

正會員 林成植\*      正會員 鄭永在\*

( A Study On The Subscriber Traffic Forecasting Method )

Seong Sik LIM\*, Yeong Jae JEONG\* Regular Members

要約 본 논문에서는 보다 정확하게 가입전화 트래픽을 예측하기 위한 회귀진단 기법을 제안하였다. 기존 방법은 다만 회귀계수만 추정하여 트래픽을 예측하였지만, 새로운 방법에서는, 각각의 측정치에서 트래픽 예측에 영향을 미치는 측정치와 영향을 미치지 않는 즉, 이용할 수 없는 자료를 식별하는 방법을 보였으며, 이 방법을 이용하여 예측하는 것이 기존 예측 방법보다 타당성이 있음이 시뮬레이션에 의하여 입증되었다.

ABSTRACT This paper proposed the method of regression diagnostics with a view to forecasting the subscriber traffic with more accuracy. Up to now, the traffic forecasting method is only used the estimated regression coefficient.

However, the new method showed how to discriminate between influential measuring data and unavailable data. And this method proved out more reasonable for forecasting than any other method through simulation.

I. 서론

$(x_{1i}, x_{2i}, \dots, x_{ki}, y_i), i=1, 2, \dots, n$ 가 있을 때, 일반적으로 중회귀 모형은

$$Y = X\beta + e \quad (1.1)$$

여기서  $X$ 는  $n \times (k+1)$  행렬이고

$$e \sim N(0, I\sigma^2)$$

이라 가정할 때,  $b$ 를  $\beta$ 의 추정벡터라 하면 최소승법에 의하여

$$b = (X'X)^{-1} X'Y \quad (1.2)$$

가 된다. 그러나 위와 같은 방법으로 회귀계수( $b$ )를 추정하여 가입전화 트래픽을 예측한다면 주어진 자료에서 보다 유용한 정보를 얻을 수가 없으므로 좀더 유용한 정보를 얻기 위하여 여러가지 진단기법이 필요하게 되었다.

각 진단기법의 종류는 다음장에서 설명을 하기로 하고, 먼저 기본적으로 알아야 할 사항을 보기로 하자. 우선, 잔차행렬을

$$e = Y - Xb = (I - X(X'X)^{-1} X')Y \quad (1.3)$$

으로 정의할 때,  $e$ 의 기대치와 분산공분산행렬은

$$E(e) = 0$$

$$V(e) = (I - X(X'X)^{-1}X')\sigma^2 \quad (1.4)$$

이 됨을 알 수 있다. 만일  $\sigma^2$ 의 추정벡터를  $S^2$ 이라 하면,  $S^2 = e'e/(n-k-1)$ 이 되며, 이때 분산공분산 행렬의 추정벡터는

$$\hat{V}(e) = (I - X(X'X)^{-1}X')S^2 \quad (1.5)$$

이 되며, 식 (1.5)에 포함되어 있는  $X(X'X)^{-1}X'$ 는 "Hat matrix"라 부르며 회귀진단에서는 매우 중요한 행렬로서

$$H = X(X'X)^{-1}X' \quad (1.6)$$

로 정의하며,  $i$ 번째 대각선상의 값을  $h_{ii}$ ,  $(i, j)$ 번째 값은  $h_{ij}$ 라고 한다

$$h_{ii} = x_i(X'X)^{-1}x_i'$$

$$h_{ij} = x_i(X'X)^{-1}x_j' \quad (1.7)$$

의 관계가 있다.

여기서 우리는 모든 가입전화 트래픽의 예측 경우가 단순회귀이므로 독립변수가 하나인 경우만을 실례로 분석하고자 한다. 이때 간단한 한 예로 단순회귀모형은 다음과 같고,

$$Y_i = \beta_1 + \beta_2 X_i + \epsilon_i, \quad i = 1, 2, \dots, n \quad (1.8)$$

$\beta_1, \beta_2$ 의 추정값을  $b_1, b_2$ 라 할 때, 잔차는  $e_i = Y_i - b_1 - b_2 X_i$ 가 되며, hat matrix는

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \quad (1.9)$$

이 된다.

잔차  $e_i$ 의 기대치와 분산은

$$E(e_i) = 0$$

$$V(e_i) = (1 - h_{ii})\sigma^2 \quad (1.10)$$

임을 알 수 있고,  $V(e_i)$ 의 추정치는

$$\hat{V}(e_i) = (1 - h_{ii})S^2 \quad (1.11)$$

이 되며, 여기서  $S^2 = \sum e_i^2 / (n-2)$ 이다.

기존 가입전화 트래픽 예측방식은 회귀계수( $b$ )를 추정하여 각 모형별로 트래픽을 예측한 결과가 예측 자료로 이용할 수 없는 측정치에 의해 영향을 받는 경우가 발생할 수 있기 때문에, 이와 같이 영향을 미치지 않는 측정치와 영향을 미치는 측정치를 식별함으로써 예측자료의 특성을 파악할 수 있으며, 결과를 분석하는데 도움을 얻을 수 있다.

다음장에서는 회귀진단에 필요한 기법들을 소개하고, 3장에서는 실제 가입전화 트래픽 자료를 회귀진단 기법에 적용하여 얻은 회귀분석 결과와 기법을 이용하지 않고 분석한 결과를 비교 검토하였다.

