



# 빅데이터 방법론 기반 연구사례

2014. 09. 29

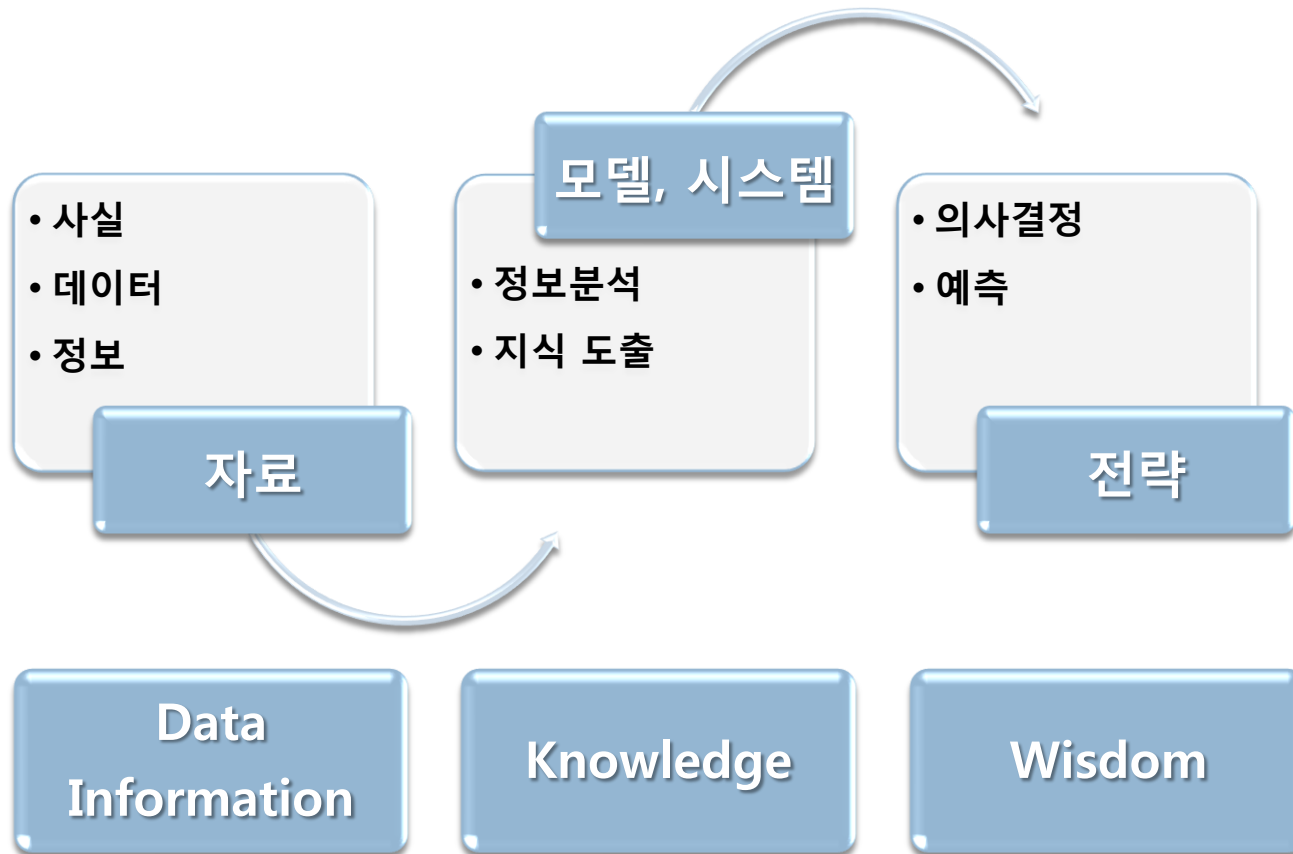
한 승 우

인하대학교 공과대학 건축공학과

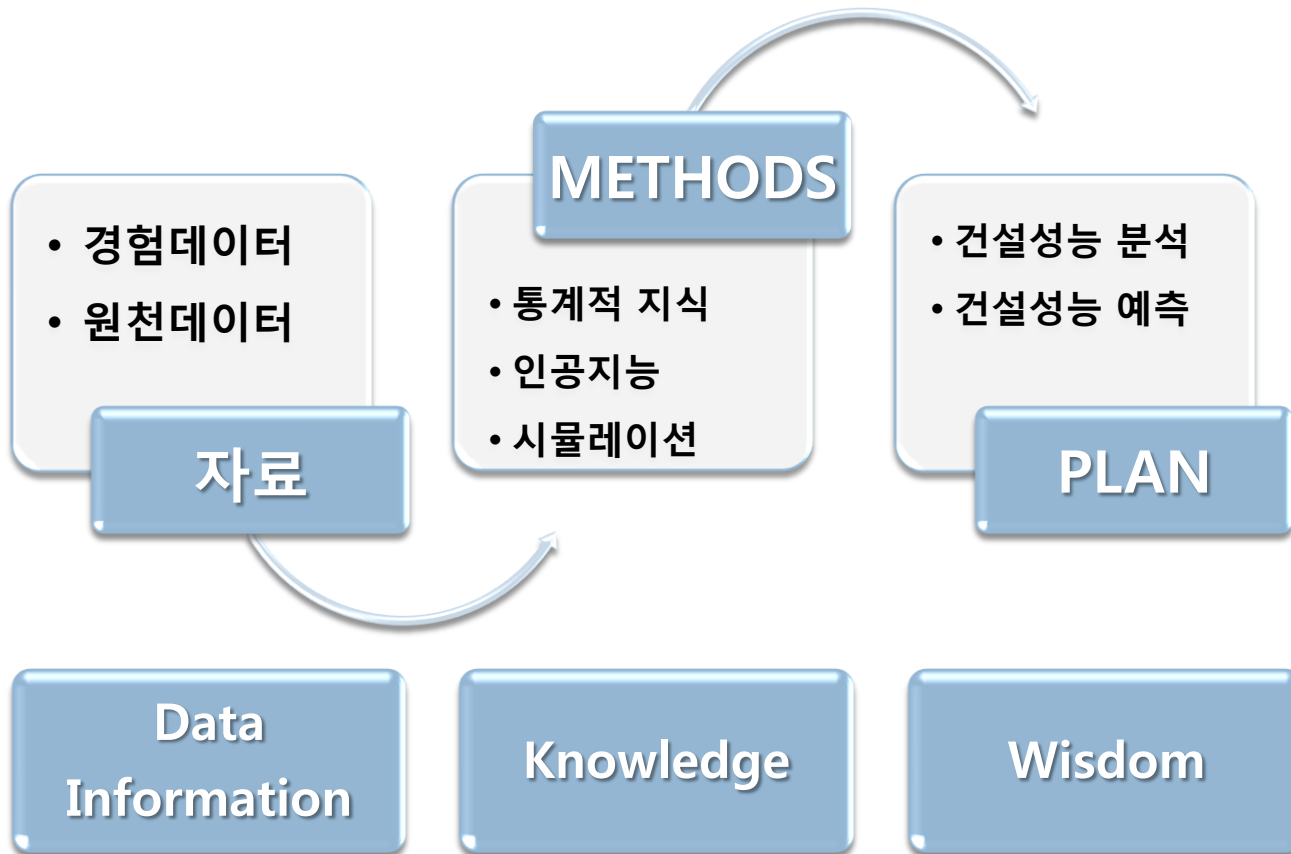


인하대학교

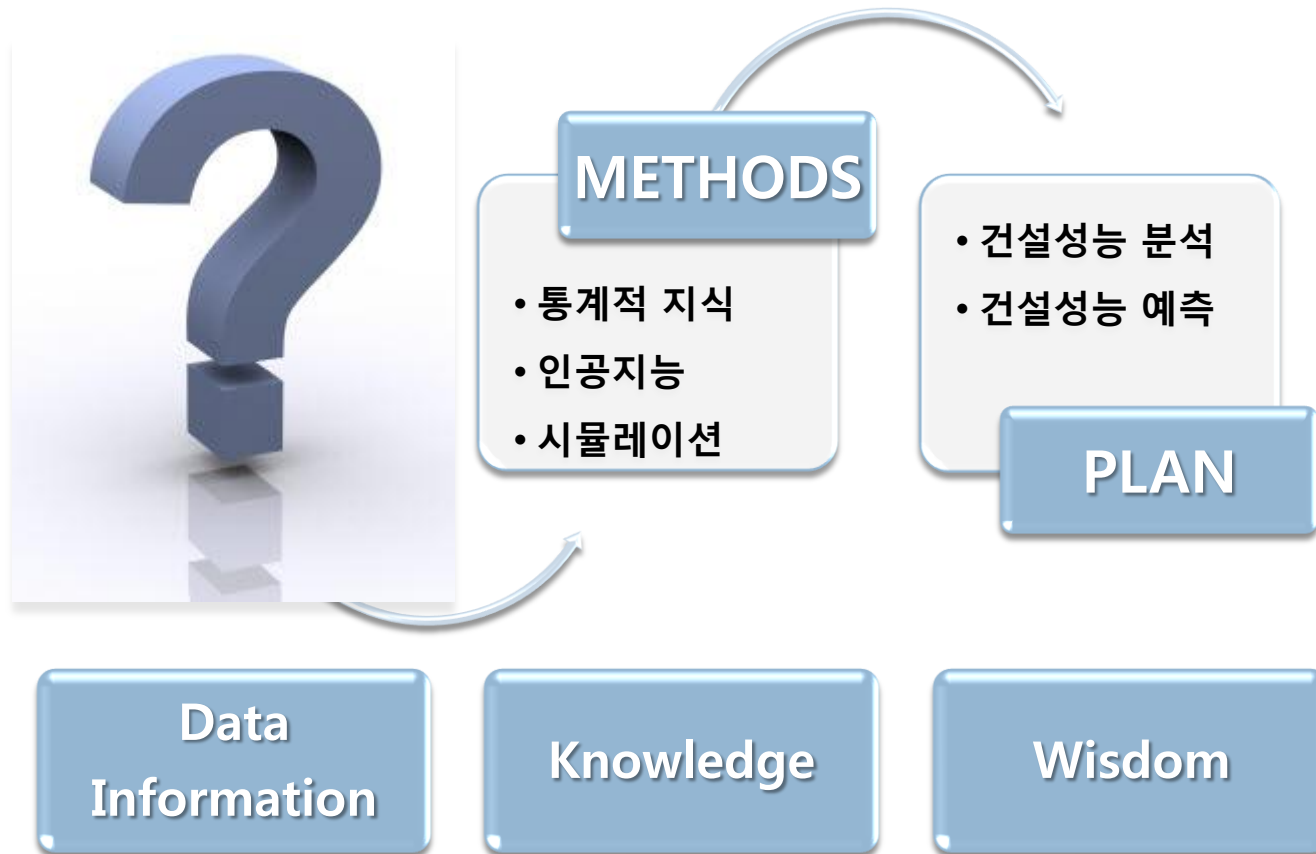
