

# 論文引用に影響を与える要因：負の二項重回帰による検討

芳鐘冬樹<sup>†</sup> 辻慶太 小野寺夏生  
筑波大学図書館情報メディア研究科  
<sup>†</sup> fuyuki@slis.tsukuba.ac.jp

抄録

論文の被引用数を論文の計量書誌学的要因から説明するモデルとして、負の2項重回帰モデルと線形回帰モデルを設定し、説明要因の有効性を調べるとともに、両モデルの当てはまりのよさを比較した。自然科学系6分野のデータを対象に分析した結果、どちらのモデルにおいてもPrice指数が全分野で有効であること、負の2項重回帰モデルが線形回帰モデルよりも観測値に対する当てはまりがよいことが示された。

## 1. はじめに

研究評価を行う際の参考データとして、論文の被引用数に基づく指標がよく用いられる。被引用数は論文の影響度を示す統計的測度としては適切であると考えられるが、これを研究評価に用いるには、その特性をよく知り、十分な注意を払う必要がある<sup>1)</sup>。たとえば、論文が掲載された雑誌のインパクトファクターを論文の評価指標に転用する例がしばしば見られるが、これは極めて不適切な使用例である<sup>2)</sup>。

更に、単純に論文間の被引用数を比較することには問題がある。論文の被引用数は、分野、掲載誌の発行国、論文の種類(原著論文、短報、レビュー等)、使用言語等、多くの要因に影響されるからである。

論文の被引用数と様々な要因との関係を検討した研究は多数報告されている。しかし、それらの結論は必ずしも一致していない。その理由の1つは、これらの研究の多くがある単一の要因に着目しているか、複数の要因をそれぞれ独立に見ていることにある。従って、ある要因が被引用数と相関があっても、それが他の要因の交絡因子である可能性を否定できない。

重回帰分析は、個々の説明変数による効果を分離できるので、考え得る多くの要因の中のどれが被引用数に対する説明力が高いかを推察することができる。この方法を用いた研究はいくつかあり<sup>3-10)</sup>、それぞれ興味深い結果が得られている。しかし、これらの研究も、特定の分野を対象にしている、研究目的が被引用数に影響を与える要因の解明ではないなどの点から、汎用的な結論を出しているとは言えない。

筆者らは、複数の分野とそれらの中の複数の雑誌について、同一年に発表された原著論文を採り、比較的長期間(6年間)の被引用数を目的変数とし、想定される多くの潜在要因を説明変数として、重回帰分析を進めている。これによって、被引用数に対する各要因の寄与を検討し、論文の被引用数に対するより精密なベースラインを与えることを目標としている。6つの分野について、9つの要因を考慮した中間的な結果を昨年報告した<sup>11)</sup>。Price指数、参考文献数、論文長が被引用数に対する有効な説明要因であり、分野を越えてある程度一貫性のある予測モデルが得られる可能性を示した。

本研究は、これに以下の2つの改善を行ったものである：説明変数に被引用論文著者の過去の業績(論文数、被引用数、活動期間)を加える、非線形回帰の1つである負の2項重回帰(NBR)モデルを採用する。目的変数である被引用数は、非負の値をとる整数であり、その分布が極めて歪度が大きいことから、線形回帰(LR)モデルよりNBRモデルが適していると考えられ、いくつかの先例<sup>7-9)</sup>がある。

## 2. データ

本研究では物性物理学、無機・核化学、電気・電子工学、生化学及び分子生物学、生理学、消化器病医学の6分野から各4誌(すべて英文誌)を選び、それぞれから2000年発表の原著論文50-60件(計1,395件)を無作為抽出して、調査対象論文とした。手順としてはまず2006年7月に、ThomsonのWeb of Scienceから、24誌の2000年における原著論文(articles)のデータ

表1 説明変数

	変数名	変数略号	変数の説明
第1著者の属性	第1著者発表実績	Publ_Prim	2000年までの第1著者の論文数
	第1著者活動期間	Age_Prim	第1著者の最初の論文発表から2000年までの期間
	第1著者引用実績	Impact_Prim	Publ_Primの全論文の2000年までの被引用数合計(自己引用を除く)
著者集合属性	著者数	Authors	対象論文の著者数
	所属機関数	Insts	Authorsの所属機関数(同一機関内の異なる部門は1機関とする)
	所属国数	Countries	Instsが属する国の数
	*最高著者発表実績	Publ_Max	2000年までの最高論文数共著者の論文数
	*最長活動期間	Age_Max	2000年までの活動期間最長の共著者の活動期間
	*最高著者引用実績	Impact_Max	2000年までに最高の引用を得た共著者の被引用数合計(自己引用を除く)
論文の属性	論文長	Length	対象論文のページ数を紙面係数で補正
	図数	Figures	対象論文中の図の数
	表数	Tables	対象論文中の表の数
	数式数	Eqs	対象論文中の数式の数
	参考文献数	Refs	対象論文中の参考文献数
	Price指数	Price	Refsのうち最近5年以内のもの(1996以降)の比率
	雑誌	1A~4A	掲載誌を示すダミー変数(0か1)

