

Data Analysis Regression with CSR Files in Data Science

Kotadi Chinnaiah¹, Akash Saxena²

¹(Research Scholar, Department Cse, Sunrise University, Alwar, Rajasthan, India)

²(Research Supervisor, Department Cse, Sunrise University, Alwar, Rajasthan, India)

Abstract:

Statistics is an important part in big data because many statistical methods are used for big data analysis. The aim of statistics is to estimate population using the sample extracted from the population, so statistics is to analyze not the population but the sample. But in big data environment, we can get the big data set closed to the population by the advanced computing systems such as cloud computing and high-speed internet. According to the circumstances, we can analyze entire part of big data like the population of statistics. But we may be impossible to analyze the entire data because of its huge data volume. So, in this paper, we propose a new analytical methodology for big data analysis in regression problem for reducing the computing burden. We call this a divided regression analysis. To verify the performance of our divided regression model, we carry out experiment and simulation.

Keywords — Divided regression analysis, statistics, population, sample, big data analysis.

1. Introduction

We believe that some components of data science and business analytics have been around for a long time, but there are significant new questions and opportunities created by the availability of big data and major advancements in machine intelligence. While the notion that analytical techniques can be used to make sense of and derive insights from data is as old as the field of statistics, and dates back to the 18th century, one obvious difference today is the rapid pace at which economic and social transactions are moving online, allowing for the digital capture of big data. The ability to understand the structure and content of human discourse has considerably expanded the dimensionality of data sets available. As a result, the set of opportunities for inquiry has exploded exponentially with readily available large and complex data sets related to any type of phenomenon researchers want to study, ranging from deconstructing the human genome, to understanding the pathology of Alzheimer's disease across millions of patients, to observing consumer response to different marketing offers in large scale field experiments. And, easy (and relatively inexpensive) access to computational capacity and user-friendly analytical software have democratized the field of data science allowing many more scholars (and practitioners) to participate in the opportunities enabled by big data.

In some ways it could be argued that the nature of inquiry has also changed, turbocharged by machines becoming a lot smarter through better algorithms, and by information technologies that enable people and things to be inherently instrumented for observation and interaction that feeds the algorithms. Increasingly, data are collected not with the aim of solely testing a human-generated hypotheses or essential record-keeping, but to the extent that data torrents are captured inexpensively, often for the possibility of testing hypotheses that have not yet been envisioned at the time of collection. When such data are gathered on a scale that observes every part of the joint distributions of the observed variables (behaviors, demographics, etc.), the computer becomes an active question asking machine as opposed to a pure analytic servant. By initiating interesting questions and refining them without active human intervention, it becomes capable of creating new knowledge and making discoveries on its own (Dhar 2013). It can, for example, discover automatically from a large swath of

healthcare system data that younger people in a specific region of the world are becoming increasingly diabetic and then conjecture and test whether the trend is due to specific habits, diet, specific types of drugs, and a

2 Arguably, for the first time in history, a machine passed the famous Turing test by defeating human champions at Jeopardy where topics are not known in advance and questions are posed in natural language of considerable complexity and nuance.

range of factors we may not have hypothesized as humans. This is powerful. As scientists, we have not seriously entertained the possibility of theory originating in the computer, and as science-fiction-like as that may sound, we are in principle already there.

New challenging problems and inquiry also lead to research on better algorithms and systems. Since the torrent of data being generated is increasingly unstructured and coming from networks of people or devices, we are seeing the emergence of more powerful algorithms and better knowledge representation schemes for making sense of all of this heterogeneous and fragmented information. Text and image processing capability are one frontier of research, with systems such as IBM's Watson being on the cutting edge in natural language processing, albeit with a long way to go in terms of their capability for ingesting and interpreting big data across the Internet.

Networks, such as those created by connections between individuals and/or products, further create significant and unique challenges at a fundamental level such as how we sample them or infer treatment effects. For example, in A/B testing, a "standard approach" for estimating the average treatment effect of a new feature or condition by exposing a sample of the overall population to it, the treatment of individuals can spill over to neighboring individuals along the structure of the underlying network. To address this type of "social interference," newer algorithms are required that support valid sampling and estimation of treatment effects (Ugander et al. 2013). This is but one example of how "relational" and "networked" data necessitates new development in algorithms. Developments may emerge not only from computer science but also from IS or other disciplines where researchers are "closer" to the problem being studied than pure methods researchers tend to be.

Finally, the Internet has fueled our ability to conduct large scale experiments on social phenomena. As of this writing, Facebook researchers conducted a massive study to determine whether the mood of users could be manipulated and found that it could (Kramer et al. 2014). By conducting controlled experiments in large numbers of people such studies can extract the causal structure among variables. While the study raised important questions about privacy and the ethical implications of conducting experiments without informed consent, the broader point is that researchers now have a medium for theory development through massive experimentation in the social, health, urban, and other sciences.