# 개요 1. 데이터기반 지식시스템



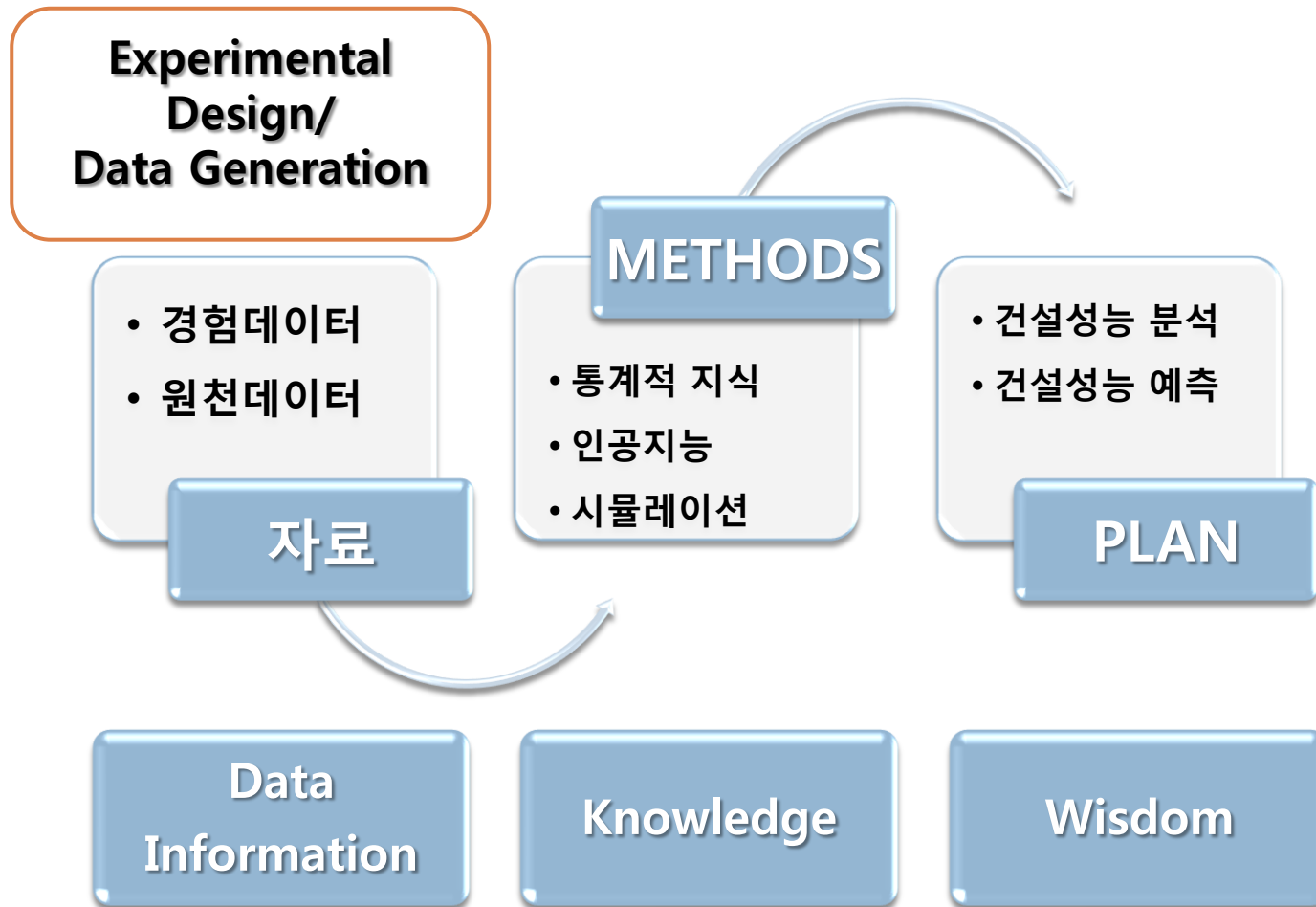
## 개요 2. 건설분야의 지식시스템



## 개요 2. 건설분야의 지식시스템

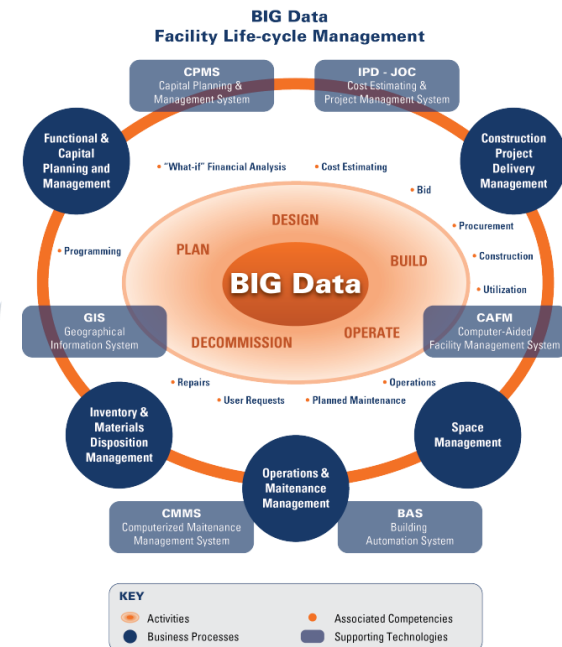
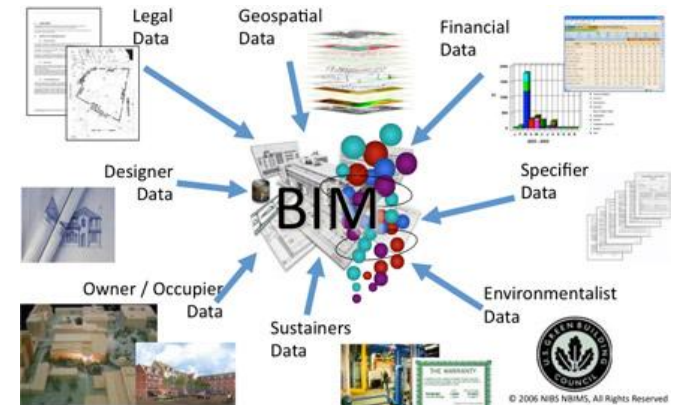
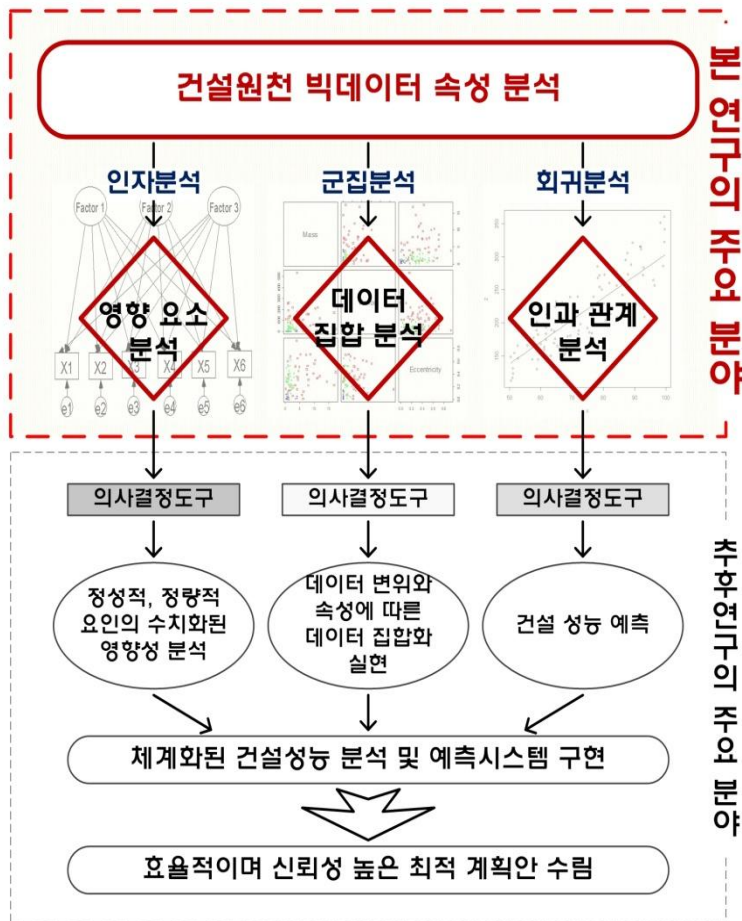