をダウンロードした後、各誌から 50-60 論文を無作為抽出して調査用標本とした(2 ページ以下の論文、会議プロシーディング論文は除外した)。これらの標本論文の書誌データと、その参考文献及び引用文献のデータを、Thomson から購入し、これから後述の説明変数及び目的変数のデータを取得した。図の数、表の数、数式の数、別途原論文に当たって取得した。

### 3. 方法

NBR モデルに基づく回帰分析では、個体  $i$  に対する目的変数  $y_i$  (論文  $i$  の被引用数) が負の 2 項分布:

$$\Pr(y_i = k) = \frac{\binom{k + \mu_i}{k}}{\binom{k + \mu_i}{k}} \left( \frac{\mu_i}{\mu_i + 1} \right)^k$$

に従うとし、その期待値  $\mu_i$  を次式で予測する:

$$\ln(\mu_i) = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip}$$

ここで、 $\beta_0, \beta_1, \beta_2, \dots, \beta_p$  及び  $\mu_i$  の推定値は、 $\{X_{i1}, X_{i2}, \dots, X_{ip}; y_i\}$  の入力データに基づいて算出される(  $\mu_i$  は  $i$  によらないとする )。

本研究では、以下の目的変数、説明変数を取り上げ、先述の 6 分野について非線形回帰分析(NBR)と線形回帰分析(LR)を行った:

#### [ 目的変数 ]

自己引用も含む被引用数 TC\_Total とした<sup>1</sup>。LR においては  $\ln(\text{TC\_Total}+1)$  に変換した。

#### [ 説明変数 ]

表 1 のものを用いた。「第 1 著者発表実績」「第 1 著者引用実績」「最高著者発表実績」「最高著者引用実績」は、完全計数法による<sup>2</sup>。本研究では基本的に 6 分野各々 4 つの雑誌を調査に用いたが、被引用論文著者の最高過去業績(表 1 中「\*」を付した「最高著者発表実績」「最長活動期間」「最高著者引用実績」)についてはその算出作業量が膨大なため、6 分野各々 1 雑誌のみを対象とした。即ち、本研究では 6 分野それぞれにおいて、(1) 1 誌だけを用い、最高過去業績の 3 変数も含める(ダミー変数「雑誌」は含めない)回帰分析と、(2) 4 誌すべてを用い、最高過去業績は含めない回帰分析の 2 種類を行った。

統計解析ソフトには R を使い、NBR については glm.nb 関数、LR については lm 関数を用いた。両者共に step 関数を用いてステップワイズ法で説明変数の選択を行った。

### 4. 結果

#### 4.1 回帰分析結果の概要

分野ごとの回帰分析結果の概要として、標本数  $n$ 、および NBR、LR 各々の AIC (赤池情報量規準) と LR の  $Rc^2$  (自由度調整済み決定係数) を表 2 に示した。AIC と  $Rc^2$  は共にモデルのよさを評価する尺度(AIC は小さいほど評価が高く  $Rc^2$  は大きいほど評価が高い)であるが、

