

텍스트 마이닝을 통한 ESG 지표 개발 및 ESG 지표-주가 상관관계 도출 연구

정재준*, 이강산**, 조은학***

*건국대학교 경영학과

**중앙대학교 기계공학부

***명지대학교 융합소프트웨어학부

jaej0321@naver.com, 97kangsan@gmail.com, 02069@naver.com

A Study For Developing ESG Indicators Through Text Mining And Deriving A Correlation Between ESG Indicators And Stock Prices

Cheong Gae Jun*, Lee Kang San**, Cho Eun Hak***

*Dept. of Business Administration, Konkuk University

**Dept. of Mechanical Engineering, Chung-Ang University

***Dept. of Convergence Software, Myongji University

요 약

기업의 전통적인 역할은 “이익 추구”였다. 현대 사회에 이르러서는 기업이 기존의 역할을 벗어나 새로운 사회적 기구가 되어야 한다는 주장이 떠오르며 CSR(기업의 사회적 책임)이 대두되었다. 최근 기업과 사회는 ESG 경영(환경, 사회, 지배구조를 고려한 지속가능경영)에 많은 관심을 보이고 있고 이는 더 이상 관심으로 그칠 수 없는 필수적인 요소가 되었다. 이에 본 연구는 텍스트 마이닝을 통해 ESG 지표를 개발하고 [ESG 지표 - 주가]의 상관관계를 도출하였다.

1. 개요

환경, 사회, 지배구조에 대한 기업의 태도가 윤리적 정당함에 직접적인 영향을 끼치고 있는 오늘날, 기업에 있어 ESG 경영은 선택이 아닌 필수가 되었다.

2014 년 대한항공의 땅콩 회항사건, 2013 년 대국민 사과 사건과 2019 년 황하나 마약 투약 사건 등의 남양유업, 2021 년 ‘남협’논란에 의한 GS25 등 많은 기업들이 ESG 경영에 실패하였고 이로 인해 대표이사직 박탈, 소비자의 불매 운동, 브랜드 평판 지수 하락과 같이 기업에 큰 악영향을 야기하였다.

실제로 여러 연구 사례에서도 ESG 가 기업 성과와 높은 관련이 있다는 것이 입증되었다. 펀드의 ESG 수준이 펀드의 성과 및 현금흐름에 미치는 영향[4], ESG 정보와 가치관련성에 관한 연구[5], 비재무지표와 기업의 시장성과 관련된 연구: ESG 지표 개발에 사용되는 사건의 시장반응 분석[6] 등 여러 기존 연구들이 ESG 와 경영성과와의 상관관계들을 입증하고 있다.

그러나 기존의 연구들은 자체적으로 개발한, 혹은 자체적으로 개발가능한 독립적인 변수들이 아닌, 이미 존재하고 있는 ‘기업평가 등급’ 혹은 ‘기존 지표’

들을 사용하였다는 점에서 한계점이 있다.

이에 본 연구는 ESG 보고서와 뉴스 데이터 마이닝 등을 통해 자체적으로 ESG 지표를 개발하고 이를 기반으로 [ESG 지표-기업성과]의 상관관계를 도출하고자 한다.

2. 본론

가. 선행연구 정의

<그림 1> 은 총 스터디 케이스 1816 개를 이용하여 메타 분석을 진행한 연구로 대부분의 스터디 케이스들이 ESG 경영과 기업 성과가 연관이 있으며 특히 양의 상관관계가 있다는 것을 보여주고 있다.

하단 <그림 2> 는 기업의 초과이익률, 영업성과, 토빈큐를 종속변수로 ESG 등급을 독립변수로 하여 Z 검정과 T 검정을 진행한 케이스로 AA 등급의 기업성과가 E 등급의 기업성과보다 유의미하게 높은 것을 보여주고 있다.

<그림 3> 은 총자산 순이익률, 유보액/납입자산 등과 같이 기업성과와 관련된 다양한 변수들을 종속변수로,

Study	Focus	Number of studies (N)	Positive	Neutral	Negative	Mixed
Arlow and Gannon (1982)	S	7	42.9%	42.9%	14.3%	
Cochran and Wood (1984)	S, E	13	69.2%	23.1%	7.7%	
Aupperle, Carroll, and Hatfield (1985)	S, E	9	55.6%	22.2%	11.1%	11.1%
Ullmann (1985)	S, E	24	54.2%	20.8%	12.5%	12.5%
Cipson, Farley, and Hoernig (1990)	S, E	14	75.0%		4.6%	4.6%
Wood and Jones (1995)	S, E	51	49.0%	21.6%	13.7%	15.7%
Pava and Krausz (1996)	S, E	21	57.1%	38.1%	4.8%	
Griffin and Mahon (1997)	S, E	50	44.0%	12.0%	22.0%	22.0%
Roman, Hayhoe, and Agle (1999)	S, E	45	60.0%	24.4%	4.4%	11.1%
Richardson, Welker, and Hutchinson (1999)	E, S	22	50.0%	45.5%	4.5%	
Margolis and Walsh (2003)	S, E	126	42.9%	22.2%	5.6%	29.4%
Salzman, Ionescu-Somers, and Steger (2005)	S, E	12	50.0%	25.0%	25.0%	
McWilliams, Siegel, and Wright (2006)	S, E	12	33.3%	25.0%	16.7%	25.0%
Gilao and Starks (2007)	G	39	35.9%	43.6%	5.1%	15.4%
Ambee and Lanoie (2007)	E	41	68.3%	22.0%	4.9%	4.9%
van Beurden and Gossling (2008)	E, S	34	67.6%	26.5%	5.9%	
Petosa (2009)	S, E	130	63.0%	22.0%	15.0%	
Bianco, Rey-Maquieira, and Lozano (2009)	E	32	71.9%	21.9%	6.3%	
Molina-Azorin et al. (2009)	E	32	62.5%	12.5%	12.5%	12.5%
Horvathová (2010)	E	44	54.7%	29.7%	15.0%	0%
Westlund and Adam (2010)	S	21	85.7%			14.3%
Love (2010)	G	45	77.8%	0%	22.2%	27.8%
Dorwall, Koozjik, and Horst (2011)	Funds	18	16.7%	33.3%	11.8%	43.7%
Günther, Hoppe, and Endrikat (2011)	E	274	44.5%			
Sjöström (2011)	E, S	21	23.8%	33.3%	14.3%	28.6%
Borventura, Santos da Silva, and Bandeira-de-Mello (2012)	S, E	58	55.2%	27.6%	10.3%	6.9%
Rathner (2013)	Funds	25	13.2%	72.0%	14.9%	0%
Schütze and Trommer (2012)	E	36	50.0%	19.4%	5.6%	25.0%
Viviers and Eccles (2012)	Funds	59	23.4%	56.2%	20.3%	
Filka (2013)	Reporting	45	53.3%	42.2%	4.4%	
Kleine, Krombauer, and Welter (2013)	E, S, G	182	30.8%	31.9%	7.7%	29.7%
Reveff and Viviani (2013)	Funds	75	24.0%	48.0%	14.7%	13.3%
Capelle-Blancard and Monjon (2014)	Funds	61	3.3%	47.5%	16.4%	32.8%
Clark, Feiner, and Viehs (2015)	E, S, G	110	85.5%	5.1%	0.9%	8.5%
Schröder (2014)	E, S	28	57.1%	7.1%	10.7%	25.0%
Total/weighted average		1,816	48.2%	23.0%	10.7%	18.0%

<그림 1> Meta Analysis from ESG and Financial Performance: Aggregated more than 2000 empirical studies [1]

ESG 기업 등급을 독립변수로 하여 F 검정을 진행한 결과이다.

이 역시 다수의 기업 성과 변수들이 ESG 와 관련 하여 높은 상관관계를 보이고 있으며, 특히 많은 성과가 ESG 와 관련하여 양의 상관관계를 가지고 있다는 점에서 유의미하다고 볼 수 있다.

<표 6> 초과이익률(Abnormal Stock Returns)

이 표는 각 등급의 초과이익률 및 등급 차이에 대하여 검증한 결과를 제시하고 있다. * : t 값을 의미하며, # : Wilcoxon Z-values를 의미함. *, **, ***는 각각 10%, 5%, 1% 유의수준에서 유의적임을 의미함

	AA 등급	E 등급	t값 [#]	Z값 [#]
전체기간	0.1885	0.0238	2.2747**	-2.114**
2008년도	0.3265	0.1756	1.5128*	-1.545
2009년도	0.2960	-0.0451	1.7407**	-1.801*
2010년도	0.0246	-0.0750	0.9291	-0.718
2008~2010년도	0.2139	0.0184	2.3415***	-2.330**
2011년도	0.1133	0.0389	0.4953	-0.041

<표 7> 영업성과(Operating Performance)

이 표는 각 등급의 영업성과 및 등급 차이에 대한 검증 결과를 제시하고 있다. * : t 값을 의미하며, # : Wilcoxon Z-values를 의미함. *, **, ***는 각각 10%, 5%, 1% 유의수준에서 유의적임을 의미함

	AA 등급	E 등급	t값 [#]	Z값 [#]
전체기간	0.1232	0.0592	5.6034***	-5.440***
2008년도	0.1225	0.0686	2.1067**	-2.343**
2009년도	0.1207	0.0696	2.4349***	-2.397**
2010년도	0.1332	0.0554	β.2055***	-3.399***
2008~2010년도	0.1255	0.0645	4.5176***	-4.459***
2011년도	0.1149	0.0432	3.2992***	-3.182***

<표 9> 토빈큐(Tobin's Q)

이 표는 각 등급의 토빈큐 및 등급 차이에 대한 검증 결과를 제시하고 있다. * : t 값을 의미하며, # : Wilcoxon Z-values를 의미함. *, **, ***는 각각 10%, 5%, 1% 유의수준에서 유의적임을 의미함

	AA 등급	E 등급	t값 [#]	Z값 [#]
전체기간	1.5139	0.9152	6.0436***	-7.563***
08년도	1.2814	0.8070	3.2936***	-3.440***
09년도	1.6213	0.9788	2.8937***	-3.887***
10년도	1.6910	0.8948	4.2150***	-4.293***
08~10년도	1.5279	0.8936	5.8277***	-6.884***
11년도	1.4460	0.9803	2.1980**	-3.209***

<그림 2> T-test and Z-test from 기업의 ESG 와 재무성과[2]

(표 10) 환경부문 평가등급에 따라 차이를 보이는 재무성과

시점	방향	구분	선택된 변수	F비	등급별 평균 ^{*)}			
					A+	A	B+	B이하
t (6개)	단조 증가	수익성	총자산순이익률	3.32**	5.64	3.93	2.4	1.24
		수익성	유보액/남입자본	21.99***	2987	2614	1756	1179
	현금흐름	영업활동CF/총자산	5.39***	7.73	6.46	5.25	3.75	
	단조 감소	수익성	금융비용대비용비율	2.23*	1.04	1.48	1.69	1.94
수익성	금융비용대부채비율	4.38***	1.75	1.96	2.05	2.34		
현금흐름	투자활동CF/총자산	3.97***	-7.36	-7.00	-5.55	-4.22		
t+1 (5개)	단조 증가	수익성	총자산순이익률	2.37*	4.3	3.03	1.28	1.19
		수익성	유보액/남입자본	24.11***	3234	3233	1825	1278
	수익성	금융비용대비용비율	2.86**	1.77	1.83	1.96	2.18	
	단조 감소	안정성	유동비율	5.32***	142.9	143.9	172.1	198.4
현금흐름	투자활동CF/총자산	2.46*	-7.06	-6.35	-4.81	-3.96		
t+2 (3개)	단조 증가	수익성	유보액/남입자본	20.78***	5124	2712	1937	1318
	단조 감소	현금흐름	영업활동CF/총자산	5.15***	7.45	6.57	6.33	4.19
안정성	영업자산대총자산	28.51***	19.06	19.97	24.85	29.15		

<그림 3>F-test from

기업의 ESG 노력과 재무성과의 선행적 관계 : 탐색적 연구[3]

나. 연구 방법

본 연구에서는 한국기업지배구조원의 '환경모범규준'[7], '사회모범규준'[8], '기업지배구조 모범규준'[9] 보고서에서 각각 LDA(Latent Dirichlet Allocation)[10] 알고리즘의 토픽 모델링[11] 방식을 통해 E(환경), S(사회), G(지배구조) 요소의 키워드를 도출하였다.

