

Article

Bayesian Count Data Modeling for Finding Technological Sustainability

Sunghae Jun ^{1,*}

¹ Department of Big Data and Statistics, Cheongju University, Chungbuk, 28503, Korea

* Correspondence: shjun@cju.ac.kr; Tel.: +82-10-7745-5677

Abstract: Technology development changes society and society demands new and innovative technology development. We analyze technology to understand society and technology itself. Many researches have been introduced in various fields. Most of them were about patent analysis. This is because detailed and accurate results of research and development are patented. In this paper, we study on new patent analysis method based on count data model and Bayesian regression analysis. Using count data model, we analyze the technological keywords extracted from the collected patent documents. We use the posterior distribution of Bayesian statistics to reflect the experience and knowledge of the relevant technological experts in the analysis model. Moreover, we apply the proposed model to finding sustainable technologies. Finding and developing sustainable technologies is an important activity for companies and research institutes to maintain their technological competitiveness. To illustrate how our modeling could be applied to real domain, we carry out a case study using the patent documents related to artificial intelligence.

Keywords: count data; Bayesian regression; technological sustainability; Poisson probability distribution; patent analysis

1. Introduction

Technology with sustainability is to keep technological Competitiveness of company [1-2]. Most companies have tried to find their sustainable areas for technological innovation and new product development. So, sustainable technology is important issue in management of technology (MOT) [3]. Many academics, research institutes, and companies have studied on sustainable technologies. Recently, Kim et al. (2018) published a statistical method for sustainable technology analysis [4]. They considered Bayesian inference and social network analysis for the proposed method, applied their research to the technology domain related to artificial intelligence (AI). Also, they used the IPC (international patent classification) codes extracted from patent documents as input data for sustainable technology analysis. The IPC is a hierarchical system of technologies for the classification of patents [5]. For example, the IPC code G06F represents the electric digital data processing technology [6]. In general, IPC codes cover a wide range of technologies. So, it is difficult for us to grasp the detailed technological structure of a specific technology field. In order to overcome this problem, we propose a technology analysis method using patent keywords. The keywords are extracted from patent documents related to specific technology by text mining techniques [7]. Therefore, the technology keyword can represent more detailed description of a specific technology field than the IPC code for sustainable technology analysis. In addition, we propose a statistical modeling using Bayesian count data analysis for understanding sustainability of given technology domain. The count of event is the number of times an event occurs [8]. In this paper, each patent keyword is an event, and we analyze the count data of patent keywords. We consider the Poisson probability distribution for the proposed statistical patent analysis model, because the count data of patent keywords are nonnegative integer values [9]. We also combine the Poisson count model with

Bayesian regression analysis to build Bayesian count data modeling for finding technological sustainability. We use the prior probability distribution to reflect the experience and knowledge of the relevant technology experts in the model. The collected patent data is represented by the likelihood function. We get the posterior distribution by multiplying prior distribution and likelihood function. Finally, these Bayesian probability distributions are applied to the count data regression for Bayesian count data model. Therefore, we carry out the Bayesian count data modeling to find technological sustainability. To show the validity of our modeling, we perform a case study using the patent documents related to AI. The remainder of this paper is organized as follows. In section 2, we show the research backgrounds related to our study. We explain the proposed modeling for finding technological sustainability in section 3. Next section illustrates the result of our case study. In the conclusions section, we conclude our research and describe our future works related to this paper.

2. Sustainable Technology and Patent Analysis

In this paper, sustainable technology means a technology that can sustain a company's technological competitiveness. So, it is important for a company to know what its sustainable technologies are. In the MOT field, many companies and institutes have tried to find their sustainable technologies. Many research results on sustainable technology analysis have been published in academia [1-2,4,10-11]. The sustainable technology analysis has been made using diverse analytical methods in various technological fields. However, research on this field is still lacking. Companies and research institutions are demanding a more sophisticated and feasible methodologies for sustainable technology analysis.

Most methods for sustainable technology analysis rely on patent analysis. This is because patents contain accurate and vast results on the research and development of technology. This is due to the exclusive right to use the technology granted to the inventor. Thus, a great deal of research has been done on patent analysis [12-16]. In patent analysis, we should transform the collected patent documents into structured data consisting of keywords, IPC codes, citations, etc. for statistical analysis. In the preprocessing process of patent data, we use R data language and its 'tm' package [7,17]. Using the structured patent data, we perform the proposed modeling for finding technological sustainability.