2. A Divided Regression Analysis

In big data era, we can get huge data closed to population, and gain an opportunity to analyze the population. But traditional statistics needs much time for data analysis, and it is focused on the sample data analysis for inference of population. To settle this burden of big data analysis, we propose an approach to divide big data into some sub data sets as follow.

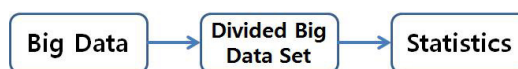


Figure 1. Divided Data Sets for BIG Data science

We divide the big data closed to population into some sub data sets with small size closed to sample. These divided data sets are proper to statistical analysis. In this paper, we select regression analysis as a statistical method for big data analysis because regression is a popular model for data analysis including big data analysis. The regression method has the computing burden because it is also one of statistical methods. To overcome the computing problem in big data regression, we propose a divided regression analysis, which splits whole data into n sub data sets. The whole data are regarded as the population in statistics, and the sub data set stands for a sample. In addition, we apply this data partitioning to estimate the parameters in regression model. In traditional regression analysis, the estimating process of the regression parameters is performed by the following figure.

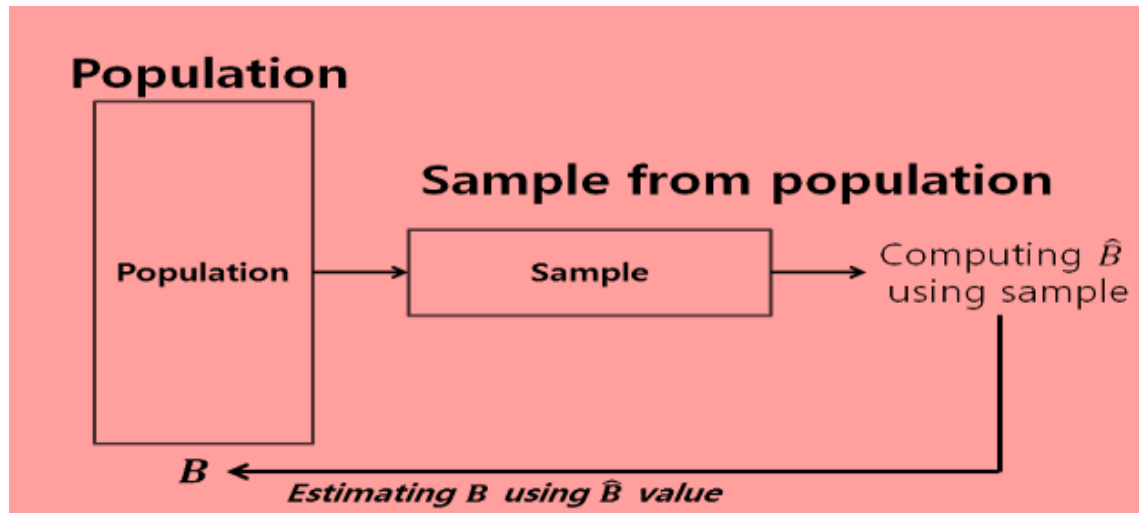


Figure 2. Traditional Regression Analysis

The multiple linear regression model is represented by the following [18].

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \varepsilon \tag{1}$$

Where Y is dependent variable, and x_1, x_2, \dots, x_k are independent variables. Also $\beta_0, \beta_1, \dots, \beta_k$ are regression parameters, and ε is an error of the model. To estimate the regression parameter vector $B = (\beta_0, \beta_1, \beta_2, \dots, \beta_k)$ for the population, we extract a sample data set from the population, and compute an estimate parameter vector $\hat{B} = (\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k)$ using the sample. Hence, we can estimate the regression function as follow.

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_k X_k \tag{2}$$

This is very standard approach in statistical analysis. But, in the big data analysis, we have new problem, which is different to the traditional statistics. This is to use and analyze whole data according to circumstances. In this paper, we consider big data to the population of statistics, and separate the population into sub-populations as follow.

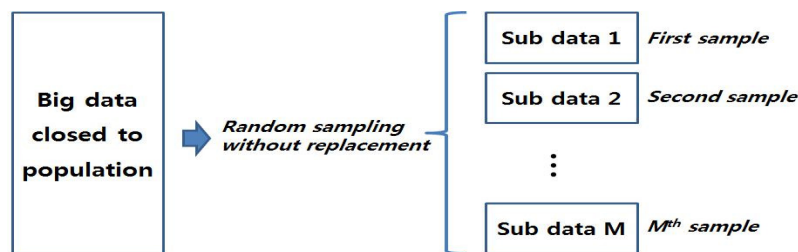


Figure 3. Dividing Big Data using Sampling Method

We use statistical sampling methods to dividing big data into sub samples. There are many sampling techniques from statistics such as simple random sampling, stratified sampling, systematic sampling, and cluster sampling [19]. These diverse sampling techniques were important to big data analysis [20, 21]. We should select sampling techniques carefully according to the aim of study and the characteristic of given data set. The main goal of all sampling methods is to get a representative sample from the population, and all sampling techniques should be based on the random sampling. In the random sampling, all elements of population have equal chances to be selected to sample. In our research, we use simple random sampling without replacement as a sampling method for dividing big data. Next figure shows the proposed method for entire regression analysis.

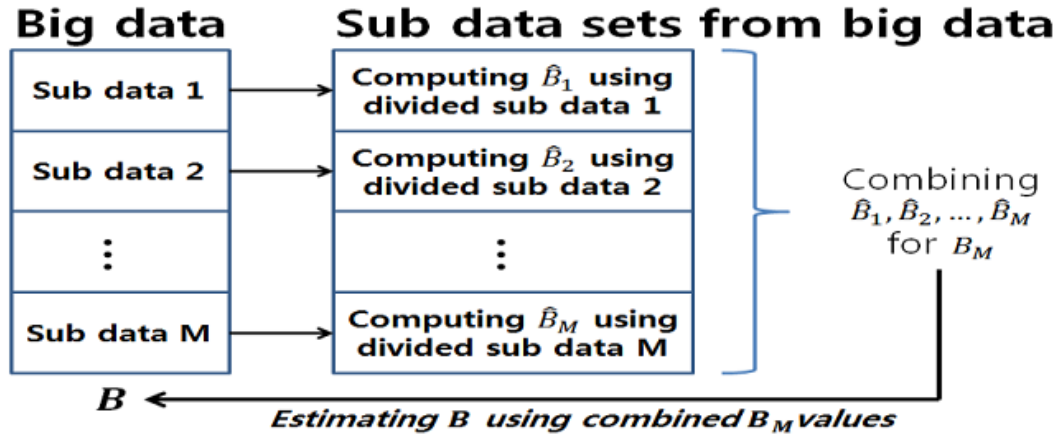


Figure 4. Divided Regression Analysis

Using all data with M sub sets, we compute the parameters which are $\hat{B}_1, \hat{B}_2, \dots, \hat{B}_M$ respectively. We combine the M parameters to decide for estimating B as follow.

Where is a combine function for combining the results from sub-data 1 to sub-data M. We can consider diverse functions for combining the sample results. In our research, we use mean value for combining the estimated parameters as follow.

$$\hat{B}_C = f_C(\hat{B}_1, \hat{B}_2, \dots, \hat{B}_M) \tag{3}$$

$$\hat{B}_C = \frac{1}{M} \sum_{i=1}^M \hat{B}_i \tag{4}$$

So, using we can get same result of population analysis. To validate the statistical significance between and , where is the regression parameter of population, we can compute the significance interval of regression parameters. If the combined parameter of our model is included in the confidence interval, the performance of the estimated parameter by our work will be validated. Also, in this paper, we use mean squared error (MSE) for verifying regression result. This is defined as follow [18].

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \tag{5}$$

The MSE measures the size and importance of forecasting error. Next we show the process of our study.

(Step1) Dividing big data

- (1.1) Separating big data (population) into M sub data sets (samples) by simple random sampling without replacement
- (1.2) Preparing sub data sets for regression analysis

(Step2) Performing multiple linear regression analysis (2.1)

Computing regression parameters for M sub data sets

(2.2) Averaging M regression parameters for estimating the regression parameters of all big data

(Step3) Evaluating model

(3.1) Comparing the regression results of M sub data sets with all big data

(3.3) Computing confidence interval of regression parameter

(3.3) Checking whether the averaged parameter is included in the confidence interval

To verify the performance our research, and to discover the potential of our model for real fields, we consider simulation and experiment in next section.

4. Experimental Results

To evaluate the proposed model, we performed simulation study, and made experiment using data set from UCI machine learning repository [9]. We also checked the confidence intervals of the regression results between divided and full data sets.

4.1 Simulation data

First we carried out an experiment using simulation data set for showing the performance of our research. Assume represents the value of dependent variable and represents the value of independent variable in the th trial. We used the following simple linear regression model.