## 개요 2. 건설분야의 지식시스템

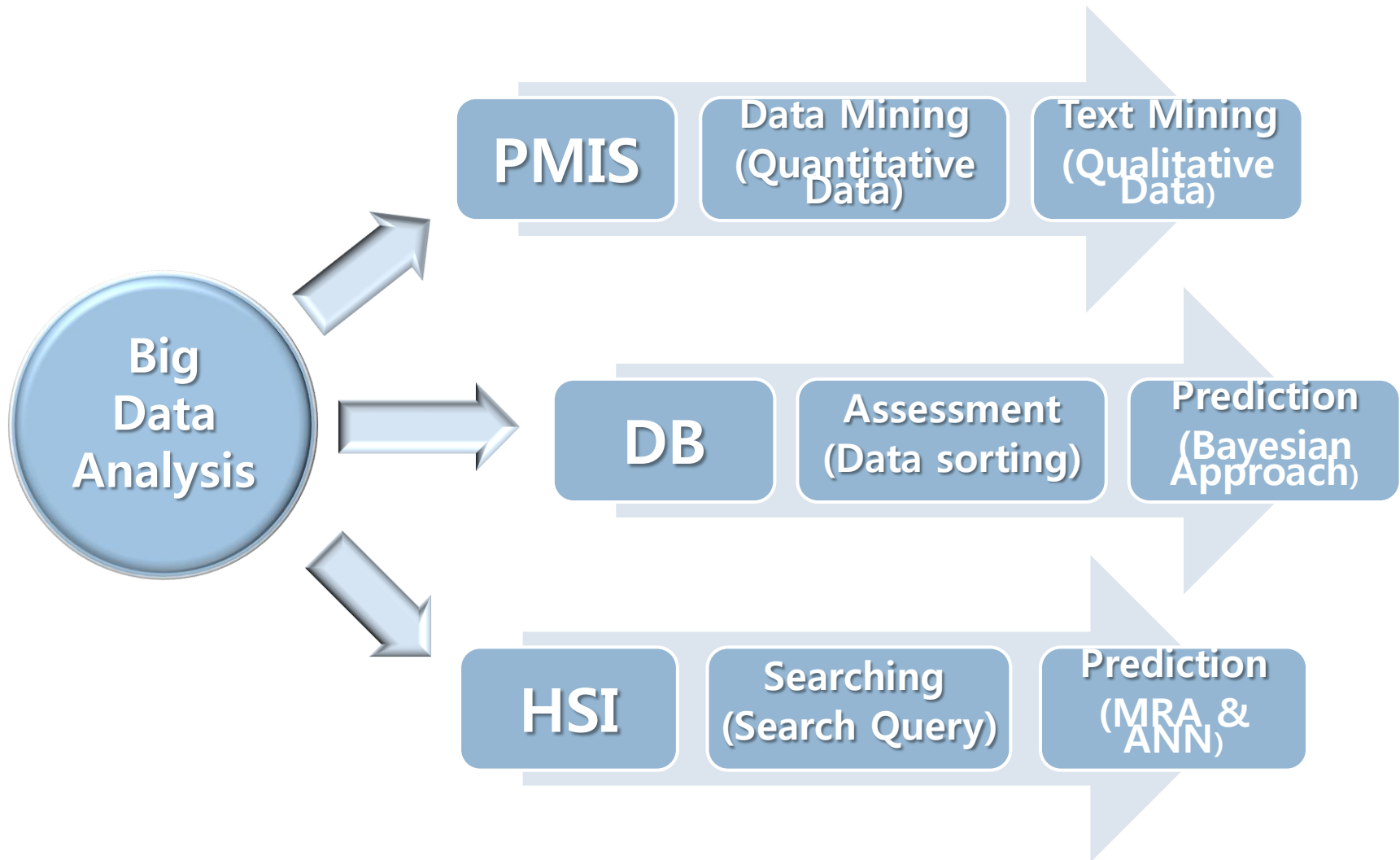


# 빅데이터 분석기반 연구 (수행)

건설원천 빅데이터 속성 분석모형 개발  
한국연구재단 중견연구자지원사업(핵심연구: 개인)  
2012.12~2015.11



# 빅데이터 개념 기반 연구과제 (수행)



## Actual Data Assessment

Daily words counting.	17	9	Calculation of average and sum.
	17	*	
	*	*	
	*	*	



# TEXT MINING ON PMIS

## Actual Data Assessment

104

703

[illegible]

# TEXT MINING ON PMIS

## Monthly and Weekly Basis Data Summation

	마감	바닥	배근작업	배근	배근sum	벽체	설치	설치sum	슬라브	자재	자재sum	미장	미장sum	방수	방수sum	작업	작업sum
2011년 4월 30일	0	0	0	0	0	0	6	6	0	0	1	0	0	0	0	0	4
2011년 5월 31일	0	0	0	0	0	0	8	34	0	0	0	0	0	0	0	2	24
2011년 6월 30일	0	0	0	0	0	0	110	128	0	3	3	0	0	0	0	36	62
2011년 7월 31일	1	0	0	0	0	0	136	158	0	0	0	0	0	0	0	37	94
2011년 8월 31일	0	1	0	0	0	0	179	213	0	2	2	0	0	0	0	9	56
2011년 9월 30일	2	19	0	0	0	0	19	19	0	4	4	0	0	0	0	42	65
2011년 10월 31일	16	10	11	0	11	14	63	65	50	4	5	2	3	0	0	121	191
2011년 11월 30일	12	20	13	0	13	8	31	34	44	37	37	2	3	0	0	128	226
2011년 12월 31일	11	47	17	0	17	27	60	60	46	25	25	5	5	0	0	206	272
2012년 1월 31일	24	66	17	0	17	20	91	99	42	13	13	1	1	0	0	247	336
2012년 2월 29일	13	88	11	0	11	29	108	110	47	13	13	14	14	0	0	263	394
2012년 3월 31일	5	99	0	23	23	25	77	77	40	36	37	35	35	5	22	300	348
2012년 4월 30일	20	120	0	39	39	28	103	103	23	19	19	61	61	0	31	347	382
2012년 5월 31일	19	66	0	26	26	12	161	161	5	20	20	77	77	12	53	408	446
2012년 6월 17일	4	7	0	8	8	13	85	85	1	11	17	45	45	5	13	247	275
2012년 6월 30일	0	0	15	1	16	28	64	103	0	14	21	0	34	0	20	0	148
2012년 7월 31일	0	1	2	0	2	75	189	247	0	53	86	2	75	13	41	10	278
2012년 8월 24일	6	5	1	4	5	33	125	209	0	27	50	4	64	13	15	41	208
2012년 8월 31일	0	6	0	1	1	5	41	105	0	5	5	10	10	11	12	119	128
2012년 9월 30일	6	17	0	2	2	38	158	283	0	41	41	23	23	32	36	304	348
2012년 10월 31일	27	29	0	3	3	101	286	310	0	26	26	17	17	42	43	531	562
2012년 11월 10일	5	18	0	0	0	23	92	95	0	14	14	4	4	11	11	153	159
2012년 11월 30일	8	35	0	0	0	72	258	266	0	6	13	15	15	12	15	464	491
2012년 12월 31일	2	79	0	0	0	67	340	349	0	16	19	10	10	2	6	497	534
2013년 1월 31일	1	7	0	0	0	15	155	163	0	11	13	7	7	0	1	247	268
2013년 2월 28일	0	3	0	0	0	0	14	14	0	1	1	1	1	0	0	26	32
	182	743	87	107	194	633	2959	2959	298	401	401	335	335	158	158	4785	##

## Monthly and Weekly Basis Data Summation

[illegible]