3. Finding Technological Sustainability using Bayesian Count Data Modeling

Bayesian modeling has been used in diverse data analysis areas such as regression and classification increasingly. This is one of two approaches to statistics. We start the Bayesian count data modeling from the following expression [18].

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)} \quad (1)$$

Where θ is model parameter, and y is response variable to be predicted. $P(\theta)$ and $P(\theta|y)$ are the prior and posterior probabilities of parameter θ respectively. $P(y|\theta)$ represents the likelihood function of y given θ respectively. Also, $P(y)$ is calculated by the following integration [19].

$$P(y) = \int P(y|\theta)P(\theta)d\theta \quad (2)$$

Using Bayesian modeling, we determine the model parameter of posterior distribution and compute the credible interval of true parameter. We are interested in the mean of the parameters in this credible interval. The $100(1-\alpha)\%$ credible interval of θ is defined as follow [19].

$$P(\theta \in C) = 1 - \alpha \quad (3)$$

Where C is an interval depended on y . In addition, the credible interval of Bayesian modeling is based on the posterior distribution. We have to select a prior distribution for beginning a Bayesian

modeling. Non-informative or informative priors can be used for the prior distribution in Bayesian modeling. In general, the result of non-informative prior is close to the maximum likelihood estimate (MLE) in frequentist statistics, because the popular non-informative prior is uniform distribution [9,19]. On the other hand, we should use the informative prior to get amend result for the parameter estimation. But, we should carry out Bayesian computing such as Markov Chain Monte Carlo (MCMC) for using the informative prior [19]. To alleviate the computational burden, we can use conjugate prior. In Bayesian count data modeling, the gamma and beta distributions are conjugate to the Poisson and binomial distributions respectively. In our research, we consider a Bayesian count data modeling with Poisson distribution for finding technological sustainability. The Poisson probability distribution is the most popular model for count data. If the random variable Y is distributed to Poisson with parameter λ , its distribution is defined as follow [9]; $Y \sim \text{Poisson}(\lambda)$.

$$P(Y = y|\lambda) = \frac{e^{-\lambda}\lambda^y}{y!}, y = 0,1,2, \dots \tag{4}$$

Where the expectation $E(Y)$ and variance of Y are equal to parameter λ . In this paper, the frequency of each patent keyword extracted from patent document data is a Poisson random variable with parameter λ_i as follow.

$$\text{frequency of keyword}_i \sim \text{Poisson}(\lambda_i), i = 1,2, \dots, m \tag{5}$$

Where m is the number of all keywords. We extract the patent keywords from the collected patent documents using text mining techniques as follow.

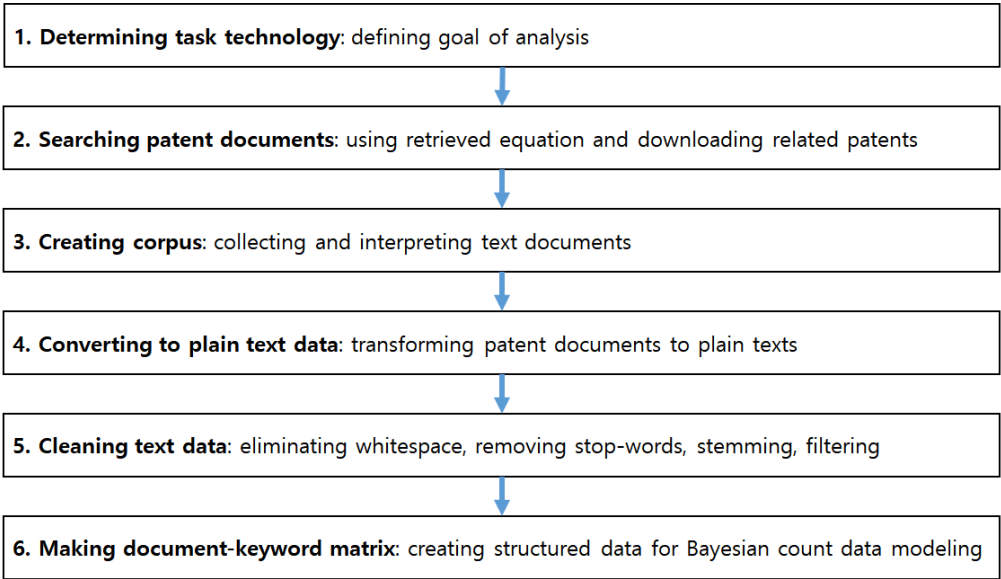


Figure 1. Text mining process for creating structure data for Bayesian count data modeling.

