



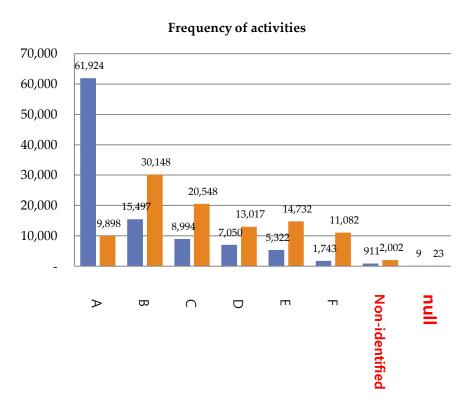
#### **Data quality for process improvement**

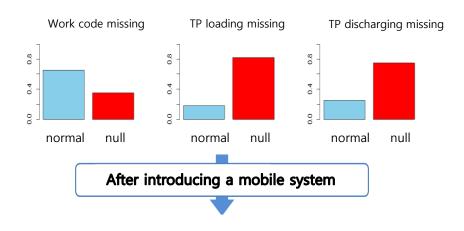
Most of the lecture slides are made by Prof. van der Aalst.

#### **Overview**



• Is technology the only issue for Industry 4.0?









#### Garbage in garbage out







#### **Data imperfection**

- Missing data
  - 자료누락
- Incorrect data
  - 잘못된 코드
- Imprecise data
  - 잘못된 측정 데이터
- Irrelevant data
  - 예측에 사용하기 힘든 자료

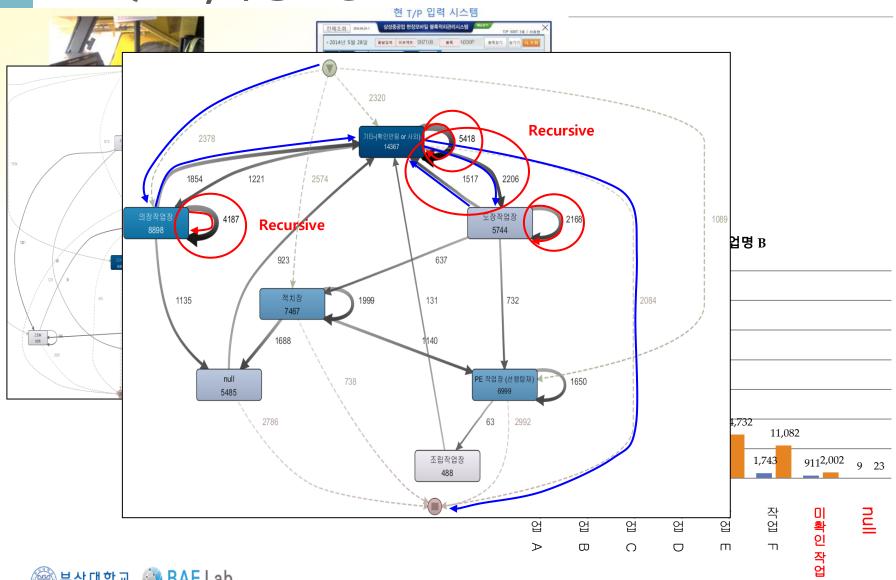
Table 1. Manifestation of quality issues in event log entities [6].

		Event log entities								
		Case	Event	Relationship	Case attrs.	Position	Activity name	Timestamp	Resource	Event attrs.
Event log	Missing data	l1	12	13	14	15	16	17	18	19
quality issues	Incorrect data	I10	l11	I12	I13	114	I15	116	117	I18
	Imprecise data			I19	120	I21	122	123	124	125
	Irrelevant data	126	127							





