



# **POORNIMA COLLEGE OF ENGINEERING JAIPUR**

A Presentation On

## **ANALYSIS & RECOMMENDATIONS OF DATA CLEANING PROCESS FOR VARIOUS TYPES OF DATA**

Group: 3CSC13 3<sup>rd</sup> SEM NSP

Submitted To :-

Mrs. SONAM GOUR Ma'am,  
**Mrs. Neelam Chaplot Ma'am,**  
NSP GUIDE,  
Professor,  
Department of CSE

Submitted By :-

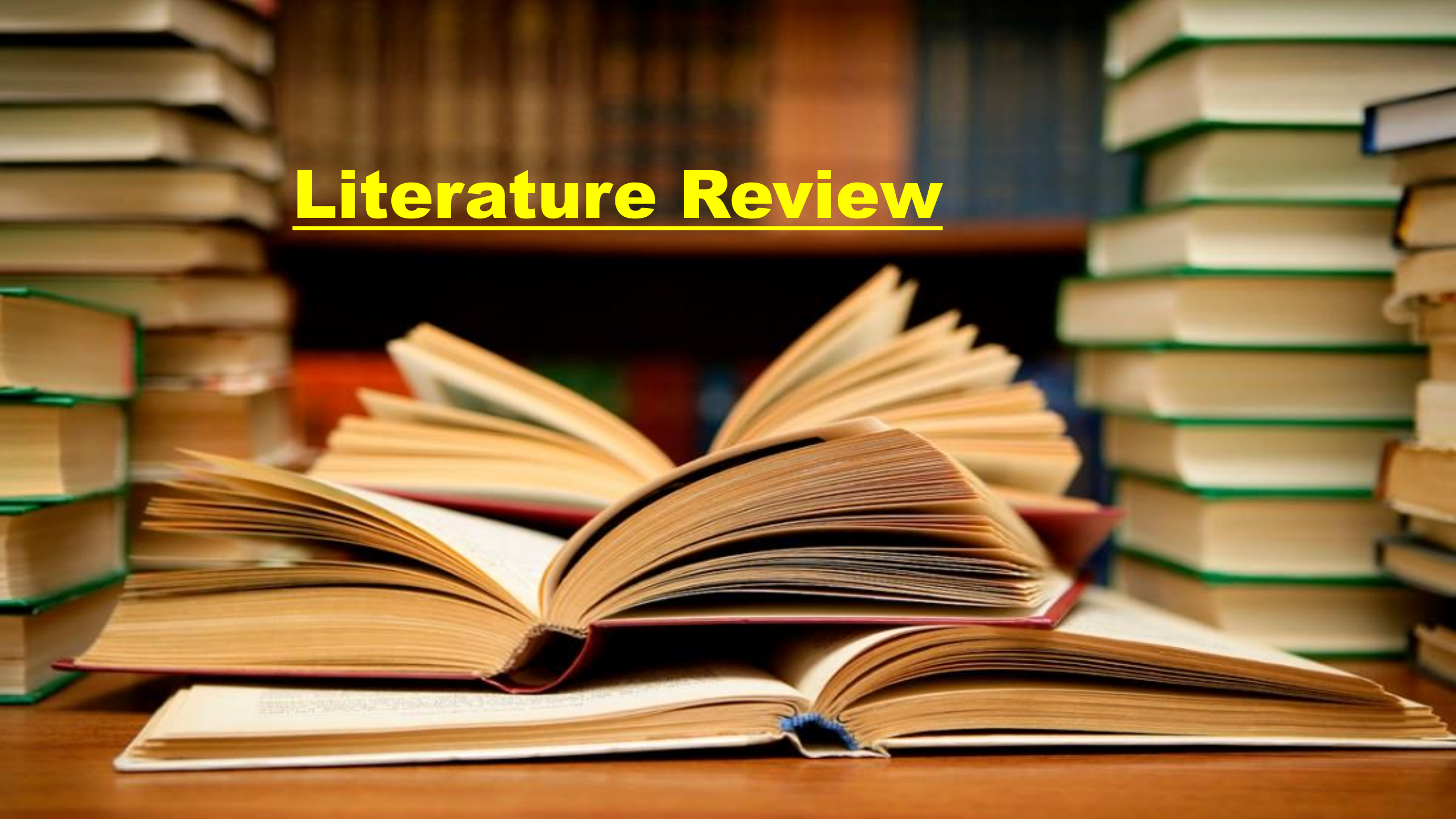
**Sejal Jain PCE20CS171, 20EPCCS171**  
**SANSKAR SHARMA PCE20CS166, 20EPCCS166**  
**SACHIN YADAV PCE20CS161, 20EPCCS161**  
**RAHUL KHANDELWAL PCE20CS152, 20EPCCS152**  
**YASH TRIPATHI PCE20CS200, 20EPCCS200**

# OUTLINES

- ABSTRACT
- SUMMARY
- APPROACH FOR DATA CLEANING
- DATA CLEANING ITERATIVE MODEL
- PROBLEM STATEMENT AND OBJECTIVE
- COMPARISON OF DATA CLEANING TOOLS
- RESEARCH DONE IN PAST
- DIFFERENT APPROACHES AND COMPARITIVE ANALYSIS
- FUTURE SCOPE
- REFERENCES



# Literature Review





# **ABSTRACT**

- The process of identifying and removing the errors
- Quality and consistency becomes significant
- Problem of data cleansing and the identification of potential errors for incorrect or inconsistent data
- False conclusion and misdirect investment

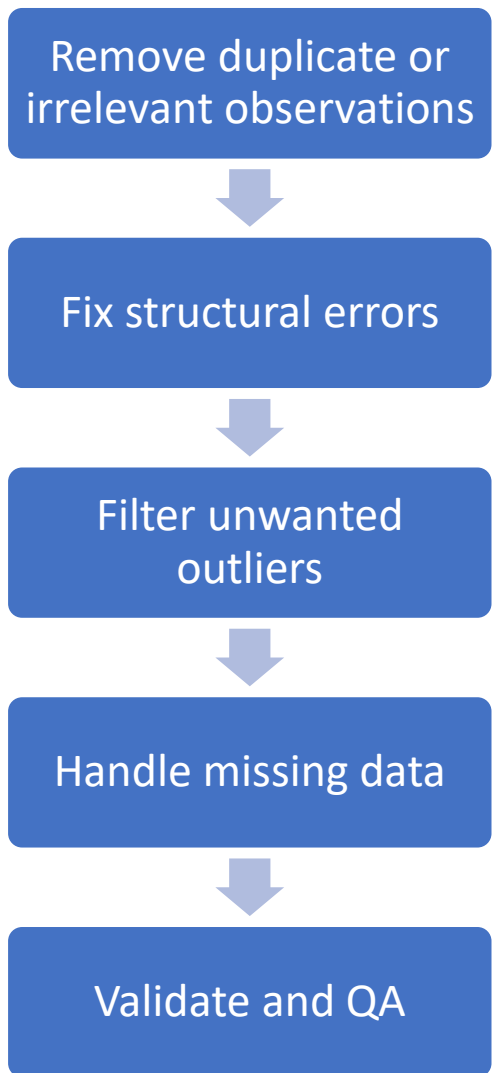


## SUMMARY

- Various Data Cleaning Algorithms And Techniