

T法における欠測データの活用に関する研究

Study of T-Method for Missing Data

有限会社 増田技術事務所 増田 雪也

1. はじめに

T法¹⁾は誰でも簡単に重回帰分析ができる手法である。数式が単純なため、MS-Excelでも容易に解析を行うことができ、様々な分野で活用され始めている。

あまり知られてはいないが、T法には従来の重回帰分析にはないメリットがある。それは、「欠測データが存在していても解析できる」というメリットである。しかし、欠測データ数が、T法における重み付けのSN比 η に悪影響を与える可能性があると考えられる。

そこで本研究では、重み付けのSN比 η を求める際、エネルギー比型SN比²⁾を用いて解析する方法を検討した。その結果、従来のSN比との差は非常に小さく、推定精度に与える影響は軽微であった。

2. 欠測データの活用方法

2.1 事例「我が家の電気使用量の解析」

T法における欠測データの活用方法を説明する。その事例として、「我が家の電気使用量の解析」を取りあげる。この事例は、図1に示すように、信号として「電気使用量（※月毎に日数が異なるため、毎月の電気使用量から1日当たりの平均使用量を求めたもの）」、測定項目として、月毎の「平均気温、平均最高気温、平均最低気温、降水量、平均風速、日照量、ガス使用量、旅行で留守にした日数」を設定した。気象データは、気象庁のWebサイトから得た。

図2に2004年9月から2011年6月までの電気使用量を示す。冬季になると電気使用量が高くなるのは、暖房としてエアコンを使用しているためである。我が家は、長野県の高原にあるため、夏期にエアコンを使用することはない。

この全データの内、2004年9月から2005年8月までの12ヶ月分のデータを、既知データとして用い、残りの2005年9月から2011年6月までの70ヶ月分のデータを、未知データとして推定することとした。

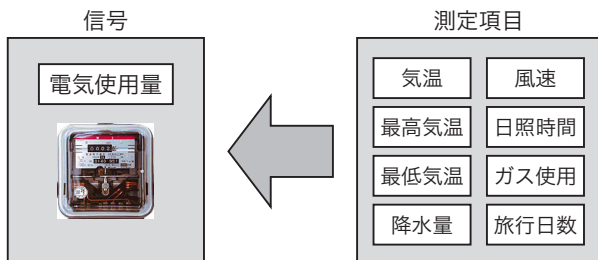


図1 我が家の電気使用量の解析（信号と測定項目）

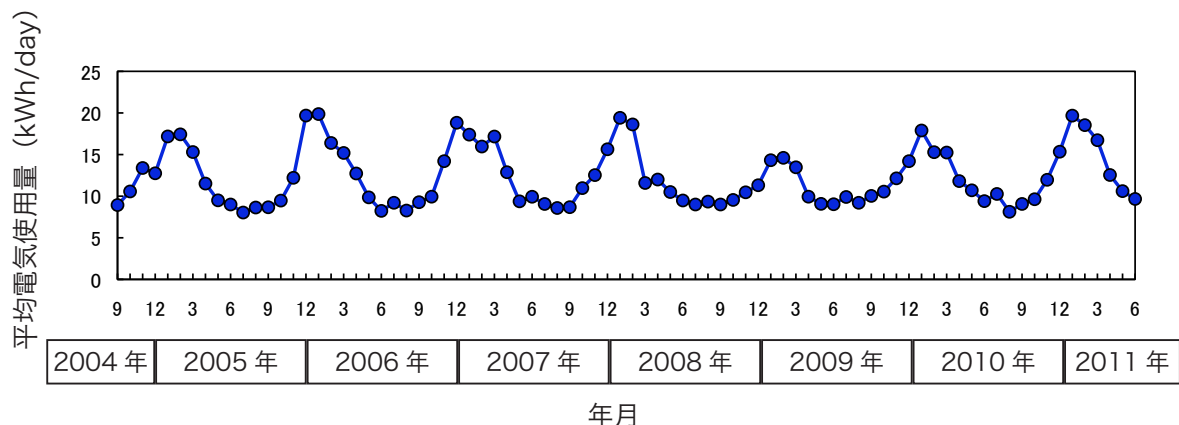
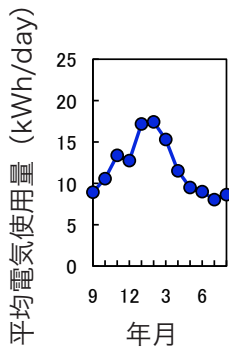


図2 2004年から2011年までの電気使用量



年月	真値	測定項目							
	M	k1	k2	k3	k4	k5	k6	k7	k8
年月	電気使用量	ガス	降水量	気温	最高気温	最低気温	風速	日照量	旅行日数
2004.9	8.93	8.5	5.1	18.8	24.2	14.9	1.6	4.7	4
2004.10	10.57	9.2	9.3	10.9	16.3	6.5	1.6	4.5	3
2004.11	13.39	12.3	1.8	6.8	13.0	1.8	1.6	6.0	3
2004.12	12.75	6.6	2.2	1.1	7.2	-3.8	1.8	6.3	18
2005.1	17.18	13.2	1.3	-3.7	1.5	-9.0	1.8	5.4	5
2005.2	17.43	14.2	1.1	-2.7	2.1	-7.5	2.0	5.6	0
2005.3	15.30	15	2.8	0.9	6.5	-4.7	2.3	6.1	2
2005.4	11.52	15	1.5	8.8	16.1	1.7	2.2	7.3	0
2005.5	9.50	9	1.8	12.5	19.6	6.3	2.0	6.9	7
2005.6	9.00	10.6	3.9	18.7	24.7	14.0	1.5	4.9	0
2005.7	8.03	8.2	5.4	21.1	26.5	16.7	1.4	4.4	7
2005.8	8.64	8.5	2.0	21.9	27.7	18.1	1.3	5.0	2

図3 T法の推定式を作るための既知データ

図3に、T法の推定式を作るための既知データ(12ヶ月分)を示す。

2.2 重回帰分析における欠測データへの対応

図4に示すように、従来の重回帰分析では、欠測値が存在する場合は、欠測値を含む行(つまり、当該月のデータ)を丸ごと削除して解析する。あるいは

年月	真値	測定項目							
	M	k1	k2	k3	k4	k5	k6	k7	k8
年月	電気使用量	ガス	降水量	気温	最高気温	最低気温	風速	日照量	旅行日数
2004.9	8.93	8.5	5.1	18.8	24.2	14.9	1.6	4.7	4
2004.10	10.57	9.2	9.3	10.9	16.3	6.5	1.6	4.5	3
2004.11	13.39	12.3	1.8	6.8	13.0	1.8	1.6	6.0	3
2004.12	12.75	6.6	2.2	1.1	7.2	-3.8	1.8	6.3	18
2005.1	17.18	13.2	1.3	-3.7	1.5	-9.0	1.8	5.4	5
2005.2	17.43	14.2	1.1	-2.7	2.1	-7.5	2.0	5.6	0
2005.3	15.30	15	2.8	0.9	6.5	-4.7	2.3	6.1	2
2005.4	11.52	15	1.5	8.8	16.1	1.7	2.2	7.3	0
2005.5	9.50	9	1.8	12.5	19.6	6.3	2.0	6.9	7
2005.6	9.00	10.6	3.9	18.7	24.7	14.0	1.5	4.9	0
2005.7	8.03	8.2	5.4	21.1	26.5	16.7	1.4	4.4	7
2005.8	8.64	8.5	2.0	21.9	27.7	18.1	1.3	5.0	2

図4 重回帰分析における欠測データへの対応

は、列方向(つまり測定項目)の欠測値が多い場合は、その列の測定項目を丸ごと削除して解析する。

図5に示すのは、極端な事例であるが、斜めに欠測値が存在する場合や、市松模様のような場合は、従来の重回帰分析では解析が難しい。

2.3 T法における欠測データへの対応

T法の特徴は、「SN比による重み付け」をすることである。図6に、そのSN比の概念を示す。信号と項目の関係に対してゼロ点を通る直線を設定し、この直線からのズレを数値化したものが『ゼロ点比例式のSN比』である。T法では、このSN比の値を「重み付け」として用いている。SN比の詳細については、参考文献¹⁾を参照されたい。