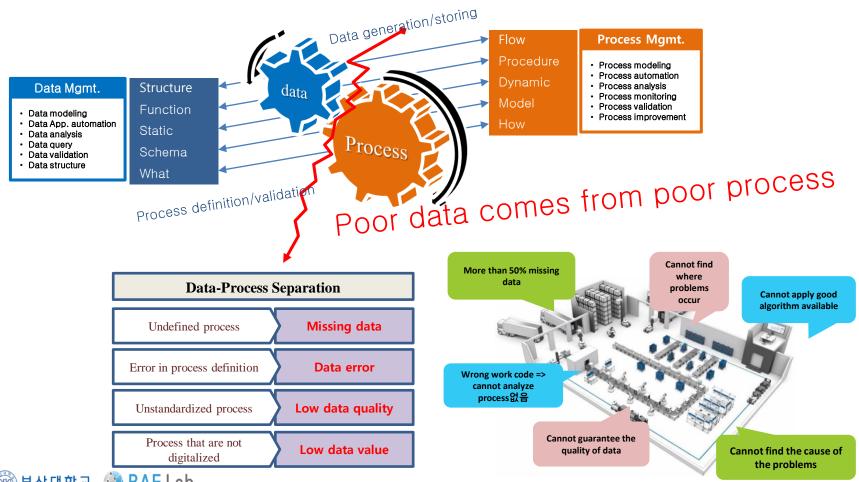
### Data Quality의 중요성







#### Big-data vs. process improvement

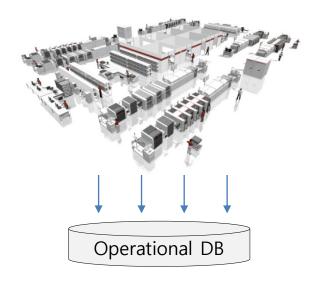


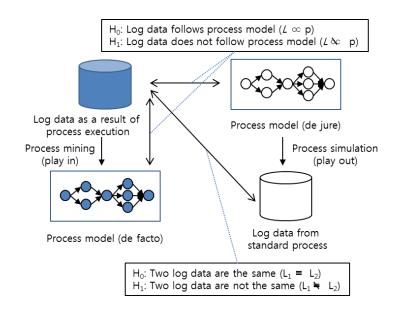




#### **Data-Process compliance**

• 만약 data-process separation을 없앨 수 있다면?









#### **Data Quality**

- **1.Consistency** logical relations
  - two similar IDs for two different employees
  - a non-existent entry in another table
- **2.Accuracy** the real state of things
  - All calculations based on such data show the true result.
- 3.Completeness all needed elements
  - lots of sensor data but there's no info about the exact sensor locations
- **4.Auditability** maintenance and control
  - data quality audits regularly or on demand will help to ensure a higher level of data adequacy
- **5.Orderliness** structure and format
  - the temperature in the oven has to be measured in Fahrenheit and can't be -14 °F.











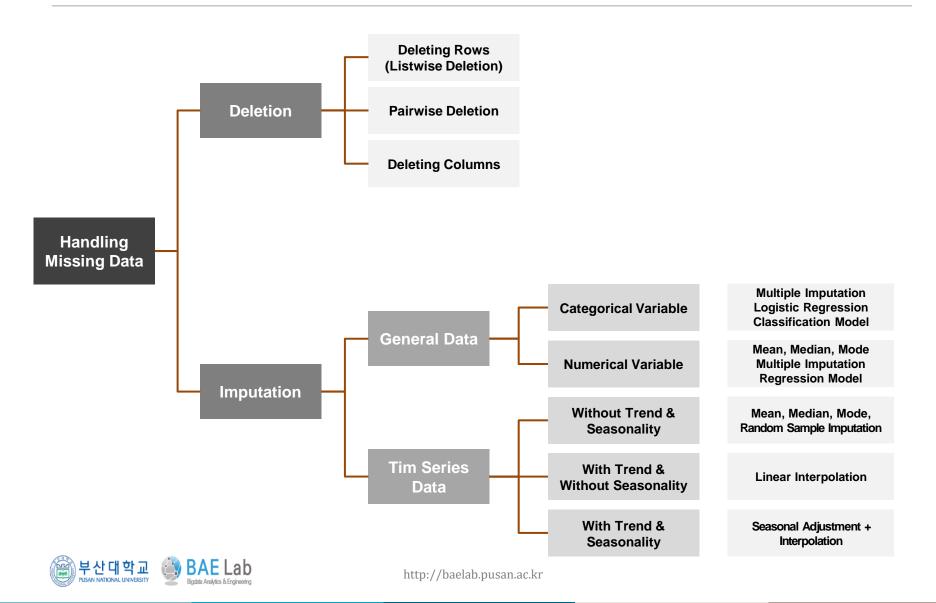
### **Data Imputation**

#### **Understanding about Missing Data**

- One of the most common problems we have faced in Data Cleaning/Exploratory Analysis is handling the missing values.
- The type of missing values was firstly classified by Rubin and the missing values have been classified into the following three kinds.
  - Missing at Random (MAR): Missing at random means that the propensity for a data point to be missing is not related to the missing data, but it is related to some of the observed data (자료 내의 다른 변수와 관련)
  - Missing Completely at Random (MCAR): The fact that a certain value is missing has nothing to do with
    its hypothetical value and with the values of other variables.
  - **Missing not at Random (MNAR)**: Two possible reasons are that the missing value depends on the hypothetical value (e.g. People with high salaries generally do not want to reveal their incomes in surveys) or missing value is dependent on some other variable's value (e.g. Let's assume that females generally don't want to reveal their ages! Here the missing value in age variable is impacted by gender variable) (Missing여부가 해당변수의 값에 의해서 결정)

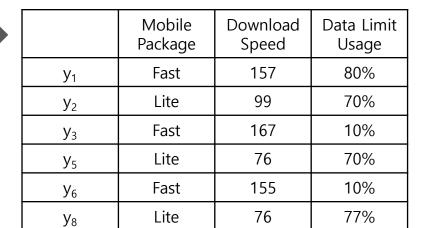






- Listwise deletion (complete-case analysis)
  - By far the most common approach to the missing data is to simply omit those cases with the missing data and analyse the remaining data.
  - This approach is known as the complete case (or available case) analysis or list-wise deletion.

	Mobile Package	Download Speed	Data Limit Usage
У1	Fast	157	80%
<b>y</b> <sub>2</sub>	Lite	99	70%
У3	Fast	167	10%
<b>y</b> <sub>4</sub>	Fast	NA	80%
<b>y</b> <sub>5</sub>	Lite	76	70%
У <sub>6</sub>	Fast	155	10%
<b>y</b> <sub>7</sub>	NA	NA	95%
y <sub>8</sub>	Lite	76	77%
<b>y</b> <sub>9</sub>	Fast	180	NA







- Pairwise deletion (available-case analysis)
  - Only the missing observations are ignored and analysis is done on variables present.
  - If there is missing data elsewhere in the data set, the existing values are used. Since a pairwise deletion
    uses all information observed, it preserves more information than the listwise deletion.

	Mobile Package	Download Speed	Data Limit Usage
У <sub>1</sub>	Fast	157	80%
<b>y</b> <sub>2</sub>	Lite	99	70%
<b>y</b> <sub>3</sub>	Fast	167	10%
y <sub>4</sub>	Fast	NA	80%
<b>y</b> <sub>5</sub>	Lite	76	70%
У <sub>6</sub>	Fast	155	10%
<b>y</b> <sub>7</sub>	NA	NA	95%
<b>y</b> <sub>8</sub>	Lite	76	77%
<b>y</b> <sub>9</sub>	Fast	180	NA



	Mobile Package	Download Speed	Data Limit Usage
y <sub>1</sub>	Fast	157	80%
<b>y</b> <sub>2</sub>	Lite	99	70%
<b>y</b> <sub>3</sub>	Fast	167	10%
y <sub>4</sub>	Fast		80%
<b>y</b> <sub>5</sub>	Lite	76	70%
y <sub>6</sub>	Fast	155	10%
<b>y</b> <sub>7</sub>			95%
y <sub>8</sub>	Lite	76	77%
<b>y</b> <sub>9</sub>	Fast	180	



- Column deletion (available-variable analysis)
  - If there are too many data missing for a variable it may be an option to delete the variable or the column from the dataset.
  - This should be the last option and need to check if model performance improves after deletion of variable.

	Mobile Package	Download Speed	Data Limit Usage
У <sub>1</sub>	Fast	NA	80%
<b>y</b> <sub>2</sub>	Lite	NA	70%
<b>y</b> <sub>3</sub>	Fast	167	10%
y <sub>4</sub>	Fast	NA	80%
<b>y</b> <sub>5</sub>	Lite	NA	70%
У <sub>6</sub>	Fast	155	10%
<b>y</b> <sub>7</sub>	Fast	NA	95%
<b>y</b> <sub>8</sub>	Lite	NA	77%
<b>y</b> <sub>9</sub>	Fast	180	80%

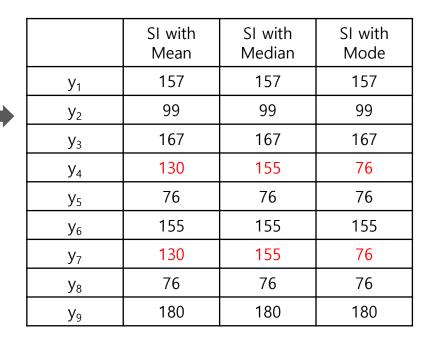


	Mobile Package	Data Limit Usage
y <sub>1</sub>	Fast	80%
<b>y</b> <sub>2</sub>	Lite	70%
<b>y</b> <sub>3</sub>	Fast	10%
y <sub>4</sub>	Fast	80%
<b>y</b> <sub>5</sub>	Lite	70%
y <sub>6</sub>	Fast	10%
<b>y</b> <sub>7</sub>	Fast	95%
y <sub>8</sub>	Lite	77%
<b>y</b> <sub>9</sub>	Fast	80%



- Simple Imputation(Mean, Median and Mode)
  - In this simple imputation technique goal is to replace missing data with statistical estimates of the missing values. Mean, Median or Mode can be used as imputation value.
  - Ex) Mean = 130, Median = 155, Mode = 200

	Mobile Package	Download Speed	Data Limit Usage
<b>y</b> <sub>1</sub>	Fast	157	80%
<b>y</b> <sub>2</sub>	Lite	99	70%
<b>y</b> <sub>3</sub>	Fast	167	10%
<b>y</b> <sub>4</sub>	Fast	NA	80%
<b>y</b> <sub>5</sub>	Lite	76	70%
У <sub>6</sub>	Fast	155	10%
<b>y</b> <sub>7</sub>	Fast	NA	95%
<b>y</b> <sub>8</sub>	Lite	76	77%
<b>y</b> <sub>9</sub>	Fast	180	80%

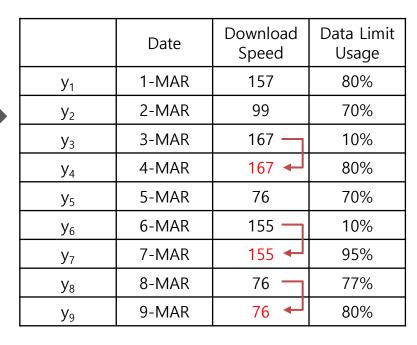






- Time Series Specific Method (LOCF, NOCB and Linear Interpolation)
  - Last Observation Carried Forward(LOCF)
    - If data is time-series data, one of the most widely used imputation methods.
    - Whenever a value is missing, it is replaced with the last observed value.

	Date	Download Speed	Data Limit Usage
y <sub>1</sub>	1-MAR	157	80%
<b>y</b> <sub>2</sub>	2-MAR	99	70%
<b>y</b> <sub>3</sub>	3-MAR	167	10%
y <sub>4</sub>	4-MAR	NA	80%
<b>y</b> <sub>5</sub>	5-MAR	76	70%
<b>y</b> <sub>6</sub>	6-MAR	155	10%
<b>y</b> <sub>7</sub>	7-MAR	NA	95%
<b>y</b> <sub>8</sub>	8-MAR	76	77%
<b>y</b> <sub>9</sub>	9-MAR	NA	80%

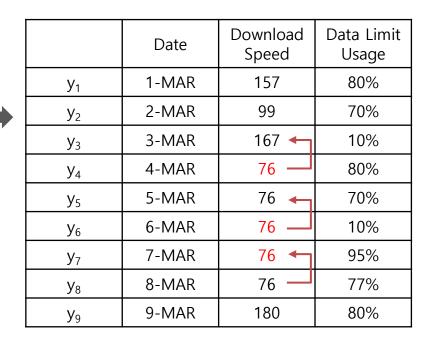






- Time Series Specific Method (LOCF, NOCB and Linear Interpolation)
  - Next Observation Carried Backward(NOCB)
    - A similar approach like LOCF which works in the opposite direction by taking the first observation after the missing value and **c**arrying it backward

	Date	Download Speed	Data Limit Usage
y <sub>1</sub>	1-MAR	157	80%
<b>y</b> <sub>2</sub>	2-MAR	99	70%
<b>y</b> <sub>3</sub>	3-MAR	167	10%
<b>y</b> <sub>4</sub>	4-MAR	NA	80%
<b>y</b> <sub>5</sub>	5-MAR	76	70%
<b>y</b> <sub>6</sub>	6-MAR	NA	10%
<b>y</b> <sub>7</sub>	7-MAR	NA	95%
<b>y</b> <sub>8</sub>	8-MAR	76	77%
<b>y</b> <sub>9</sub>	9-MAR	180	80%







- Time Series Specific Method (LOCF, NOCB and Linear Interpolation)
  - Linear Interpolation
    - Interpolation is a mathematical method that adjusts a function to data and uses this function to extrapolate the missing data.
    - The simplest type of interpolation is the linear interpolation, that makes a mean between the values before the missing data and the value after.

	Date	Download Speed	Data Limit Usage
y <sub>1</sub>	1-MAR	157	80%
<b>y</b> <sub>2</sub>	2-MAR	99	70%
<b>y</b> <sub>3</sub>	3-MAR	167	10%
<b>y</b> <sub>4</sub>	4-MAR	NA	80%
<b>y</b> <sub>5</sub>	5-MAR	76	70%
У <sub>6</sub>	6-MAR	NA	10%
<b>y</b> <sub>7</sub>	7-MAR	150	95%
<b>y</b> <sub>8</sub>	8-MAR	76	77%
<b>y</b> <sub>9</sub>	9-MAR	180	80%

	Date	Download Speed		Data Limit Usage		
y <sub>1</sub>	1-MAR		157			80%
<b>y</b> <sub>2</sub>	2-MAR		99			70%
<b>y</b> <sub>3</sub>	3-MAR		167			10%
<b>y</b> <sub>4</sub>	4-MAR	Г	121.5			80%
<b>y</b> <sub>5</sub>	5-MAR	76			70%	
y <sub>6</sub>	6-MAR		113			10%
<b>y</b> <sub>7</sub>	7-MAR		150			95%
y <sub>8</sub>	8-MAR	76				77%
<b>y</b> <sub>9</sub>	9-MAR		180			80%
·						

(167+76)/2 = 121.5

(76+150)/2 = 113





- Time Series Specific Method (LOCF, NOCB and Linear Interpolation)
  - Linear Interpolation
    - Interpolation is a mathematical method that adjusts a function to data and uses this function to extrapolate the missing data.
    - The simplest type of interpolation is the linear interpolation, that makes a mean between the values before the missing data and the value after.

	Date	Download Speed	Data Limit Usage
У <sub>1</sub>	1-MAR	157	80%
<b>y</b> <sub>2</sub>	2-MAR	99	70%
<b>y</b> <sub>3</sub>	3-MAR	167	10%
У <sub>4</sub>	4-MAR	NA	80%
<b>y</b> <sub>5</sub>	5-MAR	76	70%
<b>y</b> <sub>6</sub>	6-MAR	NA	10%
<b>y</b> <sub>7</sub>	7-MAR	150	95%
<b>y</b> <sub>8</sub>	8-MAR	76	77%
<b>y</b> <sub>9</sub>	9-MAR	180	80%

y1     1-MAR     157     80%       y2     2-MAR     99     70%       y3     3-MAR     167     10%       y4     4-MAR     121.5     80%       y5     5-MAR     76     70%		Date	D	ownloa Speed	d	D	ata Limit Usage
y <sub>3</sub> 3-MAR     167     10%       y <sub>4</sub> 4-MAR     121.5     80%       y <sub>5</sub> 5-MAR     76     70%	y <sub>1</sub>	1-MAR		157			80%
y <sub>4</sub> 4-MAR     121.5     80%       y <sub>5</sub> 5-MAR     76     70%	<b>y</b> <sub>2</sub>	2-MAR		99			70%
y <sub>5</sub> 5-MAR 76 70%	<b>y</b> <sub>3</sub>	3-MAR		167			10%
73	<b>y</b> <sub>4</sub>	4-MAR	Г	121.5			80%
	<b>y</b> <sub>5</sub>	5-MAR		76			70%
y <sub>6</sub> 6-MAR 113 10%	y <sub>6</sub>	6-MAR		113			10%
y <sub>7</sub> 7-MAR 150 95%	<b>y</b> <sub>7</sub>	7-MAR		150			95%
y <sub>8</sub> 8-MAR 76 77%	y <sub>8</sub>	8-MAR		76			77%
y <sub>9</sub> 9-MAR 180 80%	<b>y</b> <sub>9</sub>	9-MAR		180			80%

(167+76)/2 = 121.5

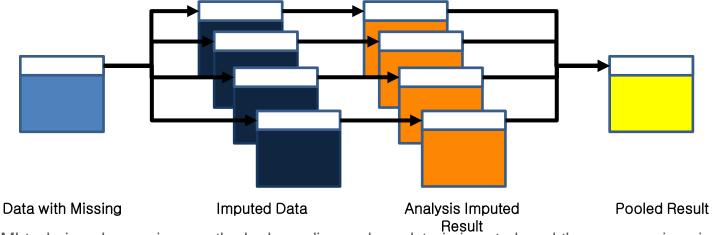
(76+150)/2 = 113





#### Multiple Imputation

- Multiple imputation (MI) is a statistical technique for dealing with missing data.
- The Multiple imputation includes the following 3 components.
  - Generate missing value: To use the distribution of the observed data to estimate a set of plausible values for the missing data.
  - Imputation and Analysis: Missing values are replaced by several estimated values, and are analyzed separately equally to obtain parameter estimates.
  - Pooling: The estimates are combined to obtain a set of parameter estimates.



 MI technique has various methods depending on how data is imputed, and there are various imputation method such as MICE (Multiple Imputation by Chain Equation), Random Forest Imputation, KNN Imputation, Expectation-Maximization Imputation.





Understanding for MICE – Single Iteration

Age	Income	Gender
33	NA	F
18	12,000	NA
NA	13,542	М



Understanding for MICE – Single Iteration

_		:1514110111	
	Age	Income	Gender
	33	NA	F
	18	12,000	NA
	NA	13,542	М

Simple
Imputation
usina mean

<u>neranon</u>		
Age	Income	Gender
33	12.771	F
18	12,000	F
25.5	13,542	М



Understanding for MICE – Single Iteration

	HOLDING	
Age	Income	Gender
33	NA	F
18	12,000	NA
NA	13,542	М

Simple Imputation using mean

HELAHOLL		
Age	Income	Gender
33	12.771	F
18	12,000	F
25.5	13,542	М

Age back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М



Understanding for MICE – Single Iteration

	1.3171111111	
Age	Income	Gender
33	NA	F
18	12,000	NA
NA	13,542	М

Simple Imputation using mean

neranon		
Age	Income	Gender
33	12.771	F
18	12,000	F
25.5	13,542	М

Age back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression Age ~ Income + Gender

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Predict Age





Understanding for MICE – Single Iteration

Age	Income	Gender
33	NA	F
18	12,000	NA
NA	13,542	М

Simple Imputation using mean

<u>lleralion</u>		
Age	Income	Gender
33	12.771	F
18	12,000	F
25.5	13,542	М

	<b></b>
Age	back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression Age ~ Income + Gender

Age	Income	Gender
33	12.771	F
18	12,000	F
35.3	13,542	М



Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М





Understanding for MICE – Single Iteration

-		$\sigma \rightarrow \sigma \rightarrow$
Age	Income	Gender
33	NA	F
18	12,000	NA
NA	13,542	М

Simple Imputation using mean

Heralion		
Age	Income	Gender
33	12.771	F
18	12,000	F
25.5	13,542	М

Age back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression Age ~ Income + Gender

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Income back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
35.3	13,542	М



Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М



Understanding for MICE - Single Ite Gender Age Income 33 NA F 18 12,000 NA 13,542 NA М

Simple Imputation

using mean

<u>Iteration</u>		
Age	Income	Gender
33	12.771	F
18	12,000	F
25.5	13,542	М

Age back to

Predict Age

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression Age ~ Income + Gender