The error terms $\epsilon_i, i = 1, 2, \dots$, are assumed to be independent random variables having a normal distribution with mean $E(\epsilon_i) = 0$ and constant variance $Var(\epsilon_i) = \sigma^2$. In this experiment, we simulated simple linear regression model having the regression parameters $\beta_0 = 5$ and $\beta_1 = 3$. In this simulation, we set the variance is one, $\sigma^2 = 1$.

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \tag{6}$$

that is, the error terms $\epsilon_i, i = 1, 2, \dots$, were all drawn from a normal distribution with mean 0 and variance 1 denoted by $N(0,1)$. In this simulation, we determined the population size N was 1,000,000. Also, we divided the population into 10 sub-populations, so the size of each sub-population was 100,000. First of all, we got the following regression model using entire simulation data.

$$\hat{Y} = 4.998612 + 2.998295X \tag{6}$$

Next we estimated the regression parameters for 10 sub-populations, respectively. Table 1 shows the simulation result of our divided regression model.

Table 1. Simulation Data Set Result for Divided Regression Model

Data set	β_0	β_1	MSE
Sub.1	5.003847	2.996239	1.004636
Sub.2	4.994696	2.998008	0.996257
Sub.3	4.997074	2.996837	1.002991
Sub.4	4.998384	3.001986	1.000697
Sub.5	5.000347	2.999503	0.995672
Sub.6	5.002670	3.001851	0.998701
Sub.7	4.998995	3.002176	0.997655
Sub.8	4.997011	2.997172	1.002975
Sub.9	4.995441	2.996445	0.999974
Sub.10	4.997664	2.992761	0.996550
Mean	4.998613	2.998298	0.999611
Overall	4.998612	2.998295	0.999609

The size of each sub-population was 100,000, and we estimated β_0 , β_1 , and MES(mean squared error). In the table, "mean" of data set column is mean value of 10 sub-populations, and overall is the result by using the population. We knew the value between mean and overall was very similar, so we verified the performance of our study. Next we show 95% confidence intervals for regression parameters. The following figure shows the 95% confidence interval of intercept β_0 .

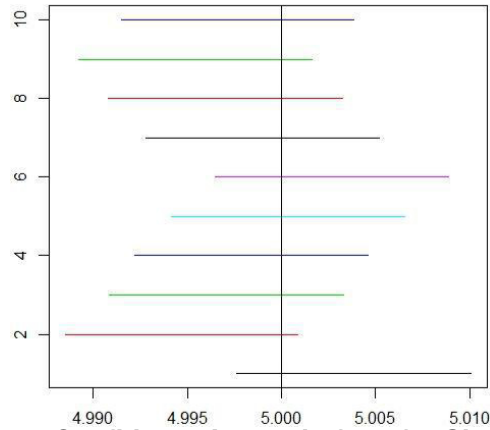


Figure 5. 95% Confidence Interval of β_0 for Simulation Set

We found all sub-populations included the parameter β_0 in their confidence intervals. Next figure shows the 95% confidence interval of the regression parameter β_1 .

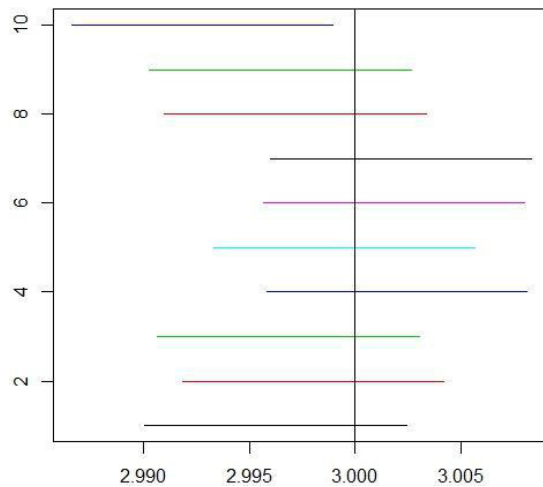


Figure 6. 95% Confidence Interval of β_1 for Simulation Set

Similar to the case of β_0 , all confidence intervals of 10 sub-populations also includes the parameter β_1 of the population. We knew that all confidence intervals of β_0 and β_1 contained their parameters, so we confirmed the validity of our research. To verify the performance of our study clearly, we performed additional simulations by 100 iterations of the previous simulation as follow.