# TEXT MINING ON PMIS

## Correlation Analysis using R

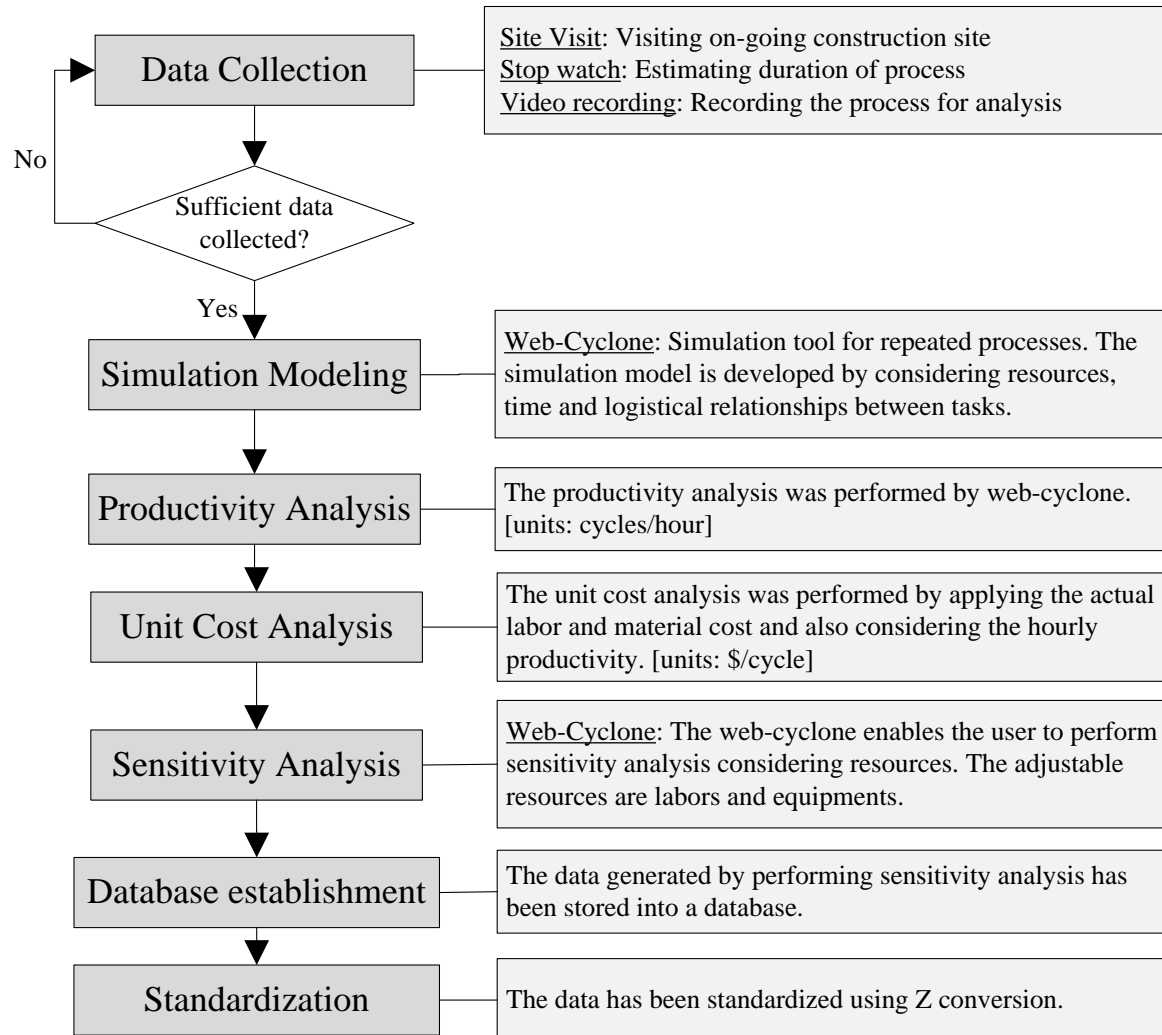
Correlation analysis between the word data and PMIS actual data has been performed.

Total 879 correlations which the correlation coefficient is higher than 0.7 has been found.

**The meaning of each correlation should be analyzed for future work.**

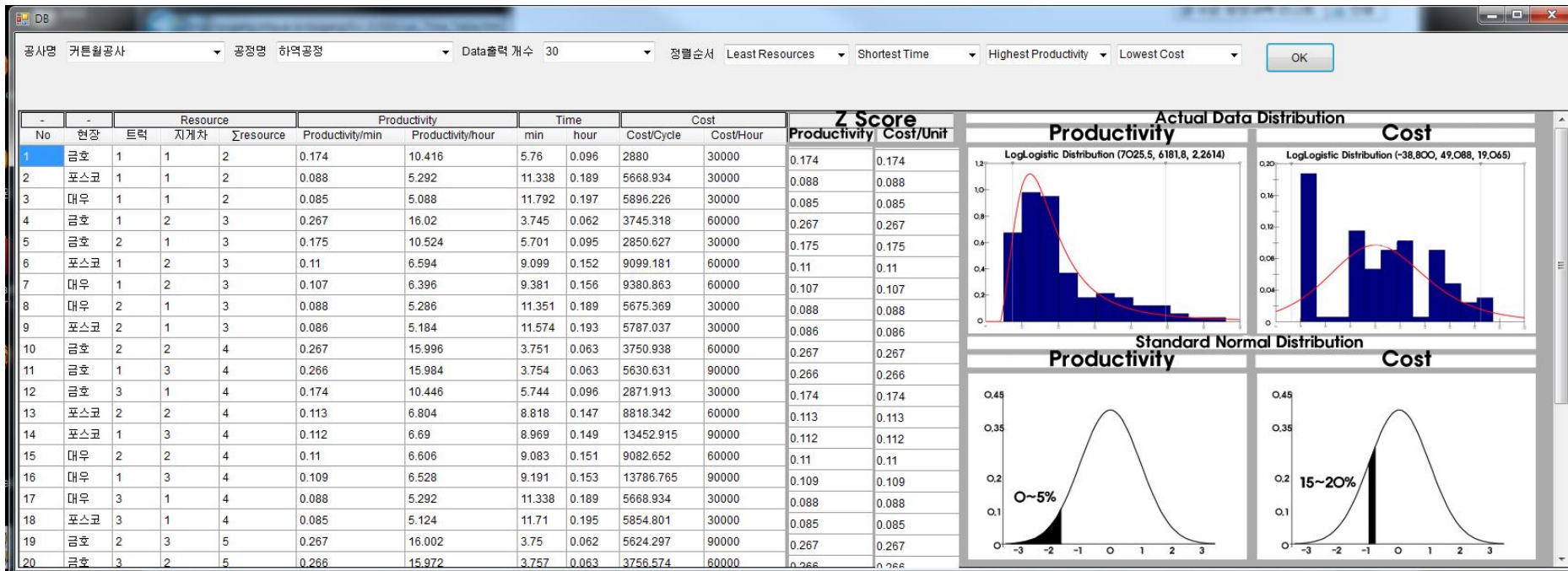
77	97	sum	미장26	0.704338	0
77	100	sum	sum	0.821362	0
77	122	sum	반장37	0.867587	0
77	139	sum	미장49	0.918647	0
77	155	sum	마감	0.720157	CHECK
77	156	sum	바닥	0.905607	CHECK
77	157	sum	배근작업	0.920721	CHECK
77	167	sum	미장sum	0.716617	CHECK
77	168	sum	방수	0.723242	CHECK
77	169	sum	방수sum	0.791864	CHECK
77	170	sum	작업	0.786195	CHECK
78	9	반장21	25-210-125	0.725001	0
78	69	반장21	유리	0.737675	0
78	70	반장21	기계설비	0.821058	0
78	71	반장21	닥트	0.894854	0
78	72	반장21	반장	0.872806	0
78	73	반장21	용접	0.855339	0
78	77	반장21	sum	0.843507	0
78	79	반장21	직원22	0.953356	0
78	80	반장21	금속	0.971942	0
78	81	반장21	sum	0.982097	0
78	90	반장21	전기설비	0.789253	0
78	94	반장21	sum	0.754725	0
78	97	반장21	미장26	0.828474	0
78	100	반장21	sum	0.797232	0
78	112	반장21	잡철물	0.762225	0
78	122	반장21	반장37	0.92403	0
78	123	반장21	수장	0.794372	0
78	125	반장21	sum	0.790979	0
78	139	반장21	미장49	0.826488	0
78	154	반장21	방수	0.885654	0
78	156	반장21	바닥	0.826675	CHECK
78	157	반장21	배근작업	0.839301	CHECK
78	168	반장21	방수	0.701581	CHECK
78	169	반장21	방수sum	0.764266	CHECK
78	170	반장21	작업	0.75813	CHECK
79	2	직원22	25-180-12	0.710533	0