図7の3つのグラフは、信号と測定項目の関係における代表的な3つのパターンである。このように、バラツキ具合をSN比によって重み付けしている。

年月	真値	測定項目							
	M	k1	k2	k3	k4	k5	k6	k7	k8
年月	電気使用量	ガス	降水量	気温	最高気温	最低気温	風速	日照量	旅行日数
2004.9	8.93	8.5	5.1	18.8	24.2	14.9	1.6	4.7	4
2004.10	10.57	9.2	9.3	10.9	16.3	6.5	1.6	4.5	3
2004.11	13.39	12.3	1.8	6.8	13.0	1.8	1.6	6.0	3
2004.12	12.75	6.6	2.2	1.1	7.2	-3.8	1.8	6.3	18
2005.1	17.18	13.2	1.3	-3.7	1.5	-9.0	1.8	5.4	5
2005.2	17.43	14.2	1.1	-2.7	2.1	-7.5	2.0	5.6	0
2005.3	15.30	15	2.8	0.9	6.5	-4.7	2.3	6.1	2
2005.4	11.52	15	1.5	8.8	16.1	1.7	2.2	7.3	0
2005.5	9.50	9	1.8	12.5	19.6	6.3	2.0	6.9	7
2005.6	9.00	10.6	3.9	18.7	24.7	14.0	1.5	4.9	0
2005.7	8.03	8.2	5.4	21.1	26.5	16.7	1.4	4.4	7
2005.8	8.64	8.5	2.0	21.9	27.7	18.1	1.3	5.0	2

解析は難しい

年月	真値	測定項目							
	M	k1	k2	k3	k4	k5	k6	k7	k8
年月	電気使用量	ガス	降水量	気温	最高気温	最低気温	風速	日照量	旅行日数
2004.9	8.93	8.5	5.1	18.8	24.2	14.9	1.6	4.7	4
2004.10	10.57	9.2	9.3	10.9	16.3	6.5	1.6	4.5	3
2004.11	13.39	12.3	1.8	6.8	13.0	1.8	1.6	6.0	3
2004.12	12.75	6.6	2.2	1.1	7.2	-3.8	1.8	6.3	18
2005.1	17.18	13.2	1.3	-3.7	1.5	-9.0	1.8	5.4	5
2005.2	17.43	14.2	1.1	-2.7	2.1	-7.5	2.0	5.6	0
2005.3	15.30	15	2.8	0.9	6.5	-4.7	2.3	6.1	2
2005.4	11.52	15	1.5	8.8	16.1	1.7	2.2	7.3	0
2005.5	9.50	9	1.8	12.5	19.6	6.3	2.0	6.9	7
2005.6	9.00	10.6	3.9	18.7	24.7	14.0	1.5	4.9	0
2005.7	8.03	8.2	5.4	21.1	26.5	16.7	1.4	4.4	7
2005.8	8.64	8.5	2.0	21.9	27.7	18.1	1.3	5.0	2

解析は難しい

図5 極端な事例

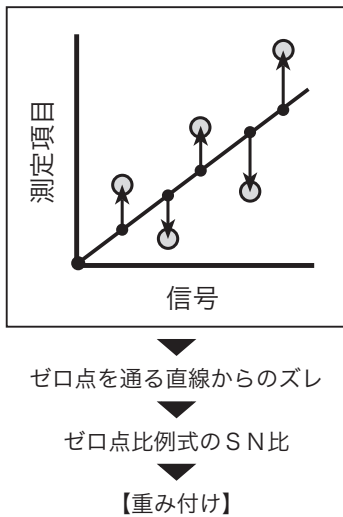


図6 【重み付け】として用いるSN比の概念

図8の左のグラフは、欠測値が無い場合である。一方、右のグラフは、左のグラフから5つのデータを欠測させた場合である。このケースでは、例えば5つの欠測が発生したとしても、全体的なバラツキの大きさは、ほぼ同程度であり、ゼロ点比例式で求められる重み付けのSN比 η も、ほぼ同程度となる。データが欠測しても、データ数が少なくなるだけで、SN比の算出に数学的な問題は無い。従ってT法では、欠測値が存在しても、重み付けのSN比を求めることができ、推定値を得ることが可能となる。ただし、推定精度の良し悪しは、ケースバイケースである。

図9に示すように、3つの測定項目において、それぞれ欠測数が異なったとしても、重み付けのSN比を求めることができ、各SN比 η と各傾き β から、推定値 \hat{M}_i を求めることができる。なお、図9の数式は、一見複雑そうに見えるが、各測定項目の重み付けのSN比で加重平均を求めているにすぎない。