Table 2. Simulation Result of 100 Iterations

iteratio n	β_0	β_1	MSE	iteratio n	β_0	β_1	MSE
1	4.99861	2.99829	0.999610	51	5.00012	3.00160	0.998443
2	3	8	7	52	9	9	9
3	4.99980	2.99792	1.001704	53	4.99924	3.00033	1.001651
4	2	4	7	54	3	4	3
5	4.99992	3.00092	0.999121	55	4.99980	2.99924	0.999774
6	5.00029	3.00089	7	56	8	7	8
7	2	3.00022	1.002215	57	4.99943	2.99889	1.000790
8	4.99905	9	5	58	3	9	5
9	3	2.99895	0.998992	59	4.99934	3.00034	0.999230
10	5.00146	8	9	60	5.00024	6	5
11	5.00202	3.00047	1.000405	61	6	2.99986	0.998855
12	2	3.00091	1.000022	62	4.99952	9	4
13	4.99975	5	2	63	9	2.99637	1.001982
14	4.99928	2.99892	0.997984	64	4.99861	9	9
15	3	6	2	65	5	3.00113	1.000260
16	5.00037	3.00072	1.002252	66	5.00073	4	1
17	5	5	8	67	4	2.99951	1.000943
18	4.99941	3.00064	1.001533	68	4.99936	6	1
19	1	2.99920	0.997547	69	7	3.00120	1.000133
20	5.00059	3	5	70	4.99935	6	6
21	9	3.00117	1.002407	71	6	3.00174	0.999657
22	5.00029	2	8	72	5.00085	5	6
23	2	2.99839	1.000664	73	7	2.99979	0.998683
24	4.99959	8	5	74	4.99906	2	3
25	4.99924	3.00229	1.001462	75	7	2.99951	0.997797
26	9	1	3	76	5.00061	2	4
27	5.00028	2.99942	1.000786	77	1	3.00063	1.002215
28	3	3.00015	9	78	4.99909	7	9
29	4.99849	9	0.99775	79	8	3.00008	0.998546
30	1	3.00081	0.999341	80	5.00084	3	9
31	5.00016	5	1	81	2	3.00034	0.999627
32	3	2.99977	1.001093	82	4.99982	5	1.001381
33	4.99986	1	2	83	4	3.00020	2
34	1	3	0.999262	84	4.99861	6	0.999283
35	4.99852	2.99951	0.998081	85	5	3.00079	8
36	2	8	2	86	5.00007	6	1.000629
37	4.99972	2.99993	0.999585	87	7	2.99875	5
38	3	1	1	88	5.00045	9	1.003816
39	5.00001	2.99754	0.999859	89	5.00076	3.00020	8
40	9	6	2	90	5.00027	5	0.997839
41	5.00052	2.99937	1.000318	91	3	3.00227	4
42	5	9	2	92	4.99982	2	0.999085
43	5.00116	3.00075	1.001588	93	4	2.99961	6
44	4.99898	9	8	94	4.99969	3	1.000650
45	1	3.00086	1.001365	95	5.00028	3.00031	5
46	4.99915	5	1	96	3	1	0.998072
47	8	3.00139	0.998801	97	5.00005	3.00014	4
48	5.00011	7	7	98	8	3	1.001486
49	6	2.99931	0.997960	99	5.00017	2.99982	8

50	5.00051	4	2	100	4.99772	1	0.998634
	9	2.99840	0.998363		7	2.99948	2
	4.99947	8	1		5.00019	9	0.999067
	7	2.99990	0.998700		1	2.99937	6
	5.00074	1	5		4.99913	9	1.000124
	3	2.99984	1.000768		5	3.00062	4
	5.00031	3	9		5.00037	7	1.000995
	8	2.99927	0.996194		4	3.00047	8
	4.99823	5	1		5.00014	8	1.003555
	9	2.99980	1.000732		6	2.99997	0.999692
	4.99893	2	8		5.00095	4	5
	7	3.00045	1.002522		7	2.99882	1.000085
	4.99928	3.00070	4		4.99863	2.99856	8
	5	4	1.003336		5	9	1.000783
	4.99963	3.00177	1		4.99939	3.00007	3
	5.0005	1	0.999906		2	2.99861	1.000637
	4.99961	3.00138	8		5.00016	2	9
	4.99937	8	0.999534		1	3.00075	1.001862
	2	3.00097	1		5.00041	2	3
	5.00153	7	1.000162		8	3.00135	0.999168
	1	2.99934	8		4.99946	7	9
	5.00001	3	1.001951		5	2.99991	0.998725
	4	3.00110	7		5.00059	6	9
	4.99963	4	1.000150		5	3.00062	1.001217
	1	2.99977	3		5.00141	7	2
	4.99976	3.00012	1.002516		1	2.99972	0.999650
	4	8	4		5.00034	9	1
	5.00003	2.99933	1.000065		6	3.00052	0.998050
	6	7	3		4.99790	1	8
	5.00009	3.00008	0.998986		7	2.99858	0.998863
	4	8	7		5.00015	6	8
	4.99901	3.00016	0.999604		6	2.99939	0.999205
	3	2	7		4.99988	8	4
	5.00009	2.99965	1.000660		4.99992	3.00034	1.001113
	7	2	4		4.99939	3	5
	4.99853	3.00096	0.998651		2	3.00011	1.000202
	1	3.00044	0.998931		4.99995	2	6
	5.00095	2.99901	2		8	2.99936	0.998768
	7	8	0.996741		5.00176	2.99900	1
	4.99822	2.99966	3		7	8	1.001441
	3	8	0.997686		4.99929	3.00004	4
	4.99893		5		9	2.99966	1.000369
	6		1.001726		4.99976	4	7
			8		7	2.99954	0.999798
			0.998129			5	2
			2			2.99739	0.997850
							2
							1.001414
							7