## II. 회귀진단 기법의 종류

분석 결과에 따라 그 결과에 영향을 주는 측정치들중에는 이용가치가 없는 자료도 있고 이용가치가 있는 자료도 있다. 전자의 경우를 이상치라 하고 후자의 경우를 영향을 크게 주는 측정치라고 말한다.

다음에서 이상치와 크게 영향을 주는 측정치를 찾는 방법을 보기로 하자.

### II.1 이상치 검출

#### 1) 표준화 잔차

$Y$ 가  $N(X\beta, I\sigma^2)$ 의 분포를 따르므로  $e$ 는  $N(0, (I-H)\sigma^2)$ 의 분포를 따름을 알 수 있다. 또한  $i$ 번째 잔차( $e_i$ )의 분포는

$$e_i \sim N(0, (1-h_{ii})\sigma^2) \quad (2.1)$$

이 되므로 따라서,  $e_i$ 를 표준 정규분포화 하면

$$Z_i = \frac{e_i}{\sigma \sqrt{1-h_{ii}}} \sim N(0,1) \quad (2.2)$$

와 같으며,  $\sigma$  대신 추정치  $S$ 를 이용하면 표준화 잔차

$$\tau_i = \frac{e_i}{S \sqrt{1-h_{ii}}} \quad (2.3)$$

가 얻어진다. 표준화 잔차  $\tau_i$ 의 절댓값  $|\tau_i|$ 가 과다하게 크면 이상치라 생각할 수 있다. 이상치 판단기준은 [별첨 1]에 있다.

2) 스트던트화 잔차

이 잔차는 표준화잔차에서  $S$ ,  $b$ 대신에  $i$ 번째 추정치  $Y_i$ 를 제외시키고  $n-1$ 개의 추정치로부터 얻어지는  $S, b$  ( $S, b$ 를 각각  $S(i), b(i)$ 로 표기하자)를 사용 한다면, 그때  $S^2$ 은

$$S^2(i) = \frac{1}{n-k-2} \sum_{k \neq i} [Y_k - X_k b(i)]^2 \quad (2.4)$$

이 되며, 위 식을 간단한 형태로 변형하면

$$(n-k-2)S^2(i) = (n-k-1)S^2 - \frac{e_i^2}{1-h_{ii}} \quad (2.5)$$

이 된다.

스트던트화 잔차는 표준화잔차에서  $S$ 대신  $S(i)$ 로 대체시킨 식으로 다음과 같이

$$\tau_i^* = \frac{e_i}{S(i) \sqrt{1-h_{ii}}} \quad (2.6)$$

로 정의되며, 자유도  $n-k-2$ 를 가진  $t$ 분포 형태를 따른다.

식 (2.6)를 좀더 간략하게 정리하면

$$\begin{aligned} \tau_i^* &= \frac{e_i}{S(i) \sqrt{1-h_{ii}}} \\ &= \tau_i \cdot \left[ \frac{n-k-2}{n-k-1-\tau_i^2} \right]^{\frac{1}{2}} \end{aligned} \quad (2.7)$$

이 된다.

여기서  $|\tau_i^*|$ 가 어느 정도 커야만  $Y_i$ 를 이상치로 볼 수 있는지를 결정하기 위한 검정방법은

$$|\tau_i^*| \geq t(n-k-2; \alpha/2) \quad (2.8)$$

이면, 유의수준  $\alpha$ 에서  $Y_i$ 를 이상치라 판정할 수 있다.

II.2. 영향을 크게 주는 추정치 검출

1) 행렬  $H$ 의 대각선 원소

최소자승법에 의하여 적합된 중회귀 모형은

$$\hat{Y} = Xb = X(X'X)^{-1} X'Y = HY \quad (2.9)$$

가 얻어지며, 예측치  $\hat{Y}_i$ 를  $Y_i$ 로 편미분하면

$$\frac{\delta \hat{Y}_i}{\delta Y_i} = h_{ii}$$

로 행렬  $H$ 의  $i$ 번째 대각선 원소가 된다.

따라서  $h_{ii}$ 의 값이 1에 가까울 수록  $Y_i$ 는  $\hat{Y}_i$ 에 영향을 크게 준다고 판단할 수 있다. 즉 이 말은  $h_{ii}=1$ 이면  $Y_i=\hat{Y}_i$ , 다시 말하면  $e_i=0$ 을 의미한다.

여기서  $H$ 의 성질을 몇가지 살펴보기로 하자.

-  $HH=H$ 가 성립하므로 멱등행렬이고,  $H$ 의 계수는

$$\text{rank}(H) = \sum_{i=1}^n h_{ii} = k+1$$

이 된다.

- 행렬  $H$ 는 양반정치(Positive Semi-definite)행렬이므로  $h_{ii} \geq 0$ 이고,  $(1-h_{ii})\sigma^2 \geq 0$ 이어야 하므로  $h_{ii} \leq 1$ 이 되어야 한다. 따라서

$$0 \leq h_{ii} \leq 1$$

의 관계가 있다. 그리고  $h_{ii}$ 의 평균치는  $(k+1)/n$ 이다.