In our text mining process, first of all we should determine the target technology for statistical analysis. Next, we use a retrieved equation to collect the patent documents related to target technology from the patent databases in the world. It is impossible to analyze the searched patent document data directly, because the data is not suitable to input data for statistical analysis including Bayesian count data modeling. So, we try to make a structured data for statistical analysis. The first step in creating structured data is to create corpus collecting and interpreting text documents. Based on the created corpus, we transform the patent documents into plain texts, and clean the text data by eliminating whitespace, removing step-words ("and", "for", "in", "is", etc.), stemming, and filtering. Finally, we make the document-keyword matrix as a structured data for Bayesian count data

modeling. The matrix consists of patent (row) and keyword (column), and its elements are the frequency (count) values of occurred keywords in each patent document.

In our modeling, we define the frequency (count) of i th occurred keyword as y_i , and represent the data set as follow [9].

$$y_i \sim \text{Poisson}(\lambda_i), E(y_i) = \text{Var}(y_i) = \lambda_i, i = 1, 2, \dots, m \quad (6)$$

Using this data set, we perform the generalized linear model (GLM) with Poisson probability distribution, no predictors, and log link function as follow.

$$\log(\lambda_i) = \beta_0 \text{ or } \lambda_i = \exp(\beta_0), i = 1, 2, \dots, m \quad (7)$$

We determine the response and predictor variables (keywords) by the results of GLM. In our study, we use variables with large coefficient values as response variables and those with small coefficient values as explanatory variables.

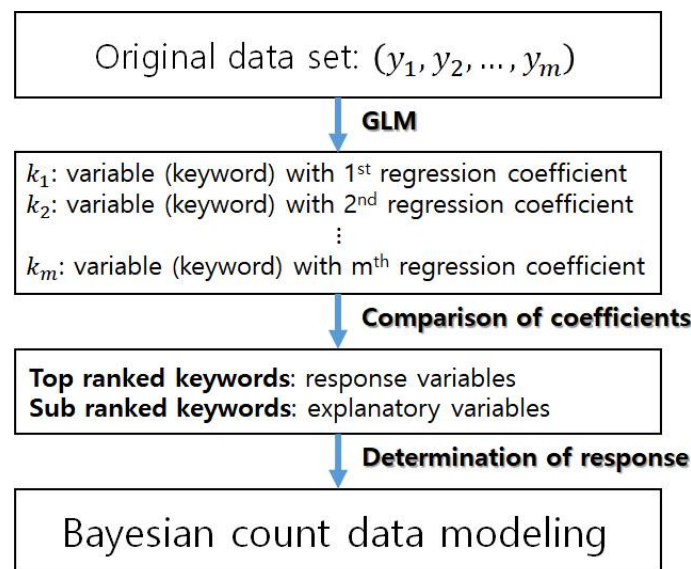


Figure 2. Determination of response variable by GLM results.

In this paper, we denote the response variable as Y , and the explanatory variables as (x_1, x_2, \dots, x_p) . Where p is the number of explanatory variables. For Bayesian count data modeling, we make the Poisson regression model with gamma distribution as prior. The Poisson regression model with λ is defined as follow [20].

$$f(Y|\lambda) = \exp(-n\lambda + n\bar{Y} \log(\lambda) - \sum_{i=1}^n \log(Y_i!)) \quad (8)$$

Where n is the number of collected patent documents. Also, an informative gamma prior for λ as follow.

$$P(\lambda) = \frac{e^{-b\lambda} \lambda^{a-1} b^a}{\Gamma(a)} \quad (9)$$

Where $\Gamma(\cdot)$ is gamma function, and $E(\lambda)$ and $\text{Var}(\lambda)$ are $\frac{a}{b}$ and $\frac{a}{b^2}$ respectively. This expression is used for the likelihood in the Bayesian count data modeling. So, using the likelihood and prior distributions, we show the posterior distribution as follow.

165
$$P(\lambda|Y) = \exp(-n\lambda + n\bar{Y} \log \lambda - \sum_{i=1}^n \log Y_i!) \times \frac{e^{-b\lambda} \lambda^{a-1} b^a}{\Gamma(a)} \tag{10}$$

166 We can ignore the terms not involving λ , so we yield the proportional result of posterior
167 distribution as follow.