# DATA ASSESSMENT USING STATISTICAL APPROACHES



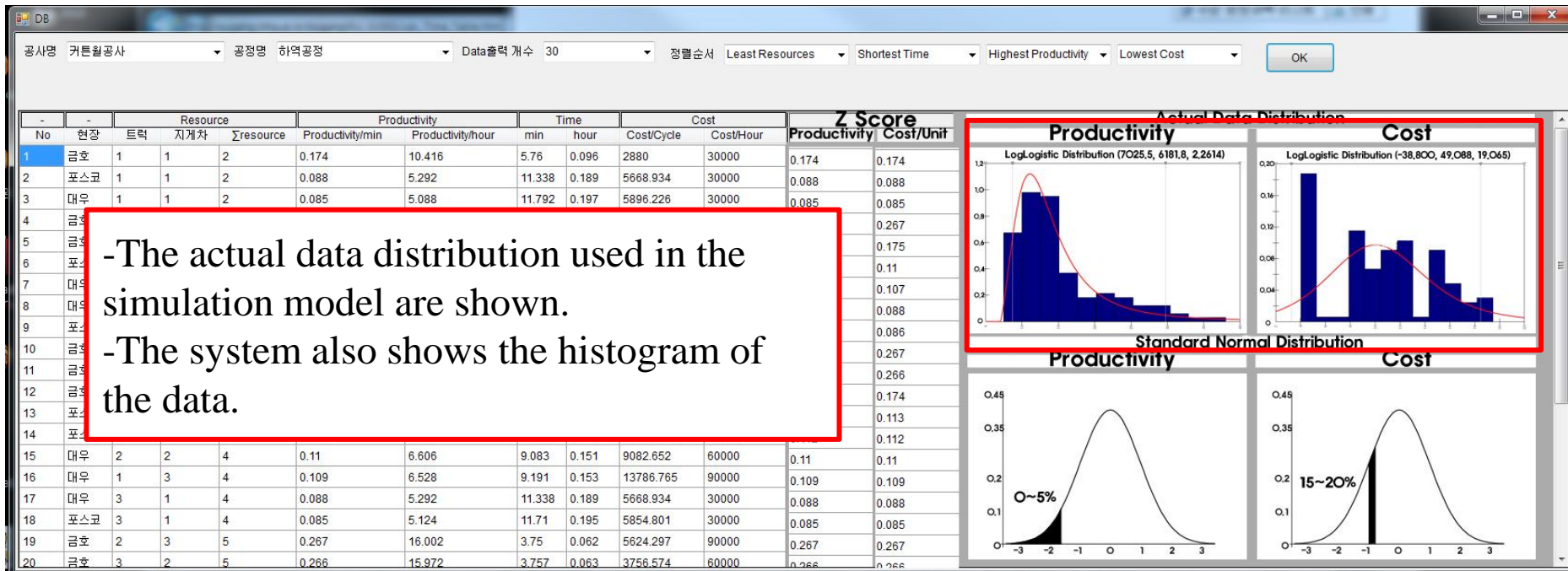
# DATA ASSESSMENT USING STATISTICAL APPROACHES

## Overviews of the DB



# DATA ASSESSMENT USING STATISTICAL APPROACHES

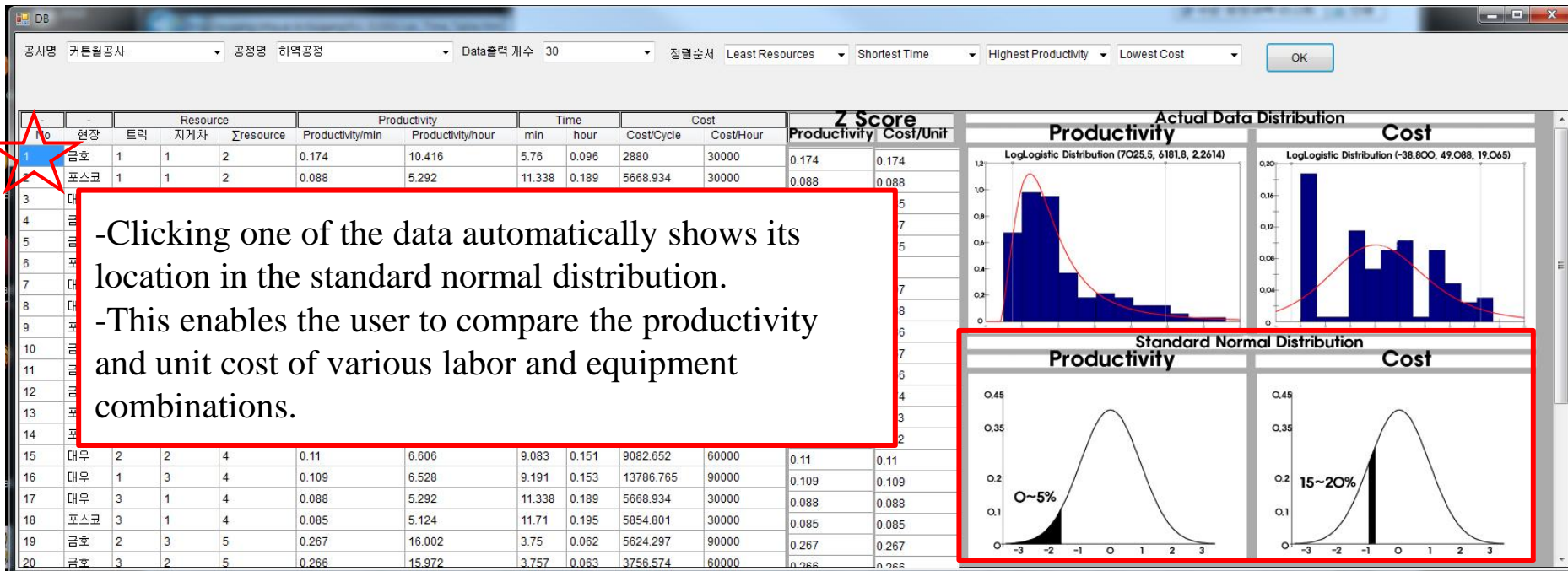
## Illustration of Data Distribution





# DATA ASSESSMENT USING STATISTICAL APPROACHES

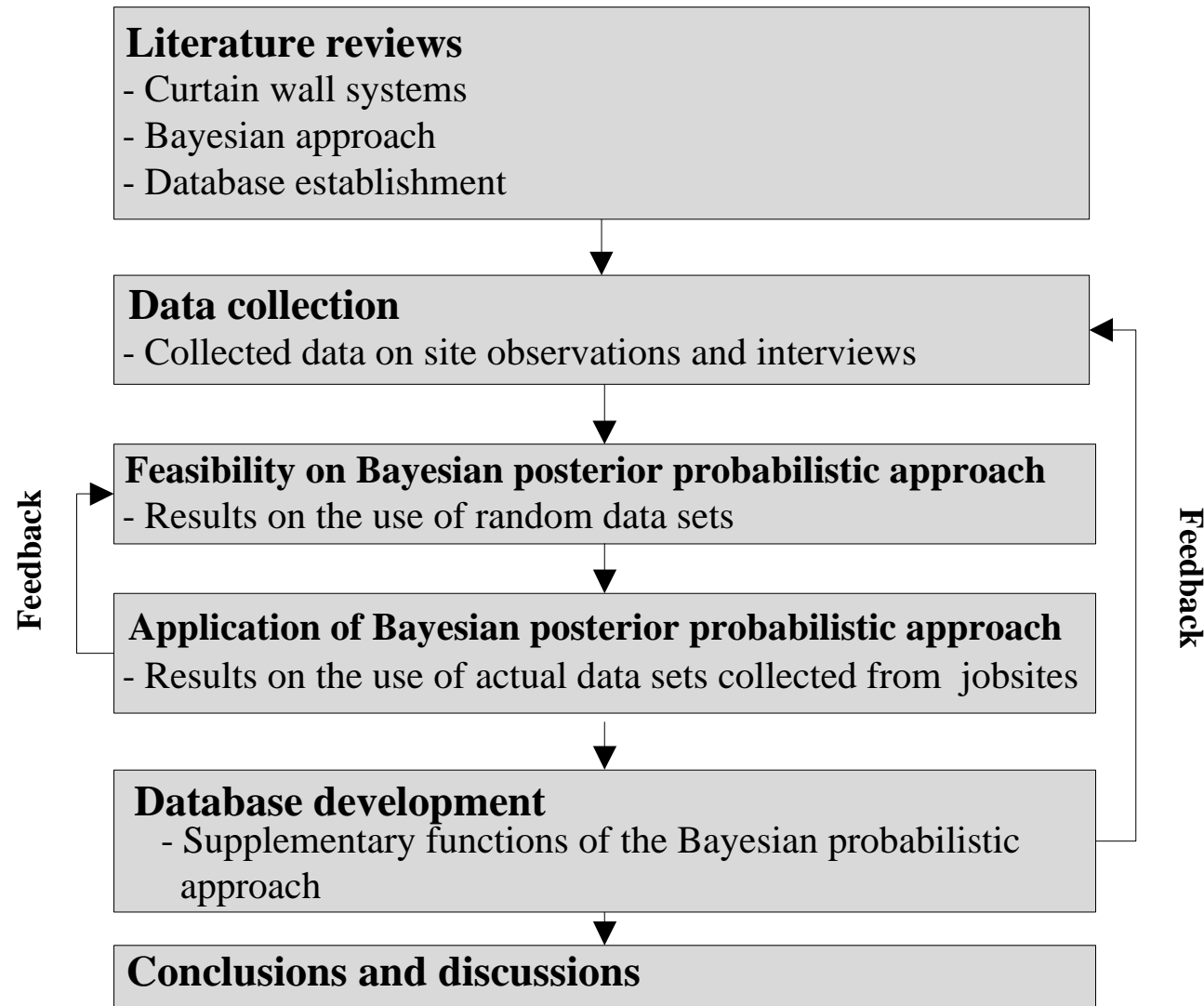
## Illustration of Statistical Position





# INFORMATION PREDICTION USING STATISTICAL APPROACHES

## Research method



# INFORMATION PREDICTION USING STATISTICAL APPROACHES

## Interface for Data Input

The screenshot shows the 'Settings' tab of the 'Interface for Data Input' application. The interface is divided into several sections:

