۴۰۱e-۰۶ ***
      ۲۵۵
---
Signif. codes:  . '***' ۰,۰۰۱ '**' ۰,۰۱ '*' ۰,۰۵ '.' ۰,۱ '' ۱

Warning message:
In leveneTest.default(resids, group = fac) : fac coerced to factor.

```

چون p-value این آزمون کمتر از ۰.۰۵ است پس فرض همگن بودن واریانس مانده ها رد می شود.

و در نتیجه با توجه به رد شدن هر سه فرض اصلی رگرسیون ، پس شرایط استفاده از رگرسیون را نداریم

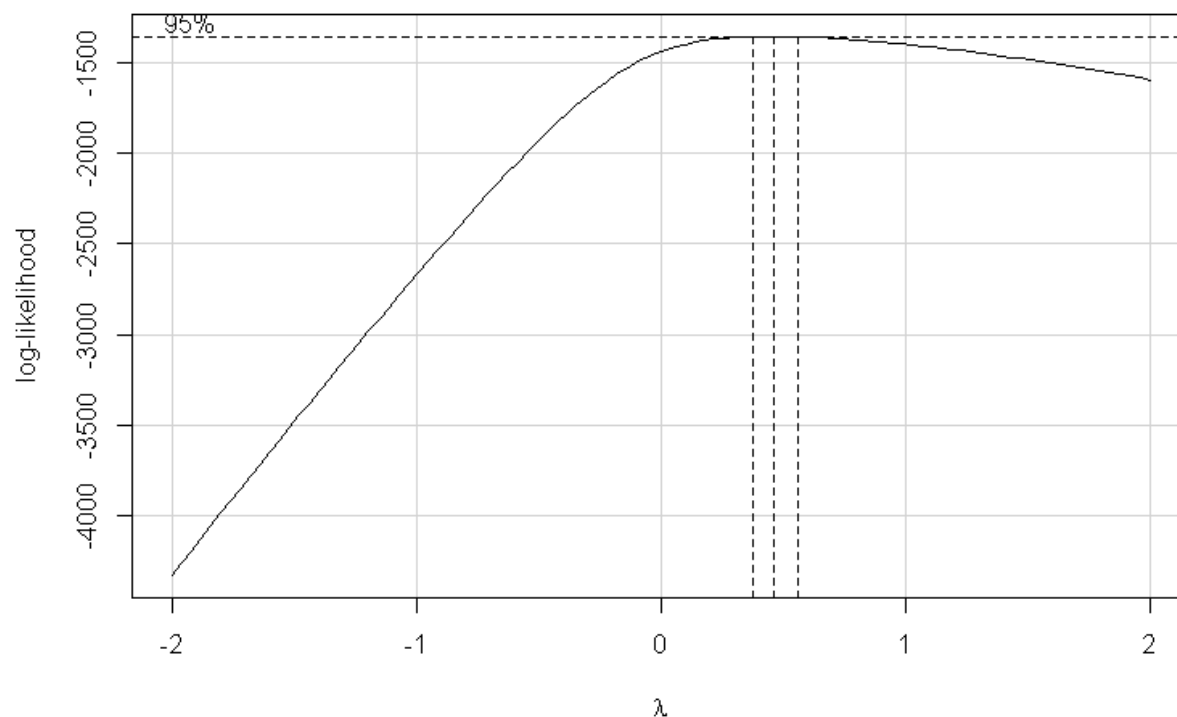
و ادامه کار را با سری های زمانی انجام می دهیم.

در ابتدا بررسی میکنیم که آیا داده های اصلی در واریانس مانا هستند یا خیر مشابه بالا این کار را انجام می دهیم
با استفاده از ۲۳ گروه ۱۲ تایی.

```
> #####variance test####
> fac = factor(rep(1:23, each = 12))
> fac = c(fac , c(23,23))
> #fac = c(fac , c(36,36,36))
> library(car)
> leveneTest(data, group = fac)
Levene's Test for Homogeneity of Variance (center = median)
      Df F value    Pr(>F)
group  22  3.3993 1.323e-06 ***
      255
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Warning message:
In leveneTest.default(data, group = fac) : fac coerced to
factor.
```

چون **p-value** این آزمون کمتر از ۰.۰۵ است پس فرض مانایی در واریانس را نداریم و از تبدیل باکس کاکس
برای مانایی در واریانس استفاده می کنیم.