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Income back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
35.3	13,542	М

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

6 .
Regression
Income ~ Age + Gender

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М





•	<ul> <li>Understanding for MIC</li> </ul>				
	Age	Income	Gender		
	33	NA	F		
	18	12,000	NA		
	NA	13,542	М		

CE – Single Iteration				
	Age	Income	Gender	
	33	12.771	F	
Simple	18	12,000	F	
Imputation using mean	25.5	13,542	М	

Age back to

Predict

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression Age ~ Income + Gender

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Income back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
35.3	13,542	М

Age	

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression	
Income ~ Age	+ Gender

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Predict Income

Age	Income	Gender
33	13.103	F
18	12,000	F
35.3	13,542	М





Understanding for MICE - Single It Gender Age Income 33 F NA 18 12,000 NA 13,542 NA Μ

Simple Imputation using mean

<u>lteration</u>		
Age	Income	Gender
33	12.771	F
18	12,000	F
25.5	13,542	М

Age back to

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression Age ~ Income + Gender

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Income back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
35.3	13,542	М

Predict Age

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Regression Income ~ Age + Gender

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Predict Income

Age	Income	Gender
33	13.103	F
18	12,000	F
35.3	13,542	М

...the same for Gender Gender back to NA Regression Predict Gender

	Age	Income	Gender
	33	13.103	F
١.	18	12,000	M
	35.3	13,542	М





• Understanding for MICE – Multiple Iterations

Age	Income		Gender	
33	12.771		F	
18	12,000 13,542		F	
NA			М	
		_	ession ~ Income + Gen	der

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Income back to NA

Age	Income	Gender
33	12.771	F
18	12,000	F
35.3	13,542	М

Predict Age

Age	Income	Gender
33	12.771	F
18	12,000	F
NA	13,542	М

Age	Income	Gender
33	NA	F
18	12,000	F
35.3	13,542	М

Predict Income

Age	Income	Gender
33	13.103	F
18	12,000	F
35.3	13,542	М

---the same for Gender Gender back to N Regression Predict Gender

	Age	Income	Gender	
	33	13.103	F	
IA	18	12,000	М	
	35.3	13,542	М	





Regression

Income ~ Age + Gender

#### Other Multiple Imputation Techniques

- Random Forest Imputation
  - Random Forest is a nonparametric replacement method that can be applied to various types of variables that work well on both randomly missing and non-randomly missing data.
  - One caveat is that random forests work best on large datasets, and using random forests on small datasets risks overconsensus.
- K Nearest Neighbor Imputation
  - k-NN replaces missing attribute values based on the nearest K neighbor and neighbors are determined by distance measurements.
  - When K neighbors are determined, the missing value is replaced by taking the mean / medium or mode of the known attribute value of the missing attribute.
- Expectation-Maximization Imputation
  - EM (Expectation-Maximization) is a type of maximum likelihood method that can be used to create a new data set, and all missing values are replaced with values estimated by the maximum likelihood method
  - The EM algorithm consists of 3 phase
    - 1) Expected phase: Estimated various parameters(e.g. variance, covariance and mean) using list-specific deletions.
    - 2) Imputation phase: Use these estimates to create a regression equation that predicts missing data.
    - 3) Maximizing phase: Uses these equations to fill in the missing data.
  - Then repeat the expected step with the new parameters. The new regression equation is determined to "fill" the missing data. Expectation and maximization steps are repeated until the system stabilizes.





Likelihood based Multiple Imputation by Event chain for Repairing Event Log

	Case	Activity	Resource	Part Desc.
$e_1$	Case <sub>1</sub>	Turning & Milling	Machine 4	Cable Head
$e_2$	Case <sub>1</sub>	Turning & Milling	Machine 4	Cable Head
$e_3$	Case <sub>1</sub>	Turning & Milling	Machine 4	Cable Head
$e_4$	Case <sub>1</sub>	Turning & Milling	Machine 4	Cable Head
$e_5$	Case <sub>1</sub>	Turning & Milling Q.C	Quality Check 1	Cable Head
$e_6$	Case <sub>1</sub>	Laser Marking	Machine 7	Cable Head
$e_7$	Case <sub>1</sub>	Lapping	Machine 1	Cable Head
$e_8$	Case <sub>1</sub>	Lapping	Machine 1	Cable Head
$e_g$	Case <sub>1</sub>	Lapping	Machine 1	Cable Head
$e_{10}$	Case <sub>1</sub>	Lapping	Machine 1	Cable Head
$e_{11}$	Case <sub>1</sub>	Round Grinding	Machine 3	Cable Head
$e_{12}$	Case <sub>1</sub>	Round Grinding	Machine 3	Cable Head
$e_{13}$	Case <sub>1</sub>	Final Inspection Q.C.	Quality Check 1	Cable Head
e <sub>14</sub>	Case <sub>1</sub>	Final Inspection Q.C.	Quality Check 1	Cable Head
e <sub>15</sub>	Case <sub>1</sub>	Final Inspection Q.C.	Quality Check 1	Cable Head
$e_{16}$	Case <sub>1</sub>	Packing	Packing	Cable Head
e <sub>17</sub>	Case <sub>2</sub>	Turning & Milling	Machine 9	Spur Gear
$e_{18}$	Case <sub>2</sub>	Turning Q.C.	Quality Check 1	Spur Gear
e <sub>19</sub>	Case <sub>2</sub>	Turning & Milling	Machine 9	Spur Gear
$e_{20}$	Case <sub>2</sub>	Turning & Milling	Machine 9	Spur Gear
	•••			