We knew this result was similar to the previous simulation result. Next two figures show the 95% confidence intervals of θ and λ from the 100 iterative simulations.

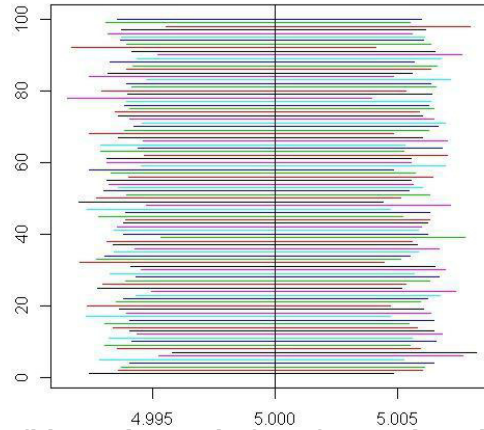


Figure 7. 95% Confidence Interval of for 100 Iterative Simulation Sets

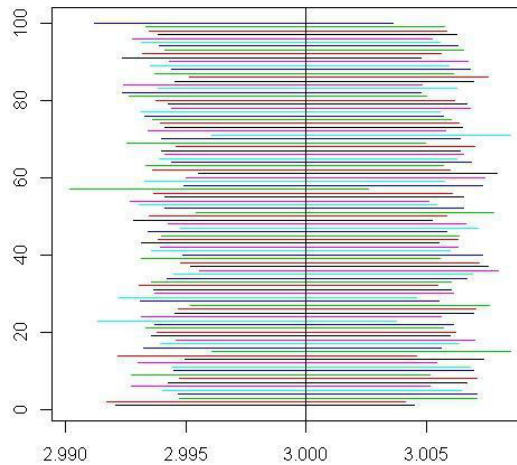


Figure 8. 95% Confidence Interval of for 100 Iterative Simulation Sets

In the two figures for 95% confidence interval, all 100 intervals included the regression parameters of 0 and 1. Therefore we showed the validity of our research. In the next section, we considered another experiment using example data set.

4.2. A case Study using the Bike Data

We made another experiment using example data from UCI machine learning repository [9]. The data set was Bike sharing data set including the time information of bike rental system, and the numbers of attributes and instances were 16 and 17,379 respectively. From the data set, we used three variables. The dependent variable was “cnt” and the independent variables were “temp” and “hum”. The “cnt” represents the number of total rental bikes, and the “temp” and “hum” show temperature and humidity respectively. So we modeled the multiple regression equation as follow.

$$cnt = \beta_0 + \beta_1 temp + \beta_2 hum + \epsilon \tag{7}$$

In this model, we found the influence of temperature and humidity to bike rental. In addition, we divided the Bike sharing data into ten sub data sets for performing our divided regression analysis. The following table shows the regression result.

Table 3. Bike Data Set Result for Divided Regression Model

Data set	β_0	β_1	β_2	MSE
Sub.1	188.8797	369.6033	-295.8221	24535.77
Sub.2	159.6991	370.4749	-250.7093	22948.67
Sub.3	181.4024	373.1526	-281.6638	23589.74
Sub.4	193.5926	353.8599	-288.9963	24233.32
Sub.5	210.175	338.5803	-296.6555	24942.94
Sub.6	180.7998	387.5178	-288.3613	24826.61
Sub.7	194.0239	349.3915	-277.8862	26746.23
Sub.8	162.5819	360.5018	-246.7723	25335.49
Sub.9	184.3552	361.4922	-275.6302	24875.45
Sub.10	185.0865	354.8126	-279.4127	24575.58
Mean	184.0596	361.9387	-278.1910	24659.18
Overall	184.2446	361.8051	-278.3578	24643.72

We knew there are the slight differences between the regression parameters of sub-populations and overall population, but mean value of the regression parameters of sub-populations were similar to the parameter value of entire population. Next we computed the confidence intervals of the regression parameters in ten sub-populations. We show the 95% confidence interval of β_0 is shown in the following figure.