```
> #####box-cox#####
> t = 1:length(data)
> library(MASS)
> library(car)
> data = data + 0.01
> boxCox(data ~ t, lambda = seq(from = -2, to = 2, by = .1))
> data.new = (data^0.5-1)/0.5
> leveneTest(data.new, group = fac)
Levene's Test for Homogeneity of Variance (center = median)
      Df F value Pr(>F)
group  22  1.1434 0.3011
      255
Warning message:
In leveneTest.default(data.new, group = fac) : fac coerced to
factor.
```



چون در داده ها روز هایی بدون فوتی هم داشتیم یعنی تعداد فوتی برابر صفر بوده و از آنجایی که در تبدیل باکس کاکس باید داده ها مثبت باشند پس به همه عدد 0.01 را اضافه می کنیم. چون $p\text{-value}$ آزمون نزدیک 0.301 است و بزرگتر از 0.05 است پس فرض مانایی در واریانس تایید می شود.

برای بررسی فرض مانایی در میانگین از آزمون ناپارامتری کروسکال والیس استفاده می کنیم که فرض توزیعی ندارد.

```

> #####mean test####
> kruskal.test(data.new, g = fac)

Kruskal-Wallis rank sum test

data: data.new and fac
Kruskal-Wallis chi-squared = 238.79, df = 22, p-value <
2.2e-16

> data.new.diff = diff(data.new)
> fac.diff = fac[-1]
> kruskal.test(data.new.diff, g = fac.diff)

Kruskal-Wallis rank sum test

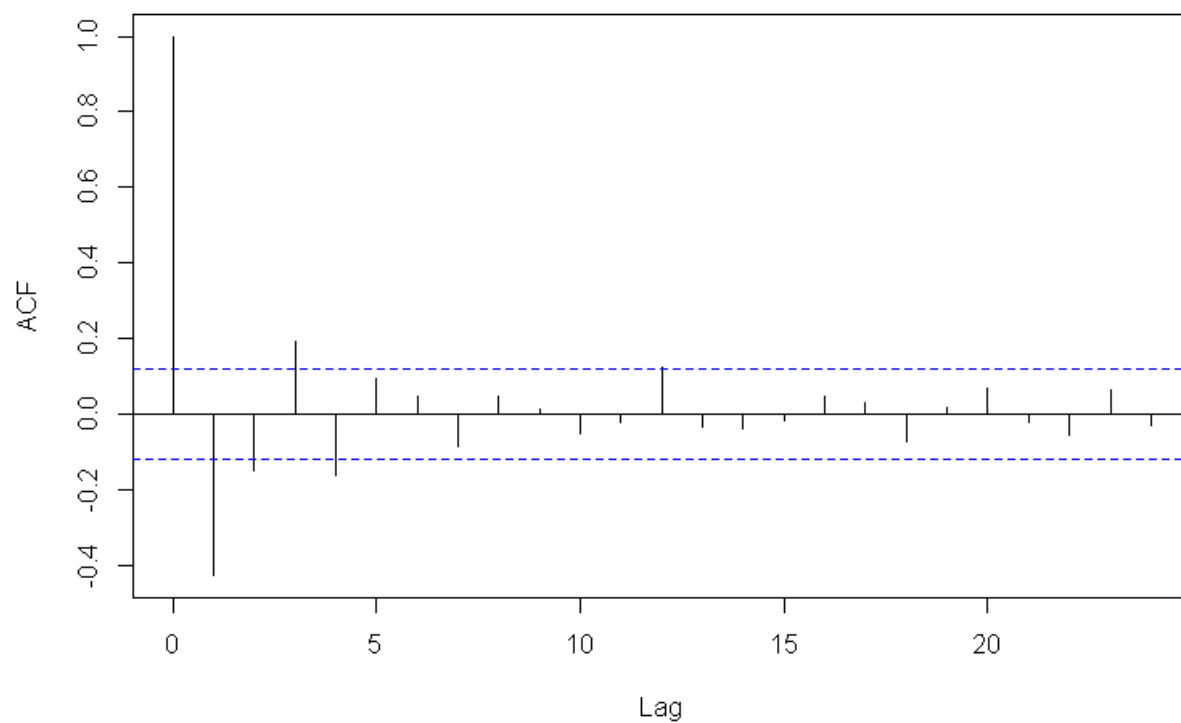
data: data.new.diff and fac.diff
Kruskal-Wallis chi-squared = 4.435, df = 22, p-value = 1

> leveneTest(data.new.diff, group = fac.diff)
Levene's Test for Homogeneity of Variance (center = median)
      Df F value Pr(>F)
group  22  1.2729 0.1895
      254
Warning message:
In leveneTest.default(data.new.diff, group = fac.diff) :
  fac.diff coerced to factor.