특정기업의 뉴스 데이터를 월 단위로 크롤링 후 ESG 요소에 대해 긍정, 부정, 중립 라벨링을 끝내고 Naive Bayesian 을 활용한 지도학습을 진행하였다. 중립을 제외한 긍, 부정뉴스 데이터에 대해 ESG 요소 키워드의 빈도수를 EG, EB, SG, SB, GG, GB (G: 긍정, B: 부정)로 계산하였다.

마지막으로 도출된 빈도수와 주가(해당 월 종가 기준)와의 OLS Regression[12]을 통해 p-value 값이 0.05 미만인 ESG 요소를 도출하여 ESG 요소와 주가의 상관관계를 분석하였다.

다. 실험 프레임 구조

가. ESG 키워드 추출 및 공부정 키워드 분류

한국기업지배구조원의 '환경모범규준', '사회모범규준', '기업지배구조 모범규준' 보고서에서 텍스트를 추출하였다.

불용어 제거를 위해 형태소 분석 라이브러리를 사용하여 명사로 토큰화를 거쳤다. TF-IDF 행렬을 만들기 위해 역 토큰화 과정을 거친 후 각각의 보고서마다 토픽 한 개당 300 개의 키워드를 추출하였다. 키워드를 GRI INDEX[14]를 참고하여 ESG 요소와 관계없는 불용어를 제거하였으며 중복된 단어의 경우 TF-IDF 를 토타로 중요도의 정도에 따라 중복 제거를 진행하였다. 단어 분류 결과의 예시는 <그림 4> 와 같다.

E_word	(‘환경’, 454.13), (‘환경성’, 89.91), (‘참여’, 37.11), (‘개선’, 32.71), (‘영향’, 32.18), (‘발생’, 28.41), (‘지속’, 27.2), (‘환경문제’, 17.58)...
S_word	(‘사회’, 177.96), (‘이슈’, 10.02), (‘소비자’, 96.85), (‘제품’, 85.45), (‘서비스’, 61.16), (‘제공’, 56.18)...
G_word	(‘감사’, 259.94), (‘커뮤니케이션’, 16.37), (‘공시’, 116.08), (‘위원회’, 108.57), (‘주주’, 94.86)...

<그림 4>ESG Word (단어 옆 숫자는 TF-IDF 점수)

나. 뉴스 데이터 크롤링, 공부정 뉴스 분류 및 빈도수 분석

뉴스 데이터를 수집하기 위해 ESG와 관련된 대표적인 1개 기업의 뉴스 데이터를 2019년부터 2021년 3월까지 월별 데이터를 크롤링하였다. 크롤링한 뉴스 데이터를 긍정, 부정, 중립 뉴스로 분류하기 위해 월별로 100개의 뉴스를 라벨링, 총 2000개의 학습데이터를 만들었다. 그 후 Naive Bayesian 분류 알고리즘을 통해 긍정, 부정, 중립 뉴스로 학습시킨 후 모든 뉴스를 중립을 제외하고 긍정, 부정으로 분류하였다.

ESG 키워드 추출 및 분류된 공부정 뉴스를 기반으로 월별(M)로 E(환경) 관련 공부정(EG, EB), S(사회) 관련 공부정(SG, SB), G(지배구조) 관련 공부정(GG, GB)의 빈도수를 계산하였다. 기업 주가 데이터는 ‘네이버 금융 국내증시’를 활용하였으며, 월별 기업의 종가(CA)를 사용하였다. 결과는 <표 1>과 같다.