169
$$P(\lambda|Y) \propto \exp(-(n+b)\lambda + (n\bar{Y} + a - 1)\log \lambda) \tag{11}$$

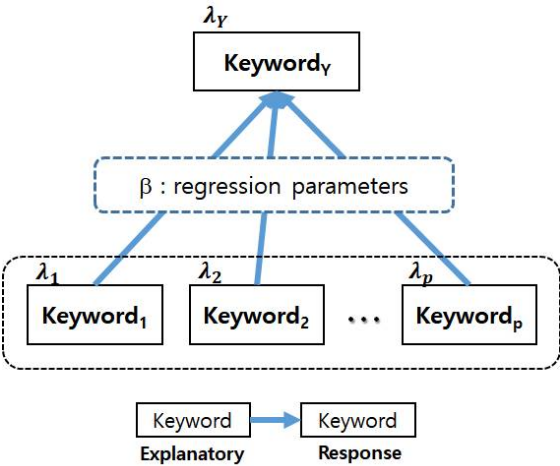
170 This expression represents the kernel of the gamma distribution with parameters $(n+b)$ and
171 $(n\bar{Y} + a)$. In addition, by the characteristic of gamma distribution, the posterior mean and variance
172 of λ are $\frac{n\bar{Y}+a}{n+b}$ and $\frac{n\bar{Y}+a}{(n+b)^2}$ respectively. In Bayesian Poisson regression case, $Y_i|x_i$ is distributed
173 Poisson with mean $\lambda_i = \exp(x_i'\beta)$, where β is the parameter vector of Poisson regression. In our
174 research, $(Y|x)$ is shown as (response keyword | explanatory keywords). Therefore, we get the
175 Bayesian count data modeling as follow.

176
$$Y_i|x_i \sim \text{Poisson}(\lambda_i) \tag{12}$$

177
$$\lambda_i = \exp(x_i'\beta) + \alpha_i, \alpha_i \sim \text{Normal}(0, \sigma^2) \tag{13}$$

178
$$P(\beta): \text{Gamma prior density of } \beta \tag{14}$$

179 Using this modeling, we make a technology structure to understand the target technology from
180 the viewpoint of sustainability in Figure 3.



188
189 **Figure 3.** Technology structure of Bayesian count data modeling.

190 In our modeling, all keywords except response are used as explanatory keywords. In Figure 3,
191 the keyword at the beginning of the arrow indicates the explanatory variable, and the keyword at the
192 end of the arrow indicates the response variable. Also, each keyword is distributed to Poisson with
193 parameter λ_i . From the final result of Bayesian count data modeling, we get the regression
194 coefficients β between response and explanatory variables (keywords). Using the β , we build a
195 technological structure of target domain for sustainable technology management. Therefore, the
196 Bayesian count data modeling is based on the following concept.

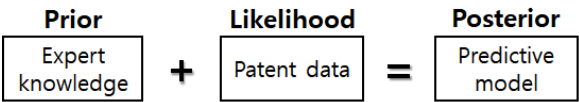


Figure 4. Concept of Bayesian count data modeling.

In this paper, we try to combine the expert’s subjective knowledge and objective result from patent data analysis. That is, the prior represents the domain knowledge of experts, and the likelihood denotes the objective data based on patent documents. The result of multiplying prior and likelihood is posterior, we use this as predictive model for finding technological sustainability.

Using this approach, we expect the improved performance of patent technology analysis for sustainable technology. We illustrate how this research could be applied to practical problem by a case study in next section.

4. Case Study

To show how this research could be applied to practical problem, we performed a case study using the patent documents related to artificial intelligence (AI) technology. We collected the patents applied and registered by 2016 from the WIPSON [21]. The total number of collected and valid patents was 11,973 cases. First of all, we consulted the experts on AI and extracted the keywords related to AI from the collected patent document data [22]. Next, using the text mining techniques, we built structured patent data for performing our case study [7,17]. Our structured patent data is shown in Figure 5.