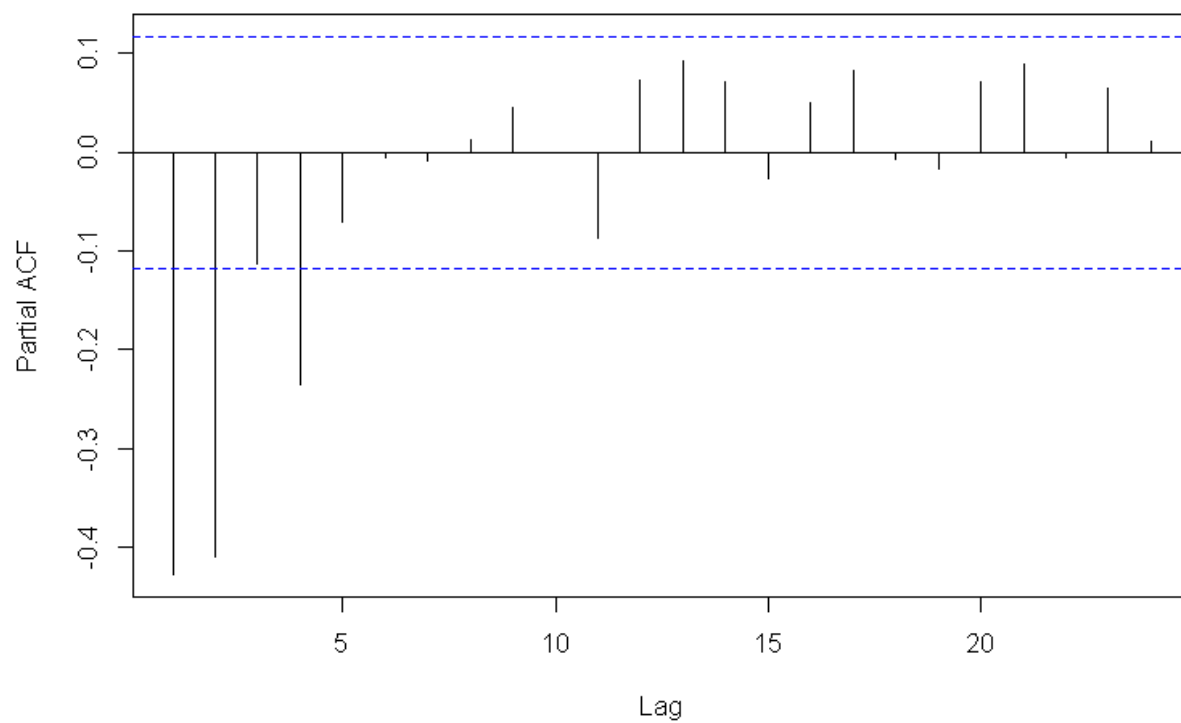
```

در ابتدا $p\text{-value} = 0.05$ است و فرض مانایی رد می‌شود همچنین پس از یک مرتبه تفاضلی کردن روی داده ها $p\text{-value}$ نزدیک یک می‌شود و داده ها در میانگین مانا می‌شوند. پس از مانایی در میانگین و واریانس نمودار های acf , $pacf$ را برای داده های جدید رسم میکنیم.

Series data.new.diff



Series data.new.diff



با توجه به نمودار های بالا برای اتورگرسیو ۲ مرتبه و برای میانگین متحرک ۳ مرتبه را در نظر می گیریم.

```
> #####model#####
> y = ts(data = data.new)
> order = c(2, 1, 3)
> library(forecast)
> fit <- Arima(y = y, order = order)
> summary(fit)
Series: y
ARIMA(2,1,3)

Coefficients:
          ar1      ar2      ma1      ma2      ma3
      0.2874  0.6051 -0.9725 -0.5688  0.6011
s.e.  0.1416  0.1340   0.1282   0.1920  0.0876

sigma^2 estimated as 1.543:  log likelihood=-451.17
AIC=914.34   AICc=914.65   BIC=936.08

Training set error measures:
MAPE      MASE      ME      RMSE      MAE      MPE
      ACF1
Training set 0.002518319 1.22874 0.9715128 -705.2714
۷۳۶.۷۲۳۲ ۰.۸۱۹۸۹۰۱ -۰.۰۳۰۰۱۲۳۵

> #####forecast#####
> plot(forecast(fit, h = 20))
```


Forecasts from ARIMA(2,1,3)