Start Time	End Time
2012-01-29 23:24	2012-01-30 05:43
2012-01-30 05:44	2012-01-30 06:42
2012-01-30 06:59	2012-01-30 07:21
2012-01-30 07:21	2012-01-30 10:58
2012-01-31 13:20	2012-01-31 14:50
2012-02-01 08:18	2012-02-01 08:27
2012-02-14 00:00	2012-02-14 01:15
2012-02-14 00:00	2012-02-14 01:15
2012-02-14 09:05	2012-02-14 10:20
2012-02-14 09:05	2012-02-14 09:38
2012-02-14 09:13	2012-02-14 13:37
2012-02-14 13:37	2012-02-14 15:27
2012-02-16 06:59	2012-02-16 07:59
2012-02-16 12:11	2012-02-16 16:12
2012-02-16 12:43	2012-02-16 13:58
2012-02-17 00:00	2012-02-17 01:00
2012-01-17 07:01	2012-01-17 11:05
2012-01-17 11:06	2012-01-17 11:15
2012-01-17 19:24	2012-01-17 20:01
2012-01-17 20:01	2012-01-17 23:43
	•••

Factory Operation Event Log Example





Likelihood based Multiple Imputation by Event chain for Repairing Event Log

data included in a case (row) are dependent.

	Mobile Package	Download Speed	Data Limit Usage
y <sub>1</sub>	NA	157	80%
<b>y</b> <sub>2</sub>	Lite	99	NA
<b>y</b> <sub>3</sub>	Fast	167	10%
У <sub>4</sub>	Fast	NA	80%

General Data Set with missing

#### VS

	Case	Activity	Resource	Part Desc.
$e_1$	Case <sub>1</sub>	NA	Machine 4	Cable Head
$e_2$	Case <sub>1</sub>	Turning & Milling	NA	Cable Head
$e_3$	Case <sub>1</sub>	Turning & Milling	Machine 4	NA
$e_4$	Case <sub>1</sub>	Turning & Milling	Machine 4	Cable Head
$e_5$	Case <sub>1</sub>	NA	NA	Cable Head
•••				

Start Time	End Time			
2012-01-29 23:24	2012-01-30 05:43			
2012-01-30 05:44	2012-01-30 06:42			
2012-01-30 06:59	2012-01-30 07:21			
2012-01-30 07:21	2012-01-30 10:58			
2012-01-31 13:20	2012-01-31 14:50			

Observations included in a case are dependent.

Event Log Structure with missing





Likelihood based Multiple Imputation by Event chain for Repairing Event Log

	Case	Activity				
$e_1$	Case <sub>1</sub>	Turning & Milling				
$e_2$	Case <sub>1</sub>	Turning & Milling				
$e_3$	Case <sub>1</sub>	Turning & Milling				
$e_4$	Case <sub>1</sub>	Turning & Milling				
$e_5$	Case <sub>1</sub>	Turning & Milling Q.C				
$e_6$	Case <sub>1</sub>	Laser Marking				
e <sub>7</sub>	Case <sub>1</sub>	Lapping				
$e_8$	Case <sub>1</sub>	Lapping				
$e_g$	Case <sub>1</sub>	Lapping				
$e_{10}$	Case <sub>1</sub>	Lapping				
e <sub>11</sub>	Case <sub>1</sub>	Round Grinding				
e <sub>12</sub>	Case <sub>1</sub>	Round Grinding				
e <sub>13</sub>	Case <sub>1</sub>	Final Inspection Q.C.				
e <sub>14</sub>	Case <sub>1</sub>	Final Inspection Q.C.				
e <sub>15</sub>	Case <sub>1</sub>	Final Inspection Q.C.				
e <sub>16</sub>	Case <sub>1</sub>	Packing				



Fitting Target
Distribution Using
Event Chain

Event

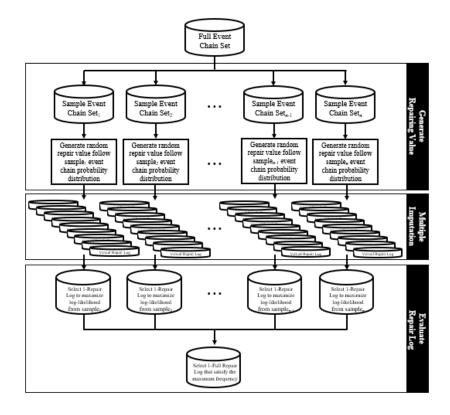
**Event Chain** 







- We developed the concept of an event chain that can reflect sequential information contained in one case for dealing with missing events.
- By converting the events included in the event into an event chain, and replacing the event value that can
  approximate the distribution of the event chain instead of missing, we developed a restoration method
  that can restore the restored event log to the original log.