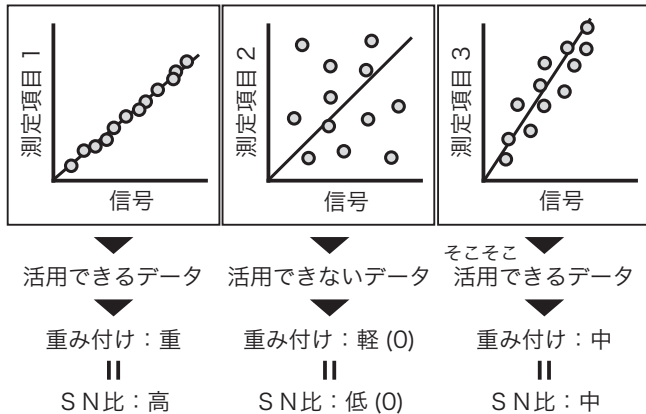


図7 T法の特徴【重み付け】

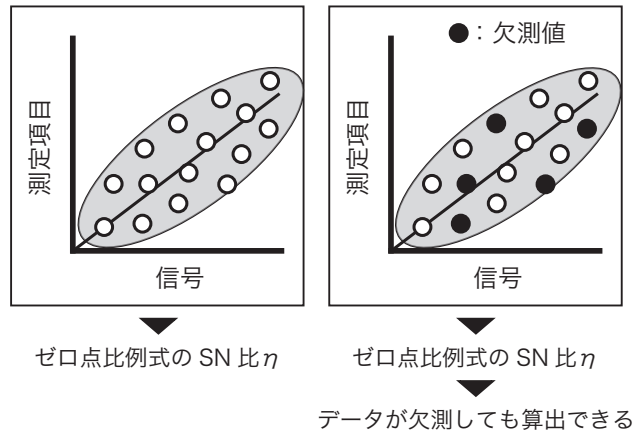


図8 欠測データが存在する場合の重み付け

2.4 故意にデータを欠測させた場合の比較

図10に、欠測値の有無および欠測の程度の違いによる、総合推定のSN比および相関係数の違いを示す。左は欠測が無い場合である。中央は故意に斜めに欠測させた場合、右は市松模様欠測させた場合である。この「我が家の電気使用量の解析」事例では、欠測値の有無および程度の違いによらず、総合推定のSN比および相関係数は、ほぼ同じ値となった。

3. 欠測データを活用する上での問題点

3.1 欠測データ数の影響

図9に示すような場合、データ数は測定項目毎に異なっている。式(1)にT法の重み付けのSN比を示す。式(1)中の V_e は、式(2)によって求められる。この V_e は、データ数 l の影響を受ける。従って、測