- Data Retrieve:** Includes dropdown menus for 'Constructor', 'Site', and 'Process', along with 'Go' and 'New' buttons.
- Result:** Contains a 'Program logo' label.
- Default Setting:** Features input fields for 'construction', 'site', and 'process'. It also includes a 'Default Set' field with the value '0' and a 'Level' dropdown menu set to 'High'.
- Display Setting:** Divided into two sub-sections:
  - Basic Values:** Lists variables: productivity, labor, equipment, material\_cost, labor\_cost, and equipment\_cost, each with an unchecked checkbox.
  - User-Defined Values:** Includes a checkbox for 'unitcost', which is currently checked.

The user can add his own formula using the variables defined above.

The screenshot shows the 'User-defined Values' section of the application. It includes a table with two columns: 'Entity1' and 'Formula1'. The first row shows 'unitcost' in the 'Entity1' column and '(labor\_cost+equipment\_cost)' in the 'Formula1' column. Below the table are '+' and '-' buttons, and a 'Save' button.

# INFORMATION PREDICTION USING STATISTICAL APPROACHES

## Interface for Data Retrieve

Data Retrieve

In data retrieve the user can insert the actual data according to the variables defined in previous steps.

NUM	SET	LEVEL	productivity
1	1	1	4.52
2	1	1	4.88
3	1	1	2.89
4	2	1	2.36
5	2	1	5.50
6	2	1	4.32
7	2	1	3.81
8	2	1	4.17
9	2	1	4.56
10	2	1	4.62
11	2	1	4.86
12	3	1	7.84
13	3	1	8.91
14	3	1	8.35
15	3	1	3.41

Detail Information

NUM  SET   
LEVEL  productivity

UPDATE

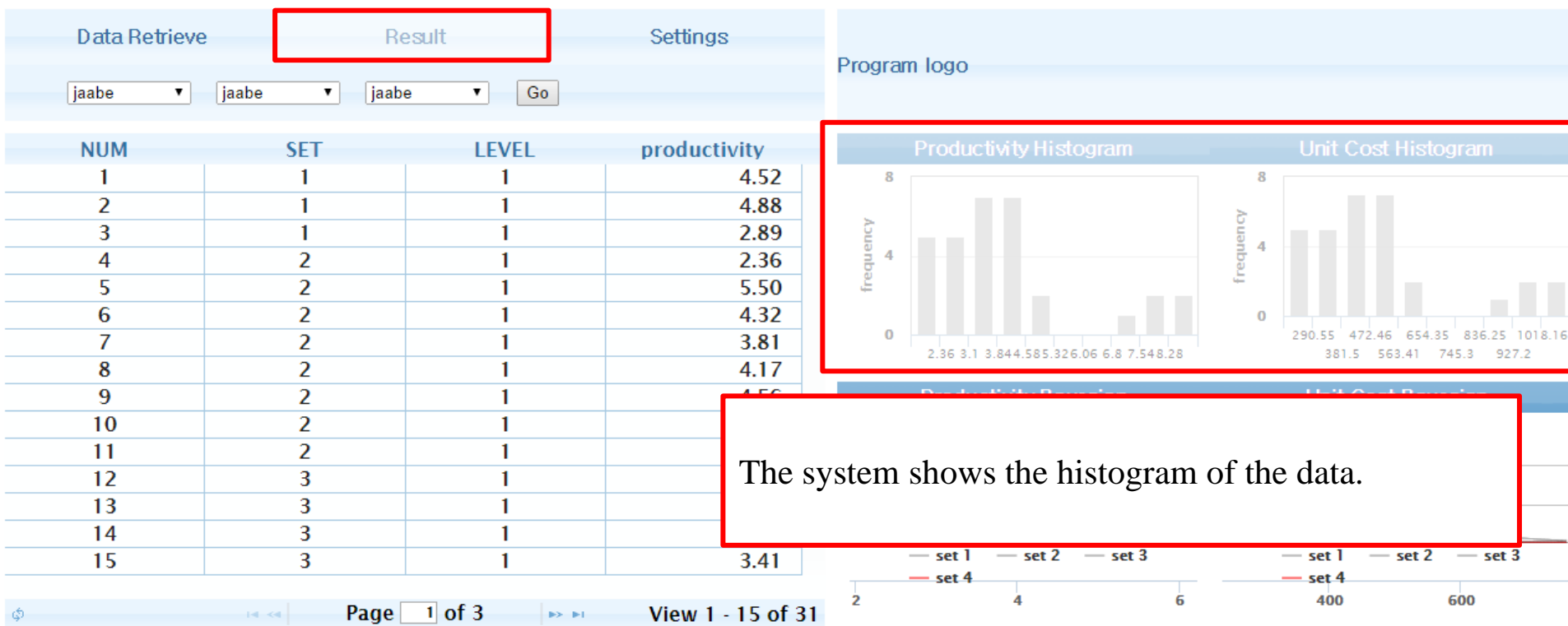
DELETE

Page  of 3

View 1 - 15 of 31

# INFORMATION PREDICTION USING STATISTICAL APPROACHES

## Illustration of the Analyzed Results



# INFORMATION PREDICTION USING STATISTICAL APPROACHES

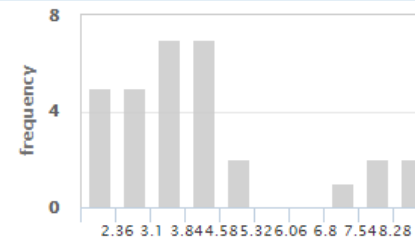
## Illustration of the Analyzed Results

Data Retrieve	Result	Settings	
jaabe	jaabe	jaabe	
<input type="button" value="Go"/>			
NUM	SET	LEVEL	productivity
1	1	1	4.52
2	1	1	4.88
3	1	1	2.80
This is the analysis results using the Bayesian probabilistic approach.			
In this example there are four data sets, therefore, it can be seen that there are four distributions.			
The red distribution represents the result.			
13	3	1	8.91
14	3	1	8.35
15	3	1	3.41

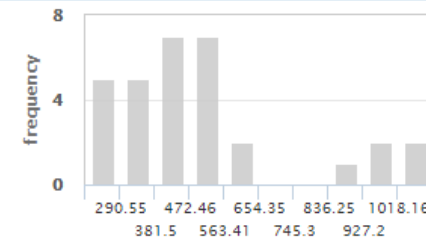
This is the analysis results using the Bayesian probabilistic approach.  
In this example there are four data sets, therefore, it can be seen that there are four distributions.  
The red distribution represents the result.