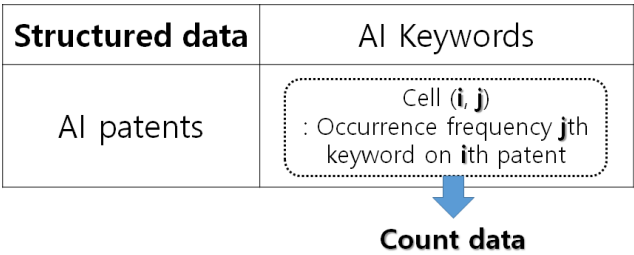


Figure 5. Structured patent data.

The row and column of this data matrix are patents and keywords related to AI, and each cell of this matrix represents occurrence frequency of each keyword on a patent. Based on the structured data, we classified the AI technology as follow.

Table 1. Hierarchical structure of AI technology

Sub-technology	Patent keywords
Learning	Learning, inference, ontology, representation, analysis, data
Behavior	Behavior, awareness, situation, sentiment, mind, spatial, collaborative
Language	Language, natural, understanding, morphological, dialogue, sentence, corpus, voice, speech, conversation, interface
Vision	Vision, figure, object, video, image
Neuro	Neuro, network, computing, feedback, pattern, recognition, cognitive

In Table 1, we denoted the AI technology to five sub-technologies as follow; learning, behavior, language, vision, and neuro. In addition, we have shown the patent keywords belonging to each sub-technology. We used this technology tree to retrieve the AI patents and analyze them. First, we estimated the Poisson parameters for the patent keywords using maximum likelihood estimator (MLE) by the frequency values of the keywords. Table 2 shows the estimates of Poisson parameters for all patent keywords.

233

Table 2. Estimates of Poisson parameters for all keywords

Keyword	λ	Keyword	λ	Keyword	λ
analysis	0.0287	image	0.2745	pattern	0.1506
awareness	0.0008	inference	0.0019	recognition	0.0211
behavior	0.0258	interface	0.0170	representation	0.0062
cognitive	0.0001	language	0.0426	sentence	0.0103
collaborative	0.0001	learning	0.0058	sentiment	0.0016
computing	0.0010	mind	0.0004	situation	0.0060
conversation	0.0035	morphological	0.0004	spatial	0.0966
corpus	0.0123	natural	0.0011	speech	0.5527
data	0.8316	network	0.2668	understanding	0.0005
dialogue	0.0009	neuro	0.0014	video	0.5114
feedback	0.0287	object	1.3510	vision	0.0123
figure	0.0007	ontology	0.0044	voice	0.0153

234

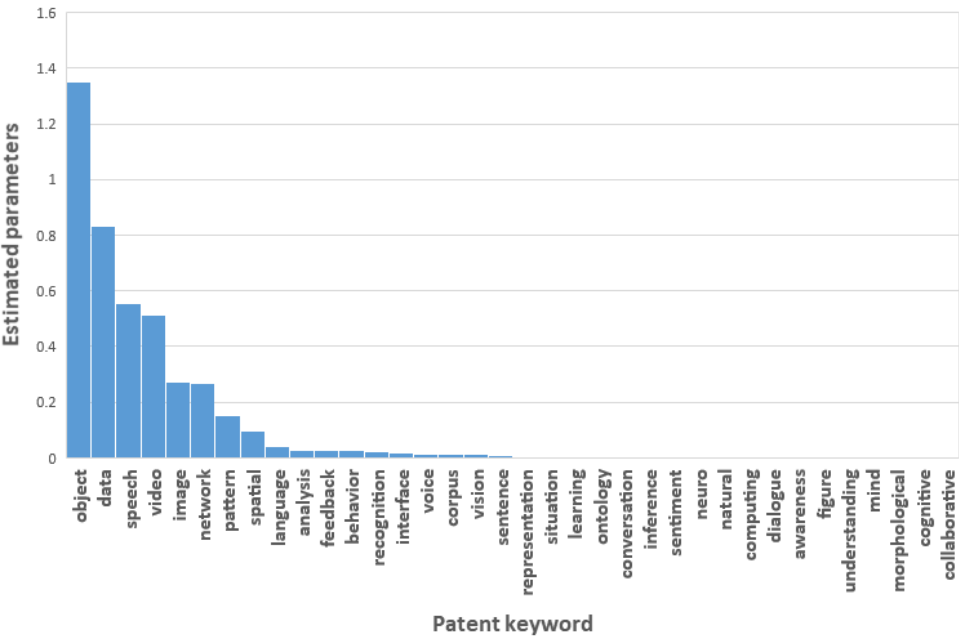
235

236

237

238

We can compare the relative frequency between the patent keywords in Table 2. These estimates are the MLE (maximum likelihood estimate) for Poisson parameters of AI keywords [18]. Figure 6 illustrates the MLEs for the Poisson parameters of all patent keywords.