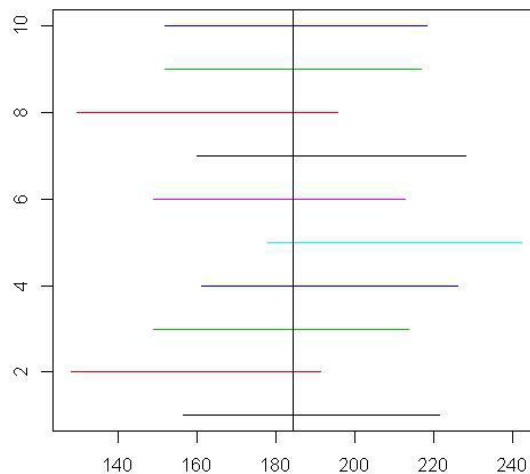


Figure 9. 95% Confidence Interval of for Bike Data Set

All confidence intervals of ten sub-populations contained the regression parameter β_0 of the population. So, we confirmed the validity of our research. Next two figure show the 95% confidence interval of β_1 and β_2 .

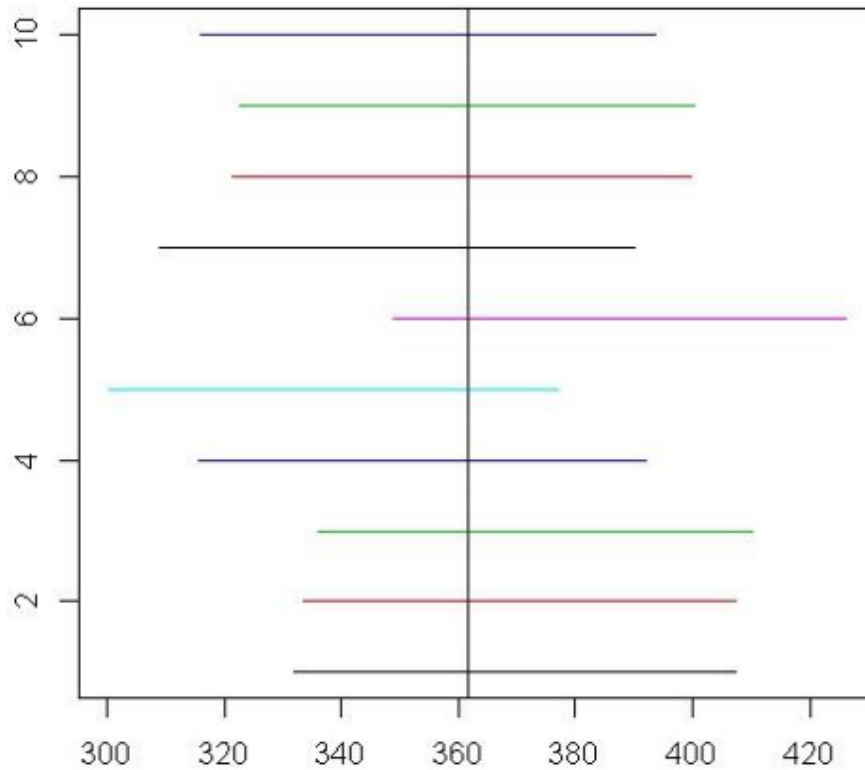


Figure 10. 95% Confidence Interval of for Bike Data Set

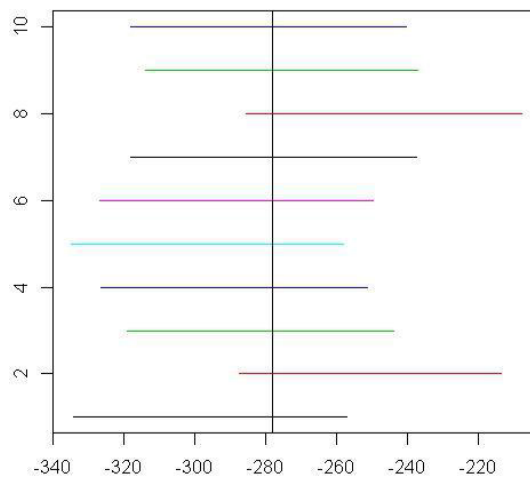


Figure 12. 95% Confidence Interval of for Bike Data Set

All confidence intervals of sub-populations include the regression parameters of the population for 1 and 2. Therefore, we can verify the performance of our regression approach.