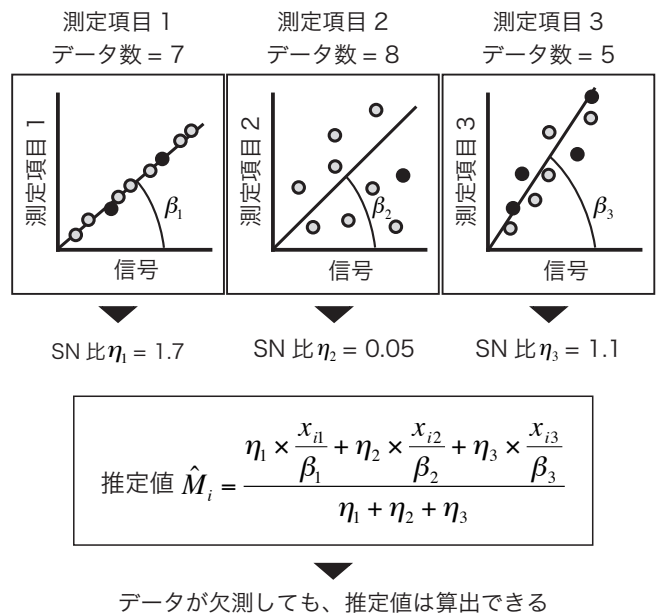


図9 欠測データが存在する場合の推定値の算出

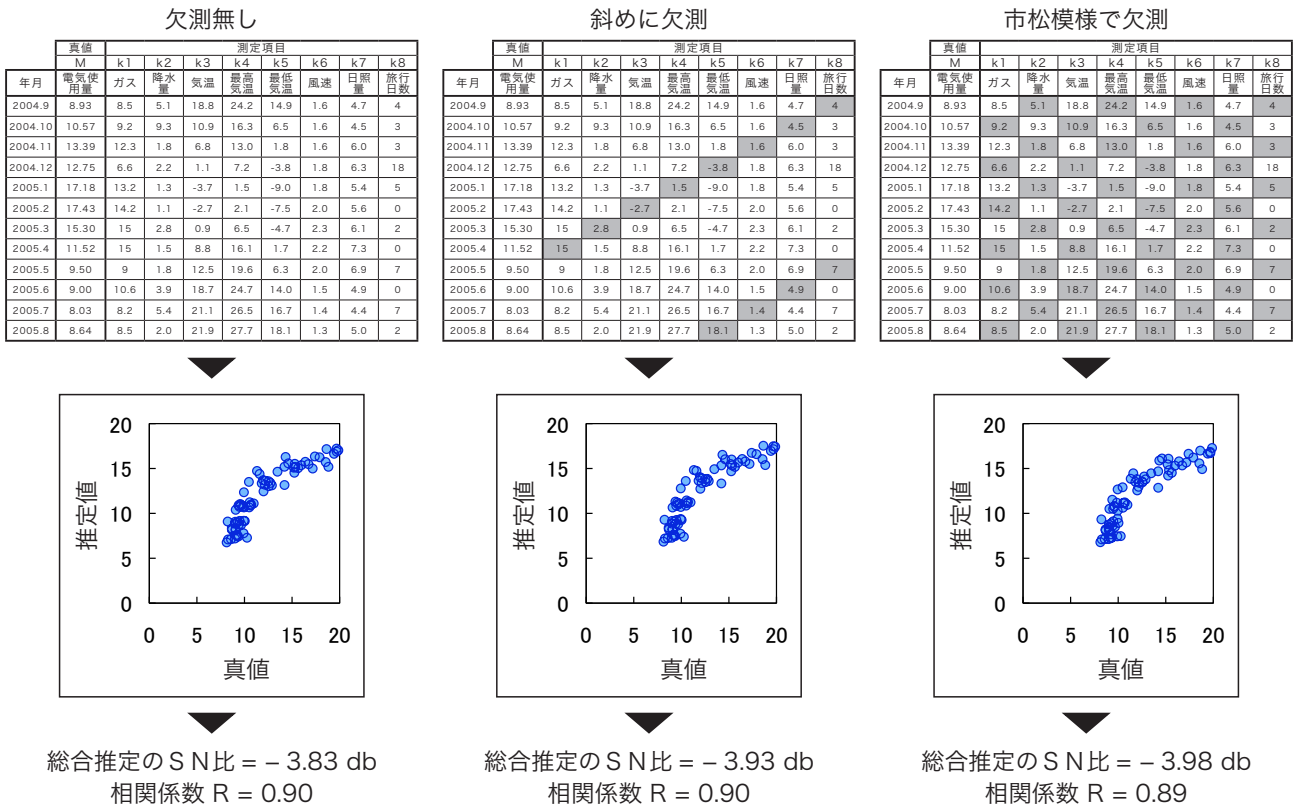


図10 故意にデータを欠測させた場合の比較

$$\eta = \frac{1}{r} \frac{(S_\beta - V_e)}{V_e} \quad (1)$$

$$V_e = \frac{1}{l-1} (S_T - S_\beta) \quad (2)$$

定項目毎にデータ数が異なる場合、データ数が重み付けのSN比ηに影響を与えられられる。このことが、T法において欠測データを扱う上での問題点となる。つまり、推定精度（総合推定のSN比、あるいは相関係数）を悪化させる原因になるのではないかという懸念である。そこで、電気使用量の事例を用いて、その影響度合いを検討してみたことにした。

3.2 電気使用量の事例で検証

図11で検証の手順を説明する。まず電気使用量の4つの測定項目（ガス使用量、降水量、平均気温、平均最高気温）について、それぞれ全く同じデータを用意し、項目数を2倍にする。次に、12ヶ月分のデータについても、それぞれ全く同じデータを用意し、データ数を2倍（24ヶ月分）とする。ここで、測定項目のガス使用量について詳しく説明すると、

年月	測定項目								
	M	k1	k1'	k2	k2'	k3	k3'	k4	k4'
2004.9	8.93	8.5	8.5	5.1	5.1	18.8	18.8	24.2	24.2
2004.10	10.57	9.2	9.2	9.3	9.3	10.9	10.9	16.3	16.3
2004.11	13.39	12.3	12.3	1.8	1.8	6.8	6.8	13.0	13.0
2004.12	12.75	6.6	6.6	2.2	2.2	1.1	1.1	7.2	7.2
2005.1	17.18	13.2	13.2	1.3	1.3	-3.7	-3.7	1.5	1.5
2005.2	17.43	14.2	14.2	1.1	1.1	-2.7	-2.7	2.1	2.1
2005.3	15.30	15	15	2.8	2.8	0.9	0.9	6.5	6.5
2005.4	11.52	15	15	1.5	1.5	8.8	8.8	16.1	16.1
2005.5	9.50	9	9	1.8	1.8	12.5	12.5	19.6	19.6
2005.6	9.00	10.6	10.6	3.9	3.9	18.7	18.7	24.7	24.7
2005.7	8.03	8.2	8.2	5.4	5.4	21.1	21.1	26.5	26.5
2005.8	8.64	8.5	8.5	2.0	2.0	21.9	21.9	27.7	27.7