Program logo

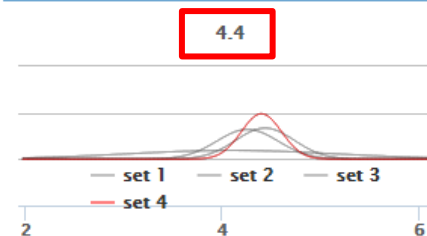
Productivity Histogram



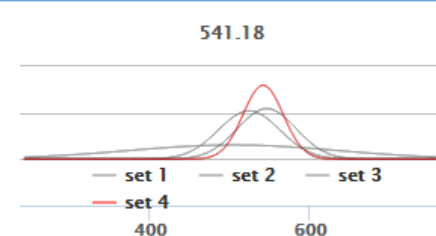
Unit Cost Histogram



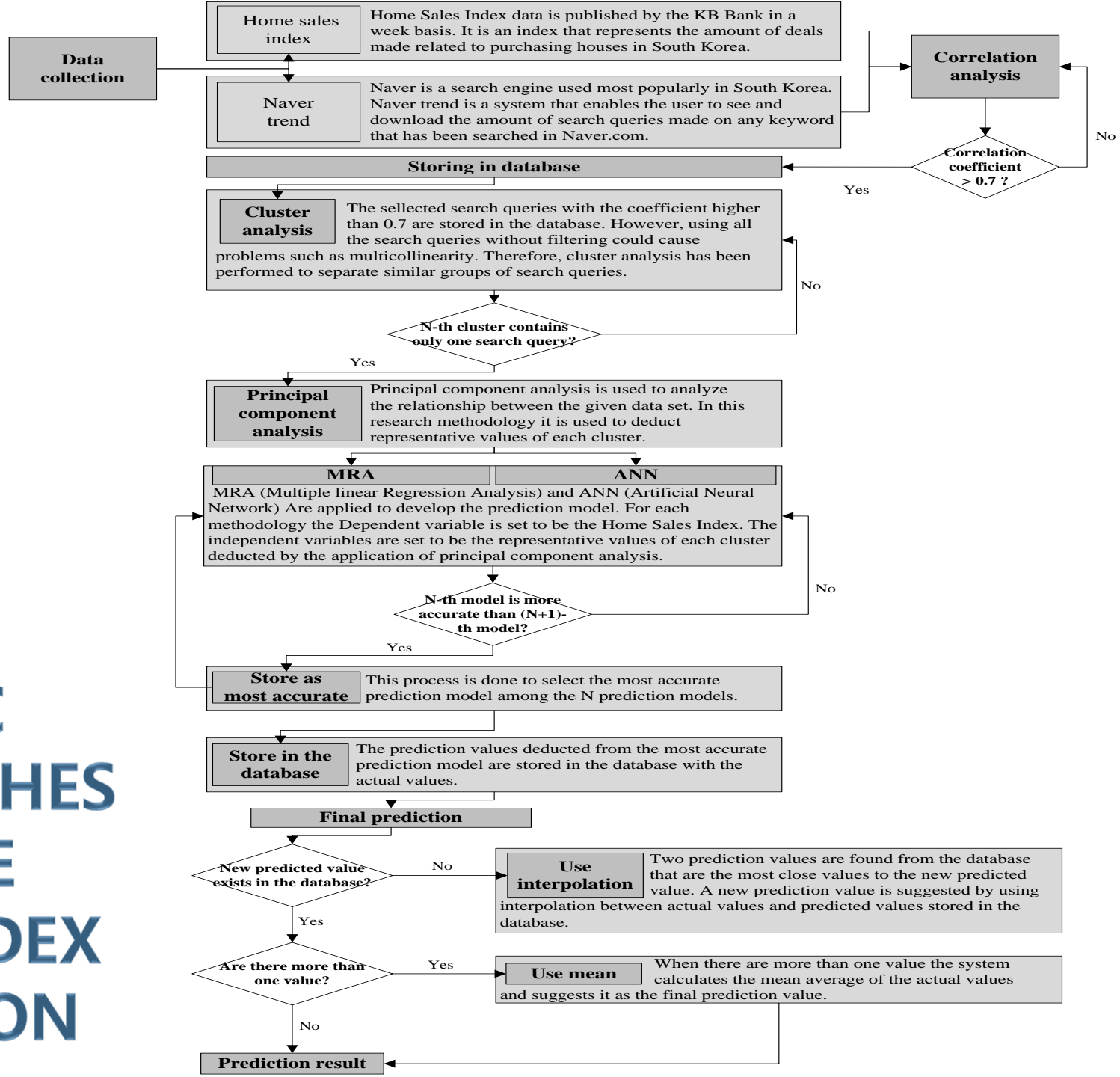
Productivity Bayesian



Unit Cost Bayesian

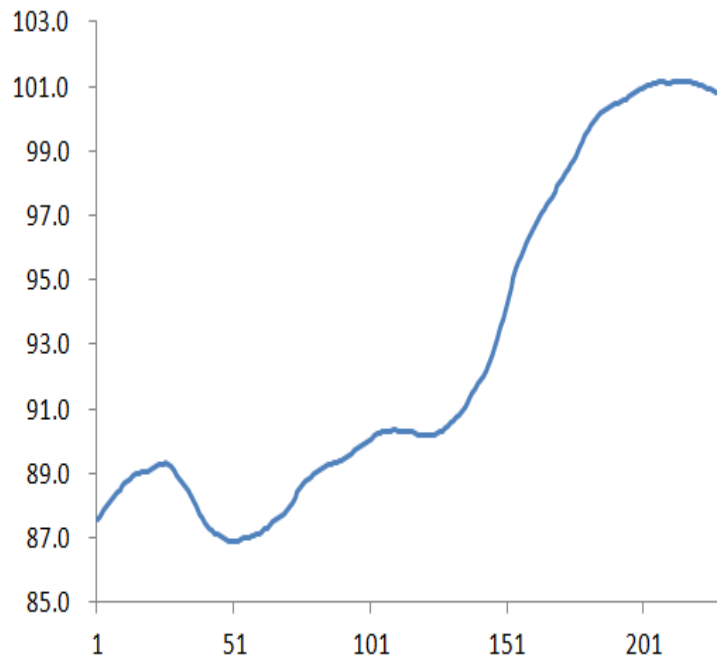


# BIG DATA ANALYTIC APPROACHES ON HOME SALES INDEX PREDICTION



# BIG DATA ANALYTIC APPROACHES ON HOME SALES INDEX PREDICTION

## HSI Data



## Naver Trend



# BIG DATA ANALYTIC APPROACHES ON HOME SALES INDEX PREDICTION

## Selected Search Queries

## Correlation Analysis

Search Queries				
2금융권	신용조회	오피스텔	국민주택채권	등기부등본
중개수수료	중고차	닥터아파트	은행대출	전세계약서
주공아파트	IBK기업은행	민원24	소형주택	주민등록증진위확인
광주집값	대구집값	동탄신도시	분양권프리미엄	신혼부부집
계약면적	등기사항전부증명서	부동산조회	매도용인감증명서	보증금반환
사다리차비용	소득증빙서류	수익형부동산		

search queries	상관계수	유의확률
1금융권	.192	.002
2금융권	-.765	.000
3금융권	.528	.000
IBK기업은행	.858	.000
계약면적	.837	.000
광주집값	.706	.000
국민주택채권	-.831	.000
닥터아파트	-.905	.000
대구집값	.778	.000
대전집값	.284	.000
동탄신도시	-.816	.000
등기부등본	-.825	.000
등기사항전부증명서	.849	.000
매도용인감증명서	.835	.000
민원24	.858	.000
보증금반환	.724	.000
부동산조회	.858	.000
부산집값	.571	.000
분양권	-.576	.000
분양권매매	.441	.000
분양권프리미엄	.719	.000
사다리차비용	.833	.000
서울집값	-.365	.000
소득증빙서류	.893	.000
소형주택	.775	.000
수익형부동산	.888	.000
수익형부동산	.888	.000
신용조회	-.898	.000
신혼부부아파트	-.245	.000
신혼부부전세	.413	.000



# BIG DATA ANALYTIC APPROACHES ON HOME SALES INDEX PREDICTION

## Cluster Analysis

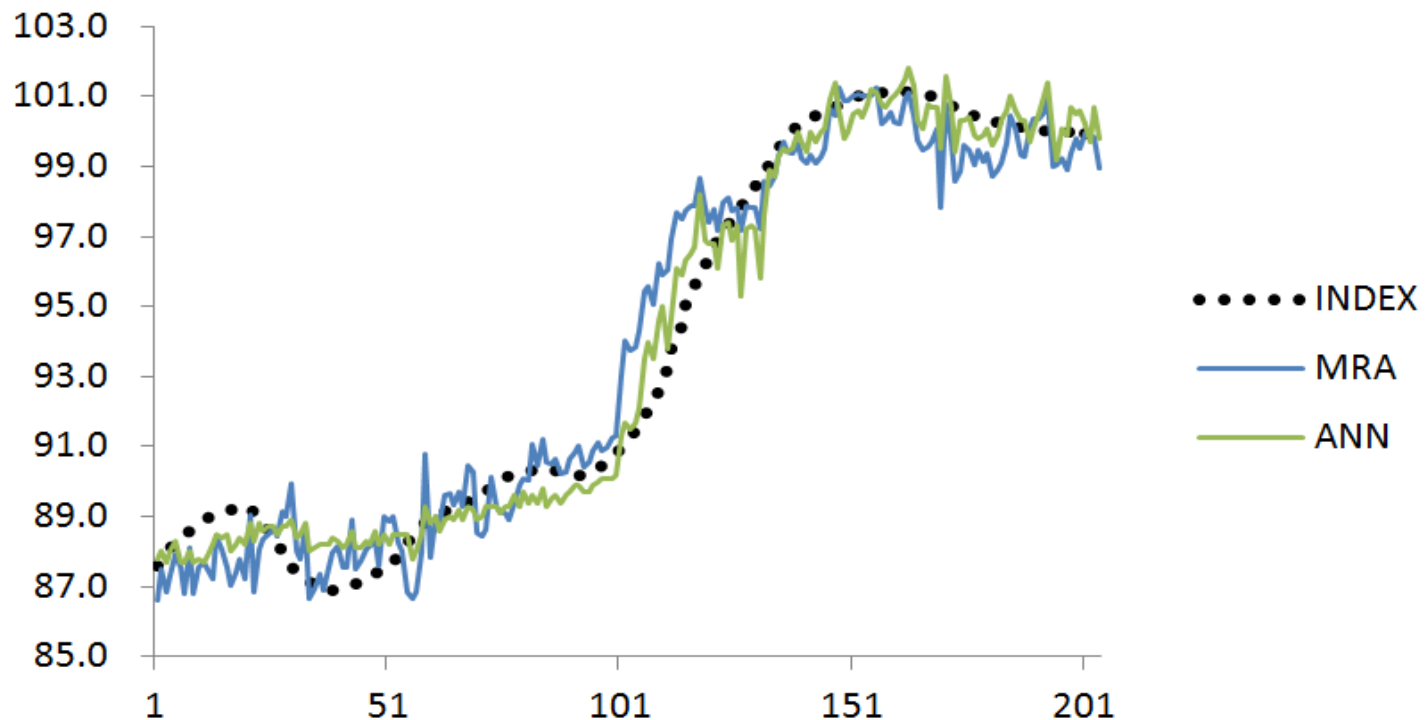
search queries	7	6	5	4	3	2
2금융권	1	1	1	1	1	1
신용조회	1	1	1	1	1	1
오피스텔	1	1	1	1	1	1
국민주택채권	5	5	4	4	1	1
등기부등본	5	5	4	4	1	1
중개수수료	5	5	4	4	1	1
중고차	5	5	4	4	1	1
닥터아파트	5	5	4	4	1	1
은행대출	5	5	4	4	1	1
전세계약서	5	5	4	4	1	1
주공아파트	5	5	4	4	1	1
IBK기업은행	2	2	2	2	2	2
민원24	2	2	2	2	2	2
소형주택	2	2	2	2	2	2
주민등록증진위 확인	2	2	2	2	2	2
<b>계약면적</b>	<b>3</b>	<b>3</b>	<b>3</b>	3	3	2
광주집값	4	4	2	2	2	2
대구집값	4	4	2	2	2	2
동탄신도시	4	4	2	2	2	2