<sup>1</sup> 自己引用を除いても結果に大きな差異はなかった。

<sup>2</sup> 調節計数法を用いても結果に大きな差異はなかった。

表2 分野ごとの回帰分析結果

分野	対象	n	NBR				LR				
			AIC	相対残差平方和	残差平方和	Rc <sup>2</sup>	AIC	相対残差平方和	残差平方和		
物性物理学	4誌	230	1499.0	<b>181.2</b>	<b>(0.79)</b>	24074.4	0.262	-50.5	<b>468.5</b>	<b>(2.04)</b>	27366.6
	1誌	55	408.2	<b>34.1</b>	<b>(0.62)</b>	11453.9	0.231	-7.2	<b>98.2</b>	<b>(1.79)</b>	12629.5
無機・核化学	4誌	227	1466.9	<b>142.1</b>	<b>(0.63)</b>	18625.9	0.320	-107.5	<b>332.0</b>	<b>(1.46)</b>	20156.8
	1誌	54	354.4	<b>19.3</b>	<b>(0.36)</b>	2696.3	0.326	-30.4	<b>36.5</b>	<b>(0.68)</b>	3044.8
電気・電子工学	4誌	229	1201.8	<b>290.8</b>	<b>(1.27)</b>	13940.7	0.271	-86.4	<b>1068.6</b>	<b>(4.67)</b>	15097.5
	1誌	59	368.0	<b>39.9</b>	<b>(0.68)</b>	3950.5	0.306	-26.9	<b>80.4</b>	<b>(1.36)</b>	3746.2
生化学及び分子生物学	4誌	240	1798.6	<b>104.6</b>	<b>(0.44)</b>	59498.3	0.475	-194.7	<b>172.8</b>	<b>(0.72)</b>	60364.3
	1誌	60	488.7	<b>21.4</b>	<b>(0.36)</b>	18394.6	0.143	-64.8	<b>31.4</b>	<b>(0.52)</b>	19273.6
生理学	4誌	236	1600.9	<b>122.4</b>	<b>(0.52)</b>	28237.3	0.607	-198.1	<b>251.4</b>	<b>(1.07)</b>	29027.7
	1誌	58	444.7	<b>18.0</b>	<b>(0.31)</b>	8196.3	0.285	-61.1	<b>29.2</b>	<b>(0.50)</b>	8848.4
消化器病	4誌	233	1805.8	<b>161.5</b>	<b>(0.69)</b>	144454.1	0.497	-75.9	<b>385.3</b>	<b>(1.65)</b>	152235.8
医学	1誌	59	554.7	<b>34.3</b>	<b>(0.58)</b>	83866.5	0.192	-32.1	<b>93.0</b>	<b>(1.58)</b>	94743.4

どの分野で評価が高いか、また、4誌すべてを用いて最高過去業績の3変数は含めない場合と、1誌だけを用いて最高過去業績も含める場合のどちらが評価が高いかは、NBRとLRとで異なり、かつ同じLRであってもAICとRc<sup>2</sup>とで異なり、一貫した傾向は認められなかった。

ステップワイズ法による変数選択において、6分野のうち2分野以上で選ばれた変数を表3に示した。NBR, LR共に、また、4誌の場合も1誌の場合も共に、Price指数が最もよく選択されており（特に4誌の場合は全分野で選択された）分野を越えて被引用数に対して説明力を持つ変数であると確認できる。一方、被引用論文著者の過去業績（第1著者の業績と全共著者を考慮した最高業績）については、選択された分野数はすべて2以下に留まった。

4誌すべてを用いる場合に関しては、掲載誌を示すダミー変数が多くの分野で選択された。掲載誌も、被引用数に対する有効な説明要因であることが示唆される。

#### 4.2 NBRとLRの当てはまりのよさの比較

前述のRc<sup>2</sup>はNBRに適用できず<sup>3</sup>、AICもNBRとLRという異なるモデル間の比較には必ずしも適さない。そこで、モデルの当てはまりのよさについてNBRとLRを比べる尺度として、相対残差平方和SSを用いることとした。

$$SS = \sum \left( \frac{y_i - \mu_i}{\mu_i} \right)^2$$

相対残差平方和は、目的変数の観測値 $y_i$ と予測値（期待値） $\mu_i$ のずれを示す尺度で、両者の差を期待値 $\mu_i$ で規格化した値の平方和である。NBRにおいても $\log(y)$ を目的変数とするLRにおいても、残差は $\mu_i$ に比例して大きくなりやすいため、より合理的に比較できるよう、 $\mu_i$ で規格化して残差の均等化を図った。表2に各々の相対残差平方和（括弧内はnで規格化した相対残差平方平均）を示した。参考までに、規格化していない残差平方和も併せて示した。

相対残差平方和と相対残差平方平均は、どの分野においても、また対象が4誌の場合も1誌の場合も、NBRの方がLRよりも値が小さい、即ち当てはまりがよいことが示された<sup>4</sup>。

次に、期待値ではなく、分布全体の当てはまりを見るために、被引用数ごとの論文数の観測値と予測値の分布を調べた。NBRに関しては、各論文について被引用数 $\Pr(y=k)$ の確率分布を計算し、それを積算することで予測値の分布を求めた。図1は、物性物理学の4誌を例に、被引用数ごとの分布を示したものである。図から、特に低被引用数においてNBRがLRよりも当てはまりがよいことを見て取れる。

#### 5. まとめ

<sup>3</sup> NBRに適用できるPseudo R<sup>2</sup>という尺度はあるが、一般的でなく、またLRの決定係数と定義が異なるため比較に適さない。