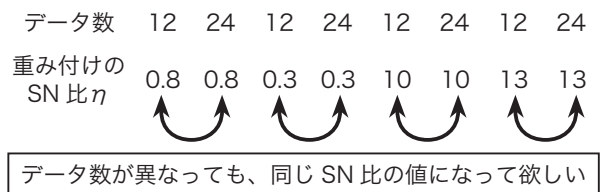


図11 電気使用量の事例で検証

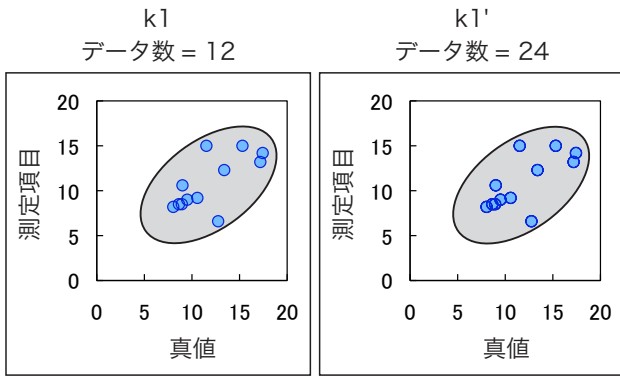


図12 k1とk1'をグラフで説明

k₁およびk_{1'}は、共にデータ数=24の同じデータであるが、k₁の下半分(灰色:12ヶ月分)のデータを欠測値として設定することにする。

降水量k₂およびk_{2'}、平均気温k₃およびk_{3'}、平均最高気温k₄およびk_{4'}についても、同様の設定をする。

k₁およびk_{1'}は、全く同じデータをダブらせたものであるから、そのバラツキの大きさ、つまり重み付けのSN比ηは、同じ値になるのが理想である。

説明がわかりにくいかもしれないので、図12で同じことを説明する。

左のk₁に対して、右のk_{1'}はk₁のデータを2倍にしてダブらせたものである。k_{1'}のグラフデータは、重なって見えるが、データ数はk₁の2倍ある。両者のバラツキの大きさは全く同じなので、SN比の値は同じになるのが理想である。

図13に、従来のゼロ点比例式で算出したk_xおよびk_{x'}の重み付けのSN比ηを示す。全く同じバラツキにもかかわらず、SN比ηは異なる値となっている。これはデータ数がSN比ηに影響を与えているからである。この差が、推定精度(総合推定のSN比および相関係数)に悪影響を与える可能性がある。

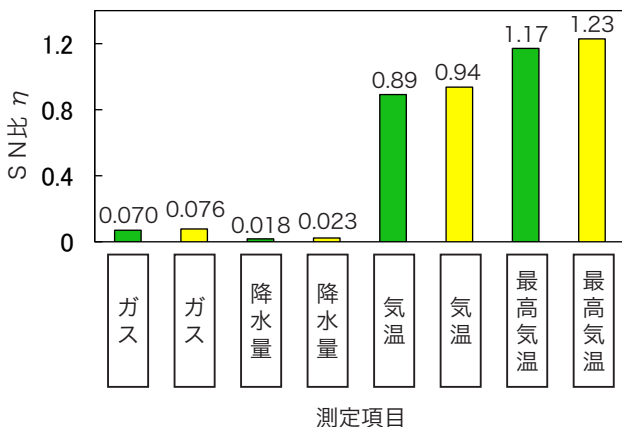


図13 従来のゼロ点比例式での検証結果

4. 欠測データ数に依存しないSN比の検討

4.1 エネルギー比型SN比で検証

欠測データ数が重み付けのSN比ηに与える影響を無くすため、清水らが提案するエネルギー比型SN比²⁾を採用してみることにした。

$$\eta_E = \frac{S_\beta}{S_\epsilon} \quad (3)$$

式(3)に示すように、エネルギー比型SN比η_Eは、データ数lに依存しない。よって、今回の欠測データにおけるT法で採用するには好都合といえる。

図14に、エネルギー比型SN比での検証結果を示す。エネルギー比型SN比で算出したk_xおよびk_{x'}の重み付けのSN比η_Eは、全く同じ値となった。

エネルギー比型SN比をT法に採用した場合、どれくらい推定精度が向上するのかを確かめるため、2つの事例で検証することにする。

4.2 適用事例その1「我が家の電気使用量の解析」

図15に、左上のデータを故意に欠測した場合の推定結果を示す。従来のSN比とエネルギー比型SN比を比較すると、ほぼ同じ推定精度となった。また、図16に示すように、右下のデータを故意に欠測した場合も、ほぼ同じ推定精度となった。