239

240

Figure 6. MLEs of all patent keywords.

241

242

243

244

245

246

247

In this figure, we found the MLEs of ‘object’, ‘data’, ‘speech’, ‘video’, ‘image’, ‘network’, ‘pattern’, ‘spatial’, ‘language’, ‘analysis’, ‘feedback’, ‘behavior’, ‘recognition’, ‘interface’, ‘voice’, ‘corpus’, ‘vision’, and ‘learning’ are relatively larger than other keywords. Using the result of Figure 6, we determine the patent keywords that affect AI technology. A keyword with a larger MLE value will have more impact on AI technology. We also carried out the Bayesian count data modeling on the structured patent data matrix. Table 3 shows the results of the modeling.

248

249

Table 3. Model parameters and weight

Keyword	Poisson	Gaussian	Weight
analysis	0.1023	0.0176	0.0599
awareness	-0.6630	-0.0013	-0.3321
behavior	0.1062	0.0103	0.0582
cognitive	-0.2173	-0.0016	-0.1094
collaborative	-0.1971	0.0000	-0.0985
computing	0.0867	0.2221	0.1544
conversation	-1.7300	-0.0009	-0.8654
corpus	-4.5122	-0.0009	-2.2566
data	-0.2869	-0.0065	-0.1467
dialogue	-0.8820	-0.0007	-0.4414
feedback	-5.6046	-0.0046	-2.8046
figure	-0.6320	-0.0002	-0.3161
image	-16.8687	-0.0159	-8.4423
inference	-1.1521	-0.0013	-0.5767
interface	0.1670	0.0284	0.0977
language	-0.2675	0.0020	-0.1327
learning	0.1336	0.0500	0.0918
mind	-0.5577	-0.0004	-0.2791
morphological	-0.5222	-0.0012	-0.2617
natural	-0.8831	-0.0016	-0.4424
network	-0.3976	-0.0099	-0.2038
neuro	-1.5059	-0.0007	-0.7533
object	-0.3076	-0.0085	-0.1580
ontology	-2.0103	-0.0028	-1.0065
pattern	0.1729	0.0155	0.0942
recognition	-0.0030	0.0214	0.0092
representation	-1.2470	-0.0080	-0.6275
sentence	-3.2219	-0.0033	-1.6126
sentiment	-1.6950	-0.0006	-0.8478
situation	-2.1258	-0.0019	-1.0638
spatial	0.1356	0.0144	0.0750
speech	-1.2228	-0.0131	-0.6180
understanding	-0.5460	-0.0037	-0.2749
video	-0.5061	-0.0072	-0.2567
vision	-3.2688	-0.0032	-1.6360
voice	-3.6855	-0.0027	-1.8441

250

251

252

253

254

255

256

We performed the Bayesian regression models by Gaussian as well as Poisson distributions. In addition, the weight in Table 3 is the average value of Poisson and Gaussian parameters. In this paper, we selected the keywords with larger weight values for finding technological sustainability in AI technology. Using the result of Table 3, we show the patent keyword ranking that influences AI technology in Figure 7.

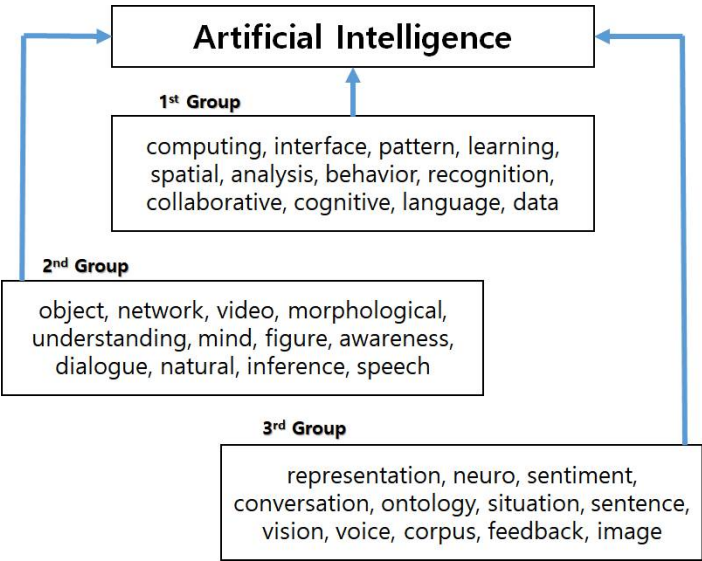


Figure 7. Patent keyword ranking that influences artificial intelligence.