$h_{ii}$ 가 어느 정도 커야만  $Y_i$ 가 영향을 크게 주는 측정치인가를 판단하는 기준은

$$\frac{n-k-1}{k} \cdot \frac{h_{ii} - (1/n)}{1-h_{ii}} \geq F(k, n-k-1; \alpha) \quad (2.10)$$

이면,  $Y_i$ 가 크게 영향을 받는 측정치라 말할 수 있다. 그런데  $k > 10$  이고  $n-k-1 > 50$  이상일 때 95%의 신뢰도에서 F 분포의 값은 2보다 작은 값이 나타난다. 그러므로 좀더 간단히 하면 다음과 같은 식으로 판단기준을 설정하여  $Y_i$ 가 영향을 주는 측정치인지 아닌지를 판단할 수 있다.

$$h_{ii} \geq (2k+1)/n \quad (2.11)$$

## 2) DFFITS

이 경우는  $b$ 와  $b(i)$ 의 차이 값이 가장 큰  $\hat{Y}_i$ 를 찾는 방법을 말한다.  $i$ 번째  $Y_i$ 를 제외시킨 상태에서의 적합된 모형의 차이 값은

$$\begin{aligned} DFFIT &= \hat{Y}_i - \hat{Y}_i(i) = X_i(b - b(i)) \\ &= \frac{h_{ii} \cdot e_i}{1-h_{ii}} \end{aligned} \quad (2.12)$$

이되며, 위 식을  $S(i)$   $h_{ii}$ 로 나누면 DFFITS(i)는

$$DFFITS(i) = \left[ \frac{h_{ii}}{1-h_{ii}} \right]^{1/2} \cdot \frac{e_i}{S(i) \sqrt{1-h_{ii}}} \quad (2.13)$$

으로 정의되며, 위 식은 스트던트화 잔차의 식을 이용하여

$$DFFITS(i) = \left[ \frac{h_{ii}}{1-h_{ii}} \right]^{1/2} \cdot \tau_i \quad (2.14)$$

으로 바꿀 수 있다.

DFFITS(i)가 어느 정도 커야만  $Y_i$ 가 영향을 크게 주는 측정치로 판단할 수 있는가를 보기 위해

여, 먼저  $|\tau_i| \geq 2$  이고,  $h_{ii} \geq (k+1)/n$  이므로, 식(2.14)에서

$$\begin{aligned} |DFFITS(i)| &\geq \left[ \frac{(k+1)/n}{1-(k+1)/n} \right]^{1/2} \cdot 2 \\ &= 2 \left[ \frac{k+1}{n-k-1} \right]^{1/2} \\ &\geq 2 \left[ \frac{k+1}{n} \right]^{1/2} \end{aligned} \quad (2.15)$$

으로 되므로 위식을 이용하여 판단할 수 있다.

## 3) Cook의 통계량

일반적으로 모든  $b$ 에 대하여  $100(1-\alpha)\%$  신뢰구간은

$$\frac{(b - b^*)'X'X(b - b^*)}{(k+1)S^2} \leq F(k+1, n-k-1; 1-\alpha) \quad (2.16)$$

로 정의되는데 Cook의 통계량은  $b$  대신에  $b(i)$ 를 대입시킨 즉,

$$D(i) = \frac{(b - b(i))'X'X(b - b(i))}{(k+1)S^2} \quad (2.17)$$

으로 정의되며, 이는  $F(k+1, n-k-1; 1-\alpha)$ 의 분포를 가지며, 이 값을 크게 하는 것을 영향을 크게 주는 측정치로 판단할 수 있다.

$D(i)$ 을 간단히 하기 위하여 표준화잔차를 이용하면

$$\begin{aligned} D(i) &= \frac{1}{k+1} \cdot \frac{h_{ii} e_i^2}{S^2 (1-h_{ii})^2} \\ &= \frac{h_{ii}}{(k+1)(1-h_{ii})} \cdot \tau_i^2 \end{aligned} \quad (2.18)$$

와 같다. 위식에서  $\tau_i^2$ 의 값이 커지면  $D(i)$ 도 값이 커지고,  $h_{ii}$ 값이 커지면  $D(i)$ 의 값이 커짐을 알 수 있다.

영향을 크게 주는 측정치의 판단기준은

$$D(i) \geq F(k+1; n-k-1; \alpha) \quad (2.19)$$

이 되는  $Y_i$ 를 찾는 방법으로 Cook는 대략적으로

$$D(i) \geq F(k+1, n-k-1; 0.5) \quad (2.20)$$

이면 영향을 크게 주는 추정치로 판단해도 좋다고 제안하였다.

#### 4) Andrew-Pregibon의 통계량

행렬  $X$ 와 벡터  $Y$ 를 같이 고려하여, 다음의 통계량

$$AP(i) = \frac{|X^*(i)'X^*(i)|}{|X^* ' X^*|} \quad (2.21)$$

을 제안하였다. 여기서  $X^*(i) = (X(i):Y(i))$ ,  $X^* = (X:Y)$ 을 의미하는데,  $X^*(i)$ 는  $X$  행렬에서  $i$  번째 행을 제거시킨 것이다.

위식을 다시 표현하면

$$AP(i) = (1 - \frac{e_i^2}{(1-h_{ii})SSE})(1-h_{ii}) \\ = 1-h_{ii} - \frac{e_i^2}{(n-k-1)S^2} \quad (2.22)$$

이 된다. 여기서 SSE는 잔차 자승합을 의미하며,  $SSE = (n-k-1)S^2$ 으로 이미 알고 있다. 식(2.22)에서  $h_{ii}$ 가 크거나 또는  $e_i^2$ 이 크면  $AP(i)$ 값은 작아지며, 이들 값 중 가장 작은  $AP(i)$ 의 값을  $Y_i$ 가 영향을 크게 주는 추정치로 판단하면 된다.