4.3 適用事例その2「世界各国の出生率の解析」

アメリカ中央情報局(CIA)のWebサイトからデータを引用し、世界各国の出生率の解析を行った。既知データとして、15カ国(Afghanistan, Albania, Algeriaなど)のデータを用い、未知データとして133カ国の出生率を推定した。

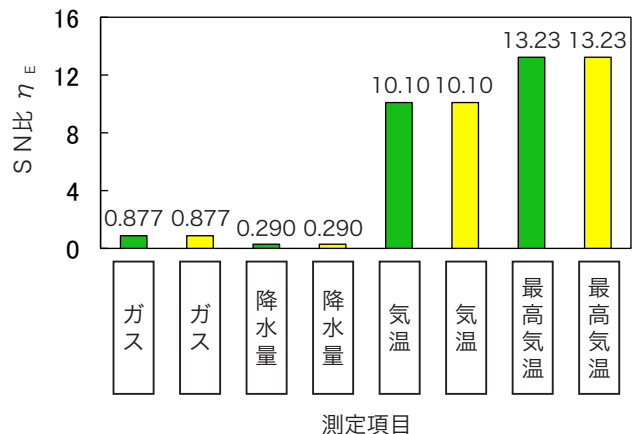


図14 エネルギー比型SN比での検証結果

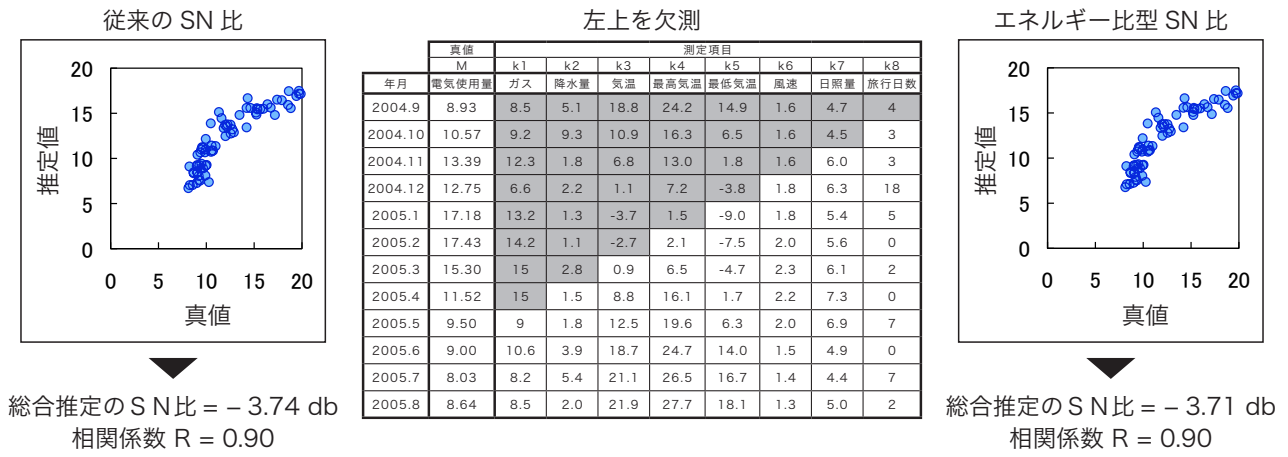


図 15 電気使用量 (故意に左上を欠測)