<sup>4</sup> 残差平方和の比較でも、1つの例外（電気電子工学の1誌）を除き、NBRの方が値が小さい。

負の2項重回帰モデルにおいても、線形回帰モデルと同様、Price 指数が被引用数の説明に有効であることが確認できた。本研究で説明変数として新たに加えた、著者の過去業績については、その有効性を明確に示すことはできなかった。負の2項重回帰モデルと線形回帰モデルの比較では、前者の方が観測値に対する当てはまりがよいことが示された。

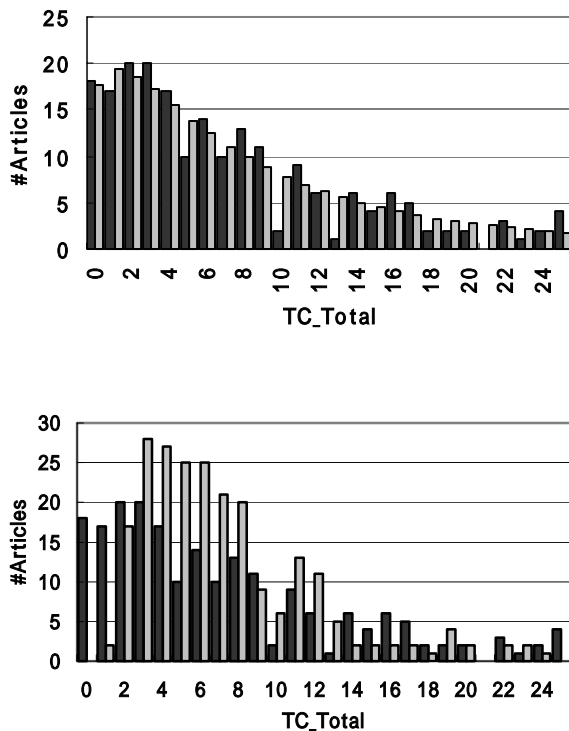


図1 被引用数ごとの分布(上:NBR,下:LR)  
(黒:観測値,グレー:予測値)

謝辞

本研究は、科学研究費補助金(基盤研究(B))により行っているものである。

引用文献

(1) 根岸正光;山崎茂明. 研究評価. 丸善, 2001.

(2) Seglen, P. O. The skewness of science. *J. Am. Soc. Inf. Sci.* 43(9), 628-638, 1992

(3) Snizek, W. E. et al. Textual and non-textual characteristics of scientific papers. Neglected science indicators. *Scientometrics.* 20(1), 25-35, 1991.

(4) Peters, H. P. F.; van Raan, A. F. J. On determinants of citation scores: A case study in chemical engineering. *J. Am. Soc. Inf. Sci.* 45(1), 39-49, 1994.

(5) Callahan, M. et al. Journal prestige, publication bias, and other characteristics associated with citation of published studies in peer-reviewed journals. *JAMA.* 287(21), 2847-2850, 2002.

(6) Basu A.; Lewison, G. Going beyond journal classification for evaluation of research outputs. *Aslib Proc.* 57(3), 232-246, 2005.

(7) van Dalen, H. P.; Henkens, K. Signals in science - On the importance of signaling in gaining attention in science. *Scientometrics.* 64(2), 209-233. 2005.

(8) Bornmann, L.; Daniel, H-D. Multiple publication on a single research study: Does it pay? The influence of number of research articles on total citation counts in biomedicine. *J. Am. Soc. Inf. Sci Technol,* 58(8), 1100-1107, 2007.

(9) Bornmann, L.; Daniel, H-D. Selecting manuscripts for a high-impact journal through peer review: A citation analysis of Communications that were accepted by *Angewandte Chemie International Edition*, or rejected but published elsewhere. *J. Am. Soc. Inf. Sci Technol,* 59(11), 1841-1852, 2008.

(10) Davis, P. M. Author-choice open-access publishing in the biological and medical literature: A citation analysis. *J. Am. Soc. Inf. Sci. Technol.* 60(1), 3-8, 2009.

(11) 小野寺夏生ほか. 論文の被引用数に影響する要因に関する統計学的研究. 第56回日本図書館情報学会研究大会発表要綱, 41-44, 奈良, 帝塚山大学, 2008年11月.

表3 よく選択された説明変数, およびそれらが選択された分野数

NBR				LR			
4誌		1誌		4誌		1誌	
Price	6	Price	5	Price	6	Price	4
Refs	5	Length	4	Refs	5	Impact Max Full	2
Age Prim	2	Authors	3	Authors	3	Length	2
Authors	2	Refs	2	Figures	3	Age Max	2
Figures	2	Tables	2	Length	3	Authors	2
Length	2	Impact Max Full	2	Impact Prim Full	2	Eqs	2
Tables	2	Eqs	2	Insts	2		
				Tables	2		
				Countries	2		