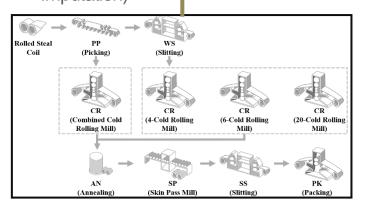






Likelihood based Multiple Imputation using Event chain for Repairing Event Log

 Case Study: Korea Steel Company Event Log Data with MIEC(Multiple Imputation by Chained Equations, Expectation-Maximizing Imputation, Random Forest Imputation, K Nearest Neighbor Imputation)



COIL_NO	PRC_CD	PRC_CD1	THK	WDT	WGT	SDT	EDT	PLNPRC_CD	ORD_NO	DRTCOI_NO	Machine
15KM12191111A	PP21	PP	5.99	1132	22900	01/04/2016 16:16:00	01/04/2016 16:38:00	PP21092	KD16090327	15KM1219111	PP2
15KM12191111A	CR21	CR	3.5	1132	22900	01/05/2016 10:10:00	01/05/2016 10:50:00	PP21092	KD16090327	15KM1219111	CR2
15KM12191111A	RC11	RC	3.5	1132	22900	01/05/2016 10:45:00	01/05/2016 11:10:00	PP21092	KD16090327	15KM1219111	RC1
15KM12191111A	WS31	WS	3.5	186	1717	01/07/2016 16:40:00	01/07/2016 17:50:00	PP21092	KD16090327	15KM1219111	WS3
15KM12191111A	PK41	PK	3.5	186	1736	01/07/2016 17:55:00	01/07/2016 17:55:00	PP21092	KD16090327	15KM1219111	PK4
15KM12191111A	PR11	PR	3.5	186	1736	01/07/2016 17:30:00	01/07/2016 17:30:00	PP21092	KD16090327	15KM1219111	PR1
15KM12191112A	PP21	PP	5.99	1132	22900	01/04/2016 16:16:00	01/04/2016 16:38:00	PP21092	KD16090327	15KM1219111	PP2
15KM12191112A	CR21	CR	3.5	1132	22900	01/05/2016 10:10:00	01/05/2016 10:50:00	PP21092	KD16090327	15KM1219111	CR2
15KM12191112A	RC11	RC	3.5	1132	22900	01/05/2016 10:45:00	01/05/2016 11:10:00	PP21092	KD16090327	15KM1219111	RC1
15KM12191112A	WS31	WS	3.5	186	1717	01/07/2016 16:40:00	01/07/2016 17:50:00	PP21092	KD16090327	15KM1219111	WS3
15KM12191112A	PK41	PK	3.5	186	1736	01/07/2016 17:55:00	01/07/2016 17:55:00	PP21092	KD16090327	15KM1219111	PK4
15KM12191112A	PR11	PR	3.5	186	1736	01/07/2016 17:30:00	01/07/2016 17:30:00	PP21092	KD16090327	15KM1219111	PR1
15KM12191113A	PP21	PP	5.99	1132	22900	01/04/2016 16:16:00	01/04/2016 16:38:00	PP21092	KD16090327	15KM1219111	PP2
15KM12191113A	CR21	CR	3.5	1132	22900	01/05/2016 10:10:00	01/05/2016 10:50:00	PP21092	KD16090327	15KM1219111	CR2
15KM12191113A	RC11	RC	3.5	1132	22900	01/05/2016 10:45:00	01/05/2016 11:10:00	PP21092	KD16090327	15KM1219111	RC1
15KM12191113A	WS31	WS	3.5	186	1717	01/07/2016 16:40:00	01/07/2016 17:50:00	PP21092	KD16090327	15KM1219111	WS3
15KM12191113A	PK41	PK	3.5	186	1743	01/07/2016 17:54:00	01/07/2016 17:54:00	PP21092	KD16090327	15KM1219111	PK4
15KM12191113A	PR11	PR	3.5	186	1743	01/07/2016 17:30:00	01/07/2016 17:30:00	PP21092	KD16090327	15KM1219111	PR1
15KM12191113B	PP21	PP	5.99	1132	22900	01/04/2016 16:16:00	01/04/2016 16:38:00	PP21092	KD16080281	15KM1219111	PP2

Missing Rate	MICE	ЕМВІ	RFI	KNNI	MIEC
5%	48.9%	51.3%	60.1%	44.4%	93.2%
10%	44.3%	45.5%	59.8%	42.1%	91.0%
15%	34.7%	36.9%	50.3%	40.3%	88.7%
20%	28.5%	30.1%	47.7%	32.2%	80.4%

Performance Event Imputation method using MIEC (Korea Steel Company Event log)