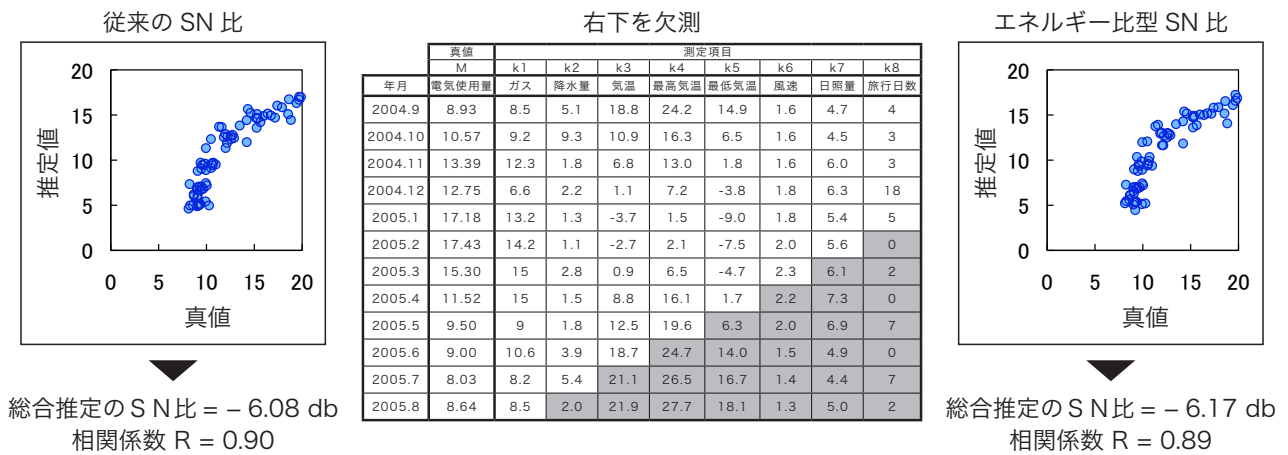


図 16 電気使用量 (故意に右下を欠測)

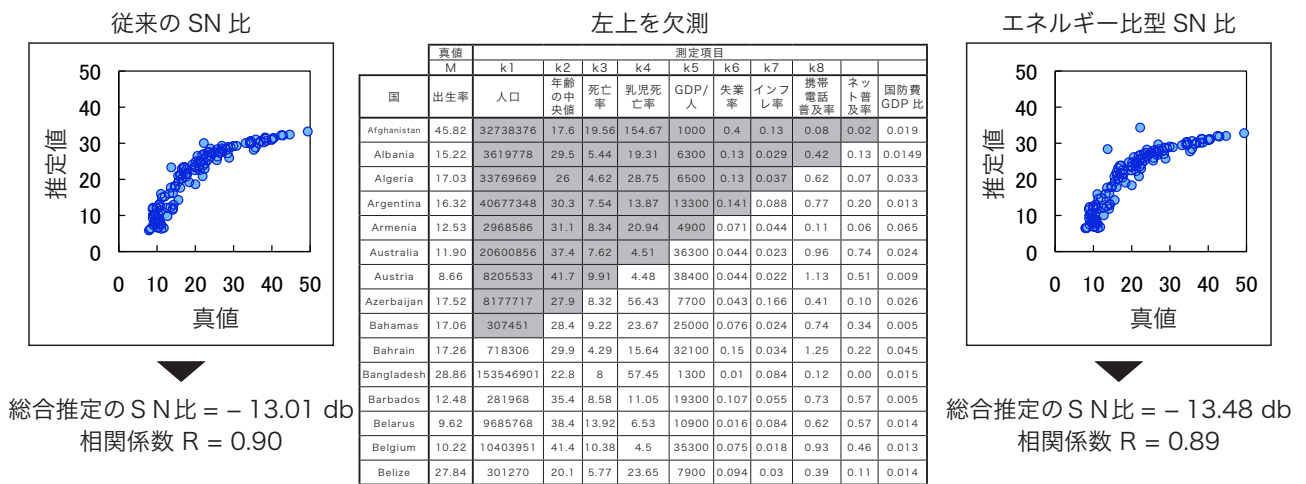


図 17 出生率 (故意に左上を欠測)

図 17 に示すように、左上のデータを故意に欠測させた場合の推定結果を示す。従来の SN 比とエネルギー比型 SN 比を比較すると、ほぼ同じ推定精度となった。また、図 18 に示すように、右下のデータを故意に欠測させた場合も、ほぼ同じ推定精度となった。

5. さいごに

今回の研究では、T 法における欠測データの活用について、重み付けの SN 比の検討を行った。その結果、従来の SN 比とエネルギー比型 SN 比の差は小さく、推定精度に与える影響は軽微であった。

今回の研究で示したように、T 法は欠測データが存在する場合でも、問題なく解析できる便利な手法



図18 出生率 (故意に右下を欠測)

であり、従来の重回帰分析にはないメリットを有している。今後、技術者を支援するツールのひとつとして、T法を普及していきたいと考えている。

謝 辞

T法における欠測データの活用方法をご教示いただいた(財)長野県テクノ財団の岩下幸廣氏に、感謝の意を表す。

■参考文献

- 1) 田口玄一：目的機能と基本機能(6), 品質工学, Vol.13, No.3, pp.309-314, 2005
- 2) 清水豊他：新SN比の研究(4), 第16回品質工学会研究発表大会論文集, (2008), pp.422-425