#### 5) COVRATIO

회귀계수 추정치의 분산공분산 행렬은

$$V(b) = (X'X)^{-1} \sigma^2 \\ V(b(i)) = \{X(i)'X(i)\}^{-1} \sigma^2 \quad (2.23)$$

이다.  $\sigma^2$  대신에 각각  $S^2, S^2(i)$ 를 대입시킨 후, 그 행렬식의 비율을 COVRATIO라 정의하고, 그 식은

$$COVRATIO(i) = \frac{|S^2(i)[X(i)'X(i)]^{-1}|}{|S^2(X'X)^{-1}|} \quad (2.24)$$

와 같다. 위 식(2.16)을 좀더 간략하게 하기 위해 먼저 다음과 같은 식

$$|X'(i)X(i)| = (1-h_{ii})|X'X|$$

을 이용하여 계산하면

$$COVRATIO(i) = \frac{|S^2(i)|^{k+1}}{(S^2)^{k+1}} \cdot \frac{1}{1-h_{ii}} \quad (2.25)$$

되는데, 식(2.5)와 (2.6)를 사용하면 식(2.25)는

$$COVRATIO(i) = \frac{1}{[\frac{n-k-2}{n-k-1} + \frac{\tau_i^2}{n-k-1}] (1-h_{ii})} \quad (2.26)$$

이 된다.

COVRATIO(i)의 값이 1에 가까우면  $Y_i$ 는 별로 영향을 못주며, 1에서 멀어질수록 영향을 크게 주는 추정치라 말할 수 있다.

$Y_i$ 가 영향을 크게 주는 추정치라고 판단하는 기준이 다음과 같으면

$$|COVRATIO(i)-1| \geq 3(k+1)/n \quad (2.27)$$

$Y_i$ 를 영향을 크게 주는 추정치라고 한다.

#### 6) FVARATIO

하나의 추정치  $Y_i$ 가 제외되었을 때,  $Y_i$ 의 분산이 어떻게 변화하는가에 대한 것으로서, 이것을 하기위해 먼저  $\hat{Y}_i$ 와  $\hat{Y}_i(i)$ 의 분산공분산 행렬은

$$V(\hat{Y}_i) = h_{ii} \cdot S^2 \\ V(\hat{Y}_i(i)) = \frac{h_{ii}}{1-h_{ii}} S^2 \quad (2.28)$$

이 되며, 이를 비율로 구하면

$$\begin{aligned}
 FVARATIO(i) &= \frac{h_{ii} \cdot S^2(i)/(1-h_{ii})}{h_{ii} \cdot S^2} \\
 &= \frac{S^2(i)}{S^2(1-h_{ii})} \quad (2.29)
 \end{aligned}$$

연계 된다.

$Y_i$ 가 영향을 크게 주는 측정치인지 아닌지의 판단기준은

$$\begin{aligned}
 FVARATIO(i) &\leq 1 - 3/n \quad \text{이거나,} \\
 FVARATIO(i) &\geq 1 + (2k+3)/n \quad (2.30)
 \end{aligned}$$

이면 영향을 주는 측정치로 볼 수 있다.

### II. 3. 측정치 측도의 비교

2절에서 이상치를 검출하는 방법으로 2 가지를, 영향을 크게 주는 측정치를 검출해내는 방법으로 6 가지를 다루었으며, 이들 간에는 상호 연관성이 있음을 알 수 있다. 이를 한번 고찰해 보기로 하자.

이상치를 찾는 2가지 방법 즉, 표준화 잔차 및 스트던트화 잔차는 모두 잔차( $e_i$ )의 함수이며,  $|e_i|$ 가 커지면 이들 값들은 커진다. 일반적으로 잔차가 매우 크면 이상치라 볼 수 있다. 또한 이들 측도들은 "hat matrix"  $h_{ii}$ 의 함수로서  $h_{ii}$ 가 커지면 이들 값들도 커진다. 따라서  $h_{ii}$ 가 크고  $e_i$ 도 크면  $Y_i$ 는 이상치라 판단할 수 다.

영향을 크게 주는 측정치를 찾는 방법 중에서,  $h_{ii}$  값은 종속변수  $Y$ 의 값과는 직접 관련이 없고  $X$ 에 대한 유관한 측도 이므로,  $h_{ii}$ 의 값이 크면  $i$  번째가 다른 측정치로부터 멀리 떨어져 있다고 판단할 수 있다. 따라서  $h_{ii}$ 의 값이 크면 그 측정치  $Y_i$ 는 영향을 크게 주는 측정치가 될 가능성이 높다.

그리고 영향을 크게 주는 측도로서  $DFFIT(i)$ ,  $D(i)$ ,  $AP(i)$ ,  $COVRATIO(i)$ ,  $FVARATIO(i)$  들은 모

두 측정된 회귀계수들의 변화를 검출해 내는 방법으로 만들어져 있어서 영향을 크게주는 측정치를 찾는 측도로 쓰이나, 이들은 모두 잔차 $e_i$ 의 함수로 나타난다.

따라서, 이들은 영향을 크게 주는 측정치로 판단될 때에는 동시에 이상치가 될 가능성도 있다.