The keywords in 1st Group have a greater impact on AI technology than 2nd Group or 3rd Group. Using the experimental results of this paper, we made the following technological structure for sustainable AI technology.

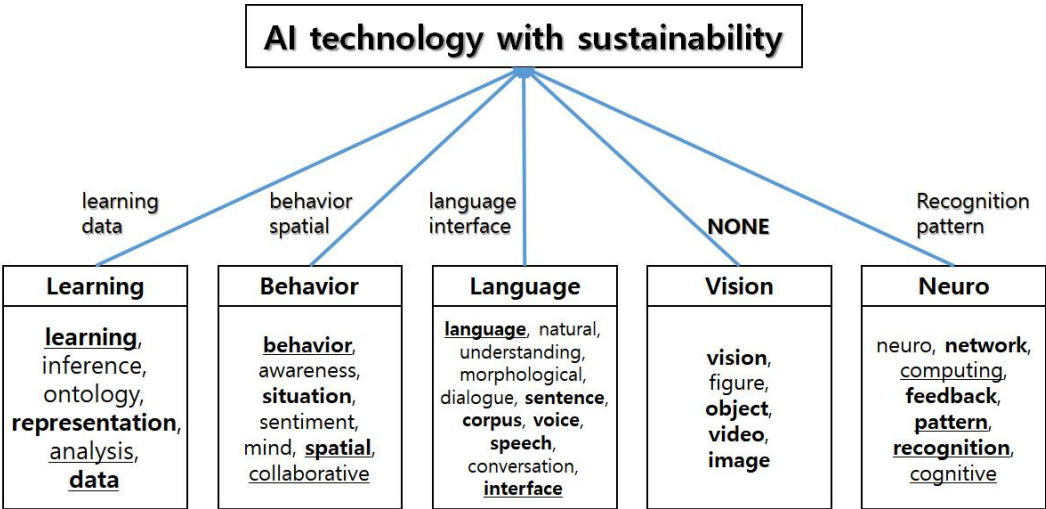


Figure 8. Technological structure for sustainable AI.

In this paper, we divided the AI technology into five sub-technologies (learning, behavior, language, vision, neuro) in Table 1. In Figure 8, each sub-technology contains the keywords that can represent its technology. For example, the keywords of learning, inference, ontology, representation, analysis, and data describe the learning technology for AI. Each keyword is represented by bold or underlined types depending on its importance from the results of Tables 2 and 3. The keywords in bold type are those that have an impact to AI technology from the results of Poisson MLEs. Also, the keywords with underlined lettering affect the AI technology by the Bayesian regression model. So, we knew that the sub-technologies related to learning, behavior, language, and neuro influence to the sustainability of AI technology. But we found that the sub-technology of vision has a relatively small effect on the technological sustainability of AI compared to other sub-technologies. In this paper, we conclude the four technologies related to ‘learning data’, ‘behavior spatial’, ‘language

interface’, and ‘recognition pattern’ are important things to continue the sustainability for AI technology.

5. Conclusions

We proposed Bayesian count data modeling to find the technological sustainability. To know the sustainable technology in given technological field is very important to improve the technological competition of company and nation. Various researches have been conducted to find sustainable technologies. Most of them carried out the statistical models not consider the characteristic of count data from patent documents. But, most structured patent data have a count data structure. In order to solve this discrepancy problem, Bayesian count data modeling is proposed in this study. In addition, we applied the domain knowledge of experts to prior distribution of model parameter in Bayesian regression model. To show the validity of proposed modeling and illustrate how our approach could be applied to practical problem, we carried out a case study using the patent documents related to AI technology. In the case study, we found the sub-technologies for the sustainable technologies of AI. They were learning from data, spatial behavior, interface of language, and pattern recognition technologies. Therefore, we should concentrate our research and development on these sub-technologies to keep the sustainability of AI technology.

Our research can be applied to the research and development-related planning of companies and research institutes. Also, this research will contribute to diverse technological fields as well as AI technology. In this research, we considered only the patent keywords extracted from patent documents for patent technology analysis using statistical modeling. Our future work will use more diverse elements, as well as keywords, such as citations and claims, to find sustainable technologies for specific technological domain.