## Principal Component Analysis

Num	index	1 cluster	2 cluster	4 cluster	5 cluster
1	87.6	1.04977	-1.18136	1.58591	-1.20381
2	88.1	0.717	-1.14796	1.32595	-1.26092
3	88.5	1.08544	-1.26073	1.67132	-0.92794
4	88.8	0.43747	-1.291	1.27459	-1.20788
5	89.0	1.21434	-1.24213	1.04835	-1.3545
6	89.2	1.67226	-1.05924	1.07667	-1.07133
7	89.3	1.3468	-1.16443	0.91849	-1.18628
8	89.1	1.65456	-1.07098	1.1521	-1.38944
9	88.6	1.33961	-1.12444	0.59324	-1.17678
10	88.0	0.47392	-0.98518	0.68877	-1.13361
11	87.4	1.1357	-1.07778	1.03508	-1.054
.....					
200	100.0	-1.11092	1.11627	-1.0566	1.45215
201	100.0	-1.23976	1.02335	-0.96124	0.68083
202	100.0	-1.3686	0.69108	-1.37334	0.56646
203	99.9	-1.16822	0.82175	-1.02179	1.32825

# BIG DATA ANALYTIC APPROACHES ON HOME SALES INDEX PREDICTION

Multiple Linear Regression & Artificial Neural Network



# BIG DATA ANALYTIC APPROACHES ON HOME SALES INDEX PREDICTION

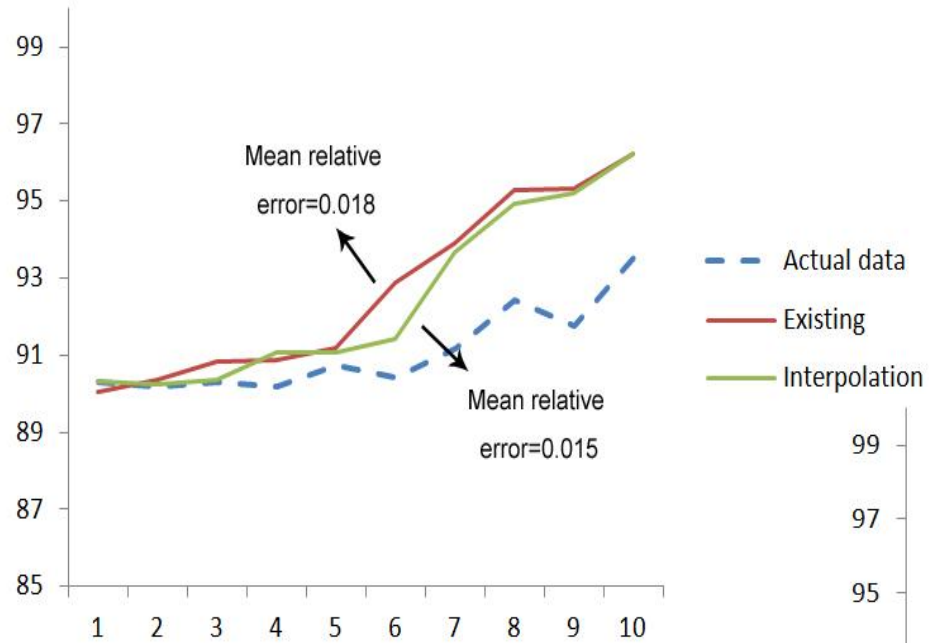
## Estimation using Interpolation

Actual data	MRA				ANN			
	Normal	Error	Interpolation	Error	Normal	Error	Interpolation	Error
90.32	90.04	0.0030	90.31	0.0001	89.40	0.0101	90.30	0.0002
<b>90.19</b>	<b>90.37</b>	<b>0.0020</b>	<b>90.26</b>	0.0008	89.40	0.0088	90.30	0.0012
90.30	90.83	0.0058	90.38	0.0010	90.30	0.0000	90.21	0.0010
90.18	90.88	0.0077	91.07	0.0098	89.40	0.0087	90.30	0.0013
90.74	91.19	0.0049	91.06	0.0035	90.30	0.0048	90.21	0.0058
90.40	92.88	0.0274	91.41	0.0112	90.40	0.0000	90.22	0.0020
91.18	93.92	0.0301	93.66	0.0272	90.40	0.0086	90.22	0.0105
92.42	95.28	0.0309	94.91	0.0269	92.40	0.0002	92.23	0.0021
91.75	95.33	0.0390	95.21	0.0376	100.20	0.0921	100.20	0.0921
93.52	96.24	0.0291	96.23	0.0290	93.50	0.0002	93.50	0.0002

Calculation:  $\frac{(90.33 - 90.25)}{(90.46 - 90.21)} \times (90.37 - 90.21) + 90.21 = 90.26$

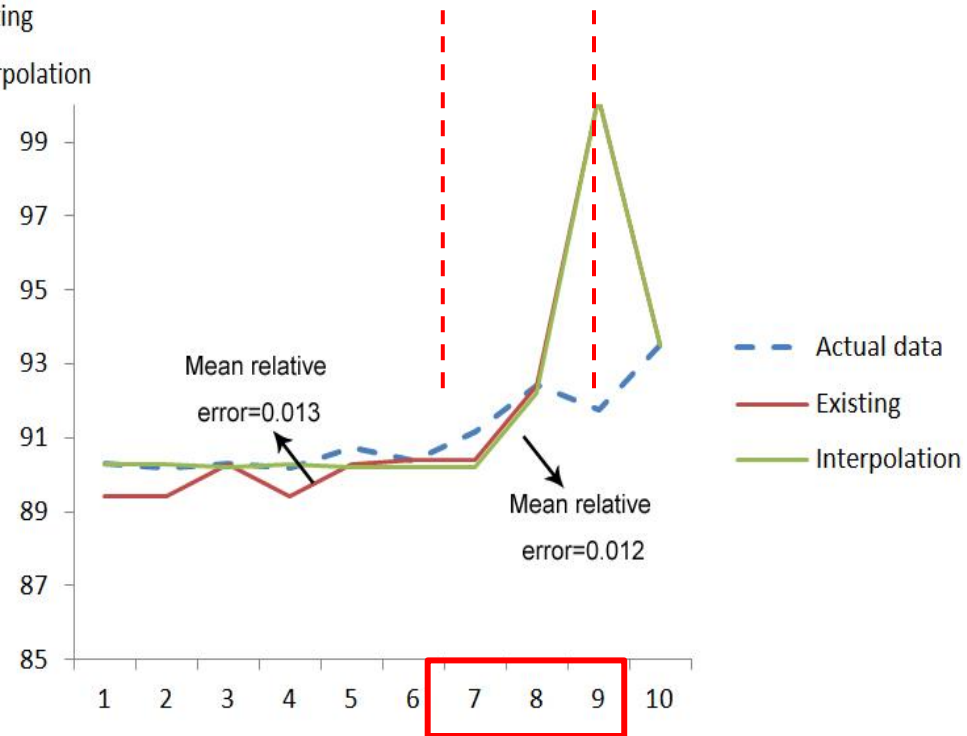
# BIG DATA ANALYTIC APPROACHES ON HOME SALES INDEX PREDICTION

## Result Comparison



## Multiple Linear Regression

## Artificial Neural Network





인하대학교  
INHA UNIVERSITY



빅데이터 전문인력 양성 협약: 인하대학교 – 한국DB진흥원  
소프트웨어융합공학 연계전공: 인하대학교 – 삼성전자