만일, 이상치로 판정되었다고 하여 이것을 무조건 버리고 나머지 만으로 분석을 실시하는 것은 옳은 방법이 아니다. 이상치로 판정되면 원인을 규명하여야 한다. 원인이 규명되면 다시 측정하여 대체시키는 경우가 있고, 측정이 불가능한 경우는 측정치를 제거시킬 수도 있다.

또한 영향을 크게 주는 측정치로 판정이 되면, 일단 이 측정치가 어떻게 해서 가장 영향을 크게 주는가를 검토해 볼 필요가 있다. 만일 다른 측정치들로부터 멀리 떨어져 있는 것이 원인이라면, 그 중간에 몇개의 측정치를 추가하여 진단을 실시하고 분석을 하는 것이 바람직하다. 만일, 다른 측정치들로부터 멀리 떨어져 있지 않는데도 불구하고 잔차가 커서 영향을 크게주는 측정치로 판정되었다면, 이것은 이상치이므로 이상치와 동일하게 취급하면 된다.

### III. 트래픽 예측 방법의 비교

이 장에서는 서울 신촌국의 NO1A 가입 전화 실측 트래픽 자료를 이용하여 트래픽 예측 결과를 비교 설명하고자 한다.

먼저, 가입 전화 실측 트래픽 자료를 수집, 분석한 후 계절지수를 고려하여 보정트래픽을 산정한다. 신촌국의 연도별 보정 트래픽은 [표 3-1] 과 같으며, 이에 대한 분기별 보정트래픽 산정도는 [그림 3-1] 과 같다.

[표 3-1] 신촌구 보정 트래픽

년 도	분 기	보정트래픽
'84	4/4	0.0387
'85	1/4	0.0390
	2/4	0.0403
	3/4	0.0385
'86	4/4	0.0398
	1/4	0.0387
	2/4	0.0414
	3/4	0.0411
'87	4/4	0.0419
	1/4	0.0403
	2/4	0.0397
	3/4	0.0275
'88	4/4	0.0408
	1/4	0.0419
	2/4	0.0510
	3/4	0.0436
'89	4/4	0.0412
	1/4	0.0396
	2/4	0.0416
	3/4	0.0417
'90	4/4	0.0436
	1/4	0.0419
	2/4	0.0427
	3/4	0.0429

t : 시간  
 a, b : 회귀계수  
 K : 상한치  
 exp : 자연대수

III.1 기존 트래픽 예측 방법

기존 트래픽 예측 방식은 회귀계수를 추정하여 적합된 회귀모형을 이용하여 익년도 가입 전화 트래픽을 예측하기 때문에 각 모형에 대한 타당성 검토에 관계없이 F-통계량 값이 가장 큰 모형을 선정하고, 이때 선정된 모형에 의해 예측된 값을 최적 예측치로 선정하므로써 실제 타당성이 없는 모형인데도 그 모형의 예측치를 전화교환시설 설계에 적용하므로써 신뢰성이 떨어진다고 볼 수 있다.

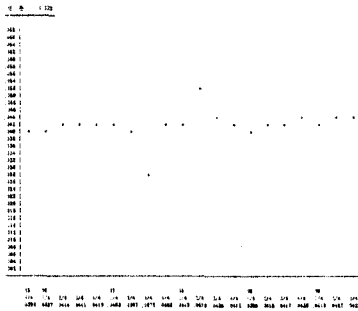
III.2 회귀진단에 의한 트래픽 예측

1절의 문제점을 보완하기 위해서 가입전화 트래픽 예측에 영향을 미치는 추정치와 영향을 미치지 않는 추정치를 검출하는 분석과 함께 영향을 미치는 추정치이면서 이상치인 경우의 자료도 예를 통하여 분석하여 보자.

[표 3-3]은 각 모형별 회귀진단방법을 이용하여 트래픽 예측 결과에 영향을 크게주는 추정치와 이상치를 검출하기 위한 표이며, 유의수준을  $\alpha=0.05$ 로 하자.

이상치를 검출하는 방법으로서 표준화잔차( $\tau_i$ )에 의한 이상치 검출방법은 [표 3-1]에서  $k=1$ 이고  $n=24$ 이므로 [별첨 1]의 표에는 기각치가 없으므로 보간법을 이용하여 구하면  $n=24$ 에서 기각치가 2.868 이 된다. 따라서 각 모형별  $|\tau_i|$ 의 값이 2.868보다 큰 분기는 87년 3/4분기 자료이므로, 이 분기의 트래픽 0.0275 Er1은 이상치임을 알수있다.

다음으로 스트던트회잔차( $\tau_i^*$ )에 의한 이상치를 검출하는 방법으로서, 자유도가 21이므로, 이상치로 판단하는 기준인 기각치는  $t(21;0.025)=2.080$ 이므로  $|\tau_i^*|$ 가 이값 보다는 큰 것들은 이상치



[그림 3-1] 신촌구 보정트래픽 그래프

또한 가입 전화 트래픽 예측에 사용되는 추세모형은 [표 3-2]와 같다.

[표 3-2] 추세모형

모 형	산 출 식
직 선	$Y_t = a + bt$
지 수	$Y_t = a \cdot \exp(bt)$
수정지수	$Y_t = K - a \cdot \exp(bt)$
성장곡선	$Y_t = K / (1 + a \cdot \exp(bt))$

주)  $Y_t$  : 가입전화 보정트래픽





면 12,14번째 측정치 자료는 이상치이면서 영향을 크게 주는 측정치로 판단이 되며, 1,2,23,24 번째 측정치는 영향을 크게주는 측정치로 판단할 수 있으며, 위와같이 12,14번째 이상치 자료는 보다 세밀한 분석·검토를 요한다.