5. Conclusions

In this paper, we proposed an approach to overcome the computing burden in big data analysis because most statistical methods were focused on small sample data. Also in big data analysis, we should analyze entire data which are considered as population in statistics, and this data set is so huge. Our research divided the big data closed to population into sub data set like sample for solving the computing cost in big data analysis. In addition, we applied this approach to regression problem in statistics. We

applied the divided method of big data to multiple regression analysis, and used simple random sampling for big data dividing. To verify the performance of our research, we used two data sets from simulation and UCI machine learning repository. In our experimental results, we knew that the regression parameters estimated by the big data were not different to the parameters by sub data sets. This research contributes to avoid the computing problem in many fields for big data analysis. We will apply our approach to more diverse methods in statistics such as factor analysis and clustering. More diverse methods of big data sampling are needed in our future works. We also will study more advanced combining methods for merging the results of sun data sets.

References

- [1] W. Hu and N. Kaabouch, "Big Data Management, Technologies, and Applications", Information Science Reference, IGI Global, (2014).
- [2] S. Jun and D. Uhm, "A Predictive Model for Patent Registration Time Using Survival Analysis", Applied Mathematics & Information Sciences, vol. 7, no. 5, (2013), pp. 1819-1823.
- [3] S. Jun, "A Technology Forecasting Method Using Text Mining and Visual Apriori Algorithm", Applied Mathematics & Information Sciences-An International Journal, vol. 8, no. (1L), (2014), pp. 35-40.
- [4] S. Jun, "Technology Forecasting by Big Data Learning", Proceedings of 2013 NIMS Hot Topic Workshops on Prediction using Fuzzy Theory, vol. 44, (2013).
- [5] J. Han, M. Kamber and J. Pei, "Data Mining: Concepts and Techniques", Third Edition, Morgan Kaufmann, (2012).
- [6] J. Manyika, M. Chui, B. Brown, J. Bughin, R. Dobbs, C. Roxburgh and A. H. Byers, "Big data: The next frontier for innovation, competition, and productivity", McKinsey Global Institute, (2011).
- [7] J. J. Berman, "Principle of Big Data", Morgan Kaufmann, (2013).
- [8] D. Vesset, B. Woo, H. D. Morris, R. L. Villars, G. Little, J. S. Bozman, L. Borovick, C. W. Olofson, S. Feldman, S. Conway, M. Eastwood and N. Yezhkova, "Worldwide Big Data Technology and Services 2012-2015 Forecast", IDC #233485, vol. 1, (2013).
- [9] UCI ML Repository, the UC Irvine Machine Learning Repository, <http://archive.ics.uci.edu/ml>, (2014).
- [10] S. M. Ross, "Introduction to Probability and Statistics for Engineers and Scientists", Elsevier, (2012).
- [11] B. Chun and S. Lee, "A Study on Big Data Processing Mechanism & Applicability", International Journal of Software Engineering and Its Applications, vol. 8, no. 8, (2014), pp. 73-82.
- [12] S. Ha, S. Lee and K. Lee, "Standardization Requirements Analysis on Big Data in Public Sector based on Potential Business Models", International Journal of Software Engineering and Its Applications, vol. 8, no. 11, (2014), pp. 165-172.
- [13] S. Jeon, B. Hong, J. Kwon, Y. Kwak and S. Song, "Redundant Data Removal Technique for Efficient Big Data Search Processing", International Journal of Software Engineering and Its Applications, vol. 7, no. 4, (2014), pp. 427-436.
- [14] J. Stanton, "An Introduction to Data Science", Ver. 3, Syracuse University, (2013).
- [15] S. M. Ross, "Introductory Statistics", McGraw-Hill, (1996).
- [16] P. Vincent, L. Badri and M. Badri, "Regression Testing of Object-Oriented Software: Towards a Hybrid Technique", International Journal of Software Engineering and Its Applications, vol. 7, no. 4, (2013), pp. 227-240.
- [17] V. Gupta, D. S. Chauhan and K. Dutta, "Regression Testing based Requirement Prioritization of Desktop Software Applications Approach", International Journal of Software Engineering and Its Applications, vol. 7, no. 6, (2013), pp. 9-18.
- [18] B. L. Bowerman, R. T. O'Connell and A. B. Koehler, "Forecasting, Time Series, and Regression", An Applied Approach, Brooks/Cole, (2005).
- [19] R. Scheaffer, W. Mendenhall III, R. L. Ott and K. G. Gerow, "Elementary Survey Sampling", 7th Edition, Duxbury, (2011).
- [20] M. Riondato, "Sampling-based Randomized Algorithms for Big Data Analytics", PhD dissertation in the Department of Computer Science at Brown University, (2014).
- [21] J. Lu and D. LiBias, "Correction in a Small Sample from Big Data", IEEE Transactions on Knowledge and Data Engineering, vol. 25, no. 11, (2013), pp. 2658-2663.