M	EG	EB	SG	SB	GG	GB	CA
0	109	577	434	893	183	559	632000
1	214	946	333	773	356	2079	627000
2	127	373	286	472	207	606	613000
3	62	4076	205	3307	89	6882	614000
4	191	764	407	1300	211	1328	571000
5	182	705	446	841	354	930	563000
6	275	662	474	758	284	896	536000
7	271	400	438	424	179	537	510000
8	153	340	428	503	170	650	491500
9	160	274	285	567	255	762	470500
10	226	318	461	531	176	607	456500
11	254	230	361	415	280	397	439500
12	193	302	528	387	390	514	411000
13	75	168	150	237	105	297	378500
14	276	194	520	285	279	420	294000
15	314	239	511	349	342	370	321000
16	416	346	805	971	753	1003	299500
17	212	162	444	486	163	341	283500
18	317	143	582	298	281	251	304000
19	305	243	825	364	283	270	272500
20	340	142	613	195	305	194	257500
21	401	97	378	147	168	145	255500
22	274	80	557	114	191	141	288000
23	649	312	1199	171	834	371	283000
24	341	726	591	878	281	1069	287000
25	390	169	1180	183	413	292	295000
26	611	183	642	163	665	253	303000

<표 1>

다. 기간별 뉴스데이터의 ESG 요소 빈도수와 주가의 관계분석

3-2의 과정을 통해 분류된 E공부정 빈도수(EG, EB), S공부정 빈도수(SG, SB), G공부정 빈도수(GG, GB)들과 이 기간의 주가와 OLS Regression을 진행하였으며 결과는 다음과 같다.

OLS Regression Results

```

=====
Dep. Variable:          CA  R-squared:          0.634
Model:                  OLS  Adj. R-squared:       0.524
Method:                 Least Squares  F-statistic:         5.768
Date:                   Wed, 29 Sep 2021  Prob (F-statistic):  0.00128
Time:                   16:39:08  Log-Likelihood:      -342.98
No. Observations:      27  AIC:                  700.0
Df Residuals:          20  BIC:                  709.0
Df Model:               6
Covariance Type:       nonrobust
=====
              coef  std err      t  P>|t|  [0.025  0.975]
-----
const      5.552e+05  6.41e+04   8.662  0.000  4.21e+05  6.89e+05
EG         -764.8348   270.768  -2.825  0.010 -1329.647 -200.023
EB          294.0537   192.181   1.530  0.142  -106.829  694.937
SG         -191.1495   122.146  -1.565  0.133  -445.942  63.643
SB           11.0235   110.119   0.100  0.921  -218.681  240.728
GG          456.4681   193.001   2.365  0.028   53.875  859.061
GB         -154.3126   105.256  -1.466  0.158  -373.873  65.248
=====
    
```

그 결과 Adj.R-squared는 0.524로 높지만 P-value가 0.05보다 낮은 유의한 변수는 EG와 GG뿐인 관계로 최적의 변수 조합을 선정(E_n = EG/EB, S_n = SB/SG, G_n = GG/GB)을 사용하여 OLS Regression을 재진행하였으며 결과는 다음과 같다.

OLS Regression Results

```

=====
Dep. Variable:          CA  R-squared:          0.667
Model:                  OLS  Adj. R-squared:       0.624
Method:                 Least Squares  F-statistic:         15.35
Date:                   Wed, 29 Sep 2021  Prob (F-statistic):  1.06e-05
Time:                   16:23:43  Log-Likelihood:      -341.70
No. Observations:      27  AIC:                  691.4
Df Residuals:          23  BIC:                  696.6
Df Model:               3
Covariance Type:       nonrobust
=====
              coef  std err      t  P>|t|  [0.025  0.975]
-----
const      3.78e+05  5.56e+04   6.793  0.000  2.63e+05  4.93e+05
E_n       -4.511e+04  2.12e+04  -2.132  0.044  -8.89e+04 -1334.510
S_n        1.145e+05  4.5e+04   2.546  0.018  2.15e+04  2.08e+05
G_n       -2.072e+04  8864.205  -2.338  0.028  -3.91e+04 -2384.494
=====
    
```

그 결과 Adj.R-squared는 0.624로 상승하였으며 세가지 변수(E_n, S_n, G_n) 모두 P-value가 0.05미만으로 유의한 변수로 도출되었다.