그러므로 보다 정확한 트래픽을 예측하기 위하여, 이상치로 판정된 측정치는 제거하거나 수정하여야 하며 그 결과는 기존 트래픽 예측 방법 보다 정확한 값을 갖는 다고 볼 수 있다. 그 이유는 F-통계량 값에 의해 모형의 타당성이 입증되기 때문인데 모형의 타당성에 대한 논의는 다음절에서 하기로하자.

### Ⅲ.3 예측 결과 비교

추정된 회귀계수를 이용하여 예측한 결과와 회귀진단을 이용하여 이상치를 제거한 후 적합된 회귀모형을 이용하여 예측한 결과에 대해 비교하고자 한다.

모형에 대한 타당성 검정은 F-통계량 값을 이용하는데, 먼저 이상치를 제거하기 전과 이상치를 제거한 후의 모형에 대한 타당성을 검토하기 위하여 다음 [표 3-4]의 F-통계량 값을 비교하여 보자.

[표 3-4] F-통계량 값 비교

모형	F - 통계량 값	
	이상치 제거 전	이상치 제거 후
직 선	2.9369	25.7884
지 수	2.8872	25.6442
수정지수	2.7292	26.0946
성장곡선	2.8817	26.0332

유의수준  $\alpha = 0.05$ 에서 모형의 타당성 검정에 대한 기각치는  $F(1,22;0.05) = 4.28$  이므로 이상치 제거전의 F-통계량 값은 기각치 보다 작은 값이기 때문에 모형이 타당하지 않으며, 이상치를 제거한 후의 F-통계량값은 기각치 보다 크기 때문에 모형

이 타당하다는 것을 [표 3-4]에서 알 수 있다.

이와 같이 이상치를 제거하고 가입 전화 트래픽을 예측하는 것이 보다 정확하고 타당성이 있는 것이므로, 회귀진단을 이용하여 불필요한 측정치를 검출해 내고 가입전화 트래픽을 예측하는 것이 유의하다고 본다

### Ⅳ.결론

본 논문에서는 가입 전화 트래픽을 예측하는데 있어서 이상치를 제거하기 전에 예측한 결과와 이상치를 제거한 후 예측결과를 비교할 때 제거하는 모형의 타당성이 유의하지만 제거전은 타당성이 없는 것으로 나타났다.

그러므로 모든 자료는 회귀진단기법을 이용하여 먼저 이상치와 영향을 크게 주는 자료를 분석검토하여 이상치는 제거를 하고 가입전화 트래픽을 예측하는 것이 신뢰성이 있으므로 이 방법을 이용하는 것이 좋다고 생각된다.

또한 기존 방법에서 탈피하여 보다 나은 시계열 분석방법을 연구하여 정확하게 트래픽을 예측함으로써 가입자에게 양질의 서비스를 제공할 수 있고 경제적인 통신투자사업을 할 수 있게 하기 위하여 계속 연구가 되어야 할 것으로 사료된다.