Author Contributions: Sunghae Jun designed this study and collected the data for the experiment. He also preprocessed the data and selected valid patents and analyzed the data to show the validity of the study and wrote the paper and performed all the research steps.

Conflicts of Interest: The authors declare no conflict of interest.

References

1. Park, S.; Jun, S. Statistical Technology Analysis for Competitive Sustainability of Three Dimensional Printing. *Sustain.* **2017**, *9*, 1142.
2. Choi, J.; Jun, S.; Park, S. A patent analysis for sustainable technology management. *Sustain.* **2016**, *8*, 1–13.
3. Roper, A.T.; Cunningham, S.W.; Porter, A.L.; Mason, T.W.; Rossini, F.A.; Banks, Forecasting and Management of Technology; John Wiley & Sons: Hoboken, NJ, USA, 2011.
4. Kim, J.; Jun, S.; Jang, D.; Park, S. Sustainable Technology Analysis of Artificial Intelligence Using Bayesian and Social Network Models. *Sustain.* **2018**, *10*, 115.
5. WIPO. World Intellectual Property Organization. Available online: www.wipo.org (accessed 2018).
6. WIPO IPC. International Patent Classification (IPC), World Intellectual Property Organization. Available online: <http://www.wipo.int/classifications/ipc/en> (accessed 2018).
7. Feinerer, I.; Hornik, K. Package ‘tm’ Ver. 0.7-5, Text Mining Package, CRAN of R Project. 2018. Available online: <https://cran.r-project.org/web/packages/tm/tm.pdf> (accessed on 1 August 2018).
8. Hilbe, J. M.; Modeling Count Data, Cambridge University Press: Cambridge, UK, 2014.
9. Cameron, A. C.; Trivedi, P. K. Regression Analysis of Count Data; Cambridge: New York, NY, USA, 2013.
10. Kim, S.; Jang, D.; Jun, S.; Park, S. A novel forecasting methodology for sustainable management of defense technology. *Sustain.* **2015**, *7*(12), 16720–16736.
11. Park, S.; Lee, S.; Jun, S. A network analysis model for selecting sustainable technology. *Sustain.* **2015**, *7*(10), 13126–13141.
12. Jun, S.; Park, S. Examining technological innovation of Apple using patent analysis. *Ind. Manage. Data Syst.* **2013**, *113*(6), 890–907.

327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346

13. Kim, J.; Jun, S. Graphical causal inference and copula regression model for Apple keywords by text mining. *Adv. Eng. Inform.* **2015**, *29*(4), 918–929.

14. Jun, S.; Park, S. Examining technological competition between BMW and Hyundai in the Korean car market. *Technol. Anal. Strateg. Manage.* **2016** *28*(2), 156–175.

15. Grimaldi, M.; Cricelli, L.; Rogo, F. Valuating and analyzing the patent portfolio: the patent portfolio value index. *European J. of Innovation Manage.* **2018** *21*(2), 174–205.

16. Kim, J.; Jun, S.; Jang, D.; Park, S. An Integrated Social Network Mining for Product-based Technology Analysis of Apple. *Ind. Manage. Data Syst.* **2017** *117*(10), 2417–2430.

17. R Development Core Team. *R: A Language and Environment for Statistical Computing*; R Foundation for Statistical Computing: Vienna, Austria. Available online: <http://www.R-project.org>, 2018.

18. Ross, S. M. *Introduction to Probability and Statistics for Engineers and Scientists*, 4th Edition; Elsevier: Seoul, Korea, 2012.

19. Gelman, A.; Carlin, J.B.; Stern, H.S.; Dunson, D.B.; Vehtari, A.; Rubin, D.B. *Bayesian Data Analysis*, 3rd Edition; Chapman & Hall/CRC Press: Boca Raton, FL, 2013.

20. Hogg, R. V.; McKean, J. M.; Craig, A. T. *Introduction to Mathematical Statistics*, 8th edition, Pearson: Upper Saddle River, NJ, 2018.

21. WIPSON. WIPS Corporation. Available online: <http://www.wipson.com>, <http://global.wipscorp.com> (accessed 2018).

22. KISTA, Korea Intellectual Property Strategy Agency, Available online: <http://www.kista.or.kr> (accessed 2018).