3. 결론

가. 연구성과

본 연구는 기존의 다양한 연구사례들이 이미 존재하고 있던 ‘ESG 기업평가등급’이나 다른 지표들을 사용한 것과는 다르게 독립적인 지표를 개발하여 사용했다는 점에서 큰 차별점을 가진다. 특히 뉴스데이터 마이닝을 활용해 모델을 개발한 점은 연구 목적뿐만 아니라 여러 기업의 현실적인 기업 성과 분석에 쉽게 활용될 수 있으며 독자적 모델의 개발의 지표가 될 수 있다는 점에서 주목할 만하다.

하지만 본 연구는 데이터의 양이 충분히 크지 못했다는 점과 ESG 지표들의 주가에 대한 영향의 연쇄성을 고려하지 못했다는 점에서 한계를 갖는데, 이점들은 향후 연구를 어떻게 진행해야 할지에 대한 방향성을 제시한다.

<표 1>를 보면 M-3 시기에 EB, SB, GB 가 고점을 찍고 M-4 시기의 CA 가 M-3 에 비해 43000 이 하락, 그 후 CA 가 지속적으로 하락하고 있는 모습에서 ESG 의 B 지표들의 영향력이 추후 시기에도 연쇄적으로 영향을 미치는 것으로 고려된다. 그렇기에 추후 본 연구의 표본의 개수를 증대하여 모델을 더 견고히 하고 지표 영향의 누적성과 연쇄성을 고려한다면 ESG 에 관심도와 중요도가 높은 현대 기업들에게 있어 매력적이고 완성도가 높은 모델이 될 것이다.

나. 연구활용

추후 실질적으로 기업이 해당 모델의 사용 혹은 서비스화를 하여 기업의 매출 성과 및 주가 예측에 활용하고 ESG 사건 발생 시 소비자가 해당 모델을 통해 매출 성과 및 주가 예측하는 등의 다양한 분야에서 활용될 수 있을 것으로 기대된다.

참고문헌

- [1] Gunnar Friede, Timo Busch & Alexander Bassen, ESG and financial performance: aggregated, Journal of Sustainable Finance & Investment, 5권, 제4호 210-233, 2015
- [2] 장승욱, 김용현, 기업의 ESG 와 재무성과, 재무관리연구, 30 권, 제 1 호, 131-152, 2013
- [3] 김범석, 민재형, 기업의 ESG 노력과 재무성과의 선후행 관계: 탐색적 연구, 한국생산관리학회지, 27 권, 제 4 호, 513-538, 2016
- [4] 위경우, 강운식, 정재만, 이재현, 펀드의 ESG 수준이 펀드의 성과 및 현금흐름에 미치는 영향, Financial Planning Review, 13 권, 제 2 호, 83-115, 2020
- [5] 나영, 임옥빈, ESG 정보와 가치관련성에 관한 연구, 한국경영교육학회 학술발표대회논문집, 한국경영교육학회, 2011, 445-459
- [6] 강원, 정무권, 비재무지표와 기업의 시장성과 간의 관계에 대한 연구: ESG 지표 개발에 사용되는 사건의 시장반응 분석, 연세경영연구, 57 권, 제 2 호, 1-22, 2020
- [7] 기업사회책임위원회, 환경모범기준, 한국기업지배구조원, 2010
- [8] 기업사회책임위원회, 사회모범기준, 한국기업지배구조원, 2010
- [9] 기업사회책임위원회, 기업지배구조 모범기준, 한국기업지배구조원, 2016
- [10] David M.Blei, Andrew Y.Ng, Michael I.Jordan, Latent Dirichlet Allocation, Journal of Machine Learning Research 3, 993-1022, 2003
- [11] David M.Blei, Andrew Y.Ng, Michael I.Jordan, Latent Dirichlet Allocation, Journal of Machine Learning Research 3, 993-1022, 2003
- [12] A. Genkin, D.Lewis, D.Madigan, Large-Scale Bayesian Logistic Regression for Text Categorization, Technometrics, 49 권, 제 3 호, 291-304, 2007

※ 본 논문은 과학기술정보통신부 정보통신창의인재양성사업의 지원을 통해 수행한 ICT멘토링 프로젝트 결과물입니다.