[부록 1] 각 용어 상용에서의 기각치  
통계수준  $\alpha = 0.05$

표준편차 n	자유도 k																
	1	2	3	4	5	6	8	10	15	20	25	30	40	50	60	70	
1	1.00																
2	2.00	1.99															
3	2.16	2.02	1.99														
4	2.24	2.09	2.05	1.97													
5	2.31	2.15	2.11	2.05	1.97												
6	2.37	2.21	2.17	2.11	2.04	1.97											
7	2.43	2.26	2.22	2.16	2.09	2.02	1.95										
8	2.48	2.31	2.27	2.21	2.14	2.07	2.00	1.93									
9	2.53	2.36	2.32	2.26	2.19	2.12	2.05	1.98	1.91								
10	2.58	2.41	2.37	2.31	2.24	2.17	2.10	2.03	1.96	1.89							
11	2.63	2.46	2.42	2.36	2.29	2.22	2.15	2.08	2.01	1.94	1.87						
12	2.68	2.51	2.47	2.41	2.34	2.27	2.20	2.13	2.06	1.99	1.92	1.85					
13	2.73	2.56	2.52	2.46	2.39	2.32	2.25	2.18	2.11	2.04	1.97	1.90	1.83				
14	2.78	2.61	2.57	2.51	2.44	2.37	2.30	2.23	2.16	2.09	2.02	1.95	1.88	1.81			
15	2.83	2.66	2.62	2.56	2.49	2.42	2.35	2.28	2.21	2.14	2.07	2.00	1.93	1.86	1.79		
16	2.88	2.71	2.67	2.61	2.54	2.47	2.40	2.33	2.26	2.19	2.12	2.05	1.98	1.91	1.84	1.77	
17	2.93	2.76	2.72	2.66	2.59	2.52	2.45	2.38	2.31	2.24	2.17	2.10	2.03	1.96	1.89	1.82	
18	2.98	2.81	2.77	2.71	2.64	2.57	2.50	2.43	2.36	2.29	2.22	2.15	2.08	2.01	1.94	1.87	
19	3.03	2.86	2.82	2.76	2.69	2.62	2.55	2.48	2.41	2.34	2.27	2.20	2.13	2.06	1.99	1.92	
20	3.08	2.91	2.87	2.81	2.74	2.67	2.60	2.53	2.46	2.39	2.32	2.25	2.18	2.11	2.04	1.97	
21	3.13	2.96	2.92	2.86	2.79	2.72	2.65	2.58	2.51	2.44	2.37	2.30	2.23	2.16	2.09	2.02	
22	3.18	3.01	2.97	2.91	2.84	2.77	2.70	2.63	2.56	2.49	2.42	2.35	2.28	2.21	2.14	2.07	
23	3.23	3.06	3.02	2.96	2.89	2.82	2.75	2.68	2.61	2.54	2.47	2.40	2.33	2.26	2.19	2.12	
24	3.28	3.11	3.07	3.01	2.94	2.87	2.80	2.73	2.66	2.59	2.52	2.45	2.38	2.31	2.24	2.17	
25	3.33	3.16	3.12	3.06	2.99	2.92	2.85	2.78	2.71	2.64	2.57	2.50	2.43	2.36	2.29	2.22	
26	3.38	3.21	3.17	3.11	3.04	2.97	2.90	2.83	2.76	2.69	2.62	2.55	2.48	2.41	2.34	2.27	
27	3.43	3.26	3.22	3.16	3.09	3.02	2.95	2.88	2.81	2.74	2.67	2.60	2.53	2.46	2.39	2.32	
28	3.48	3.31	3.27	3.21	3.14	3.07	3.00	2.93	2.86	2.79	2.72	2.65	2.58	2.51	2.44	2.37	
29	3.53	3.36	3.32	3.26	3.19	3.12	3.05	2.98	2.91	2.84	2.77	2.70	2.63	2.56	2.49	2.42	
30	3.58	3.41	3.37	3.31	3.24	3.17	3.10	3.03	2.96	2.89	2.82	2.75	2.68	2.61	2.54	2.47	
31	3.63	3.46	3.42	3.36	3.29	3.22	3.15	3.08	3.01	2.94	2.87	2.80	2.73	2.66	2.59	2.52	
32	3.68	3.51	3.47	3.41	3.34	3.27	3.20	3.13	3.06	2.99	2.92	2.85	2.78	2.71	2.64	2.57	
33	3.73	3.56	3.52	3.46	3.39	3.32	3.25	3.18	3.11	3.04	2.97	2.90	2.83	2.76	2.69	2.62	
34	3.78	3.61	3.57	3.51	3.44	3.37	3.30	3.23	3.16	3.09	3.02	2.95	2.88	2.81	2.74	2.67	
35	3.83	3.66	3.62	3.56	3.49	3.42	3.35	3.28	3.21	3.14	3.07	3.00	2.93	2.86	2.79	2.72	
36	3.88	3.71	3.67	3.61	3.54	3.47	3.40	3.33	3.26	3.19	3.12	3.05	2.98	2.91	2.84	2.77	
37	3.93	3.76	3.72	3.66	3.59	3.52	3.45	3.38	3.31	3.24	3.17	3.10	3.03	2.96	2.89	2.82	
38	3.98	3.81	3.77	3.71	3.64	3.57	3.50	3.43	3.36	3.29	3.22	3.15	3.08	3.01	2.94	2.87	
39	4.03	3.86	3.82	3.76	3.69	3.62	3.55	3.48	3.41	3.34	3.27	3.20	3.13	3.06	2.99	2.92	
40	4.08	3.91	3.87	3.81	3.74	3.67	3.60	3.53	3.46	3.39	3.32	3.25	3.18	3.11	3.04	2.97	
41	4.13	3.96	3.92	3.86	3.79	3.72	3.65	3.58	3.51	3.44	3.37	3.30	3.23	3.16	3.09	3.02	
42	4.18	4.01	3.97	3.91	3.84	3.77	3.70	3.63	3.56	3.49	3.42	3.35	3.28	3.21	3.14	3.07	
43	4.23	4.06	4.02	3.96	3.89	3.82	3.75	3.68	3.61	3.54	3.47	3.40	3.33	3.26	3.19	3.12	
44	4.28	4.11	4.07	4.01	3.94	3.87	3.80	3.73	3.66	3.59	3.52	3.45	3.38	3.31	3.24	3.17	
45	4.33	4.16	4.12	4.06	3.99	3.92	3.85	3.78	3.71	3.64	3.57	3.50	3.43	3.36	3.29	3.22	
46	4.38	4.21	4.17	4.11	4.04	3.97	3.90	3.83	3.76	3.69	3.62	3.55	3.48	3.41	3.34	3.27	
47	4.43	4.26	4.22	4.16	4.09	4.02	3.95	3.88	3.81	3.74	3.67	3.60	3.53	3.46	3.39	3.32	
48	4.48	4.31	4.27	4.21	4.14	4.07	4.00	3.93	3.86	3.79	3.72	3.65	3.58	3.51	3.44	3.37	
49	4.53	4.36	4.32	4.26	4.19	4.12	4.05	3.98	3.91	3.84	3.77	3.70	3.63	3.56	3.49	3.42	
50	4.58	4.41	4.37	4.31	4.24	4.17	4.10	4.03	3.96	3.89	3.82	3.75	3.68	3.61	3.54	3.47	

표본 크기 n = 100

표본 크기 n	확률인수 h									
	1	2	3	4	5	6	8	10	15	25
5	1.92									
6	2.07	1.93								
7	2.19	2.04	1.94							
8	2.28	2.15	2.10	1.94						
9	2.34	2.19	2.11	2.10	1.95					
10	2.42	2.27	2.21	2.22	2.11	1.93				
12	2.57	2.49	2.46	2.39	2.32	2.24	1.94			
14	2.67	2.64	2.65	2.61	2.57	2.43	2.25	1.94		
16	2.68	2.64	2.63	2.60	2.57	2.53	2.43	2.24		
18	2.71	2.72	2.70	2.68	2.65	2.62	2.63	2.44		
20	2.74	2.77	2.76	2.74	2.72	2.70	2.64	2.63	2.16	
25	2.85	2.88	2.87	2.84	2.84	2.83	2.80	2.74	2.60	
30	2.94	2.94	2.93	2.94	2.94	2.93	2.90	2.88	2.79	2.17
35	3.00	3.02	3.02	3.01	3.00	3.00	3.00	2.97	2.91	2.44
40	3.04	3.04	3.07	3.07	3.04	3.04	3.04	3.02	3.00	2.84
45	3.13	3.12	3.12	3.12	3.11	3.11	3.10	3.09	3.04	3.04
50	3.17	3.18	3.14	3.14	3.16	3.15	3.14	3.14	3.11	3.04
60	3.23	3.23	3.23	3.22	3.22	3.22	3.22	3.21	3.20	3.16
70	3.29	3.29	3.28	3.28	3.28	3.28	3.27	3.27	3.24	3.23
80	3.31	3.33	3.33	3.33	3.33	3.33	3.32	3.32	3.31	3.29
90	3.37	3.37	3.37	3.37	3.37	3.37	3.36	3.34	3.34	3.34
100	3.41	3.41	3.42	3.40	3.40	3.40	3.40	3.39	3.39	3.38

표본 크기 n = 100

표본 크기 n	확률인수 h									
	1	2	3	4	5	6	8	10	15	25
5	1.92									
6	2.17	1.99								
7	2.32	2.17	1.99							
8	2.44	2.32	2.18	1.99						
9	2.54	2.44	2.31	2.18	1.99					
10	2.62	2.55	2.45	2.33	2.18	1.99				
12	2.76	2.70	2.64	2.60	2.48	2.34	1.99			
14	2.84	2.81	2.76	2.75	2.65	2.57	2.36	1.99		
16	2.91	2.90	2.84	2.84	2.75	2.72	2.60	2.36		
18	3.02	3.00	2.97	2.94	2.90	2.85	2.74	2.60		
20	3.08	3.06	3.04	3.01	2.98	2.95	2.87	2.74	2.70	
25	3.21	3.19	3.18	3.16	3.14	3.12	3.07	3.01	2.78	
30	3.30	3.29	3.28	3.26	3.25	3.24	3.21	3.17	3.04	2.21
35	3.37	3.36	3.35	3.34	3.34	3.33	3.30	3.28	3.19	2.81
40	3.43	3.42	3.42	3.41	3.40	3.40	3.39	3.34	3.30	3.05
45	3.48	3.47	3.47	3.46	3.46	3.45	3.44	3.43	3.33	3.22
50	3.57	3.57	3.51	3.51	3.51	3.46	3.49	3.48	3.48	3.34
60	3.60	3.59	3.59	3.59	3.58	3.58	3.57	3.56	3.54	3.48
70	3.66	3.65	3.65	3.65	3.64	3.64	3.64	3.63	3.61	3.57
80	3.70	3.70	3.70	3.70	3.69	3.69	3.69	3.68	3.67	3.64
90	3.74	3.74	3.74	3.74	3.74	3.74	3.73	3.73	3.72	3.70
100	3.78	3.78	3.78	3.77	3.77	3.77	3.77	3.77	3.74	3.74

참 고 문 헌

1. Andrew, D.F. "Significance tests based on residuals." *Biometrika*. 58, pp139-148. 1971.
2. Andrews, D.F. and Pregibon, D. "Finding the outliers that matter." *Journal of the Royal Statistical Society*. B40, pp85-93. 1978.
3. Behnken, D.W., and Draper, N.R. "Residuals and their variance pattern." *Technometrics*. 14, pp101-111. 1972.
4. Bersley, D.A., Kuh, E. and Welsch, R.E. "Regression diagnostics." New York. John Wiley. 1983.
5. Cook, R.D. "Detection of influential observation in linear regression." *Technometrics*. 19, pp15-18. 1977.

6. Draper, N.R. and John, J.A. "Influential observations and outliers in regression." *Technometrics*.23, pp21-26. 1981.
7. Hoaglin, D.C. and Welsch, R.E. "The hat matrix in regression and ANOVA." *American Statistician*.32, ppi7-22. 1978.
8. Lund, R.E. "Tables for an approximate test for outliers in linear models." *Technometrics*. 17, pp473-476.1975.
9. Prescott, P. "An approximate test for outliers in linear models." *Technometrics*.17, pp 129-132. 1975.
10. Stefansky, W. "Rejecting outliers by maximum normed residual." *Ann. Statist.* 42, pp35-45. 1971.
11. Thompson, W.R. "On a criteria for the rejection of outliers and the distribution of the ratio of the deviation to the sample standard deviation." *Ann. Math. Statist.* 6, pp214-219. 1935.
12. Tietjen, G.L., Moore, R.H. and Beckmin, R.J. "Testing for a single outlier in simple linear regression." *Technometrics*.15, pp717-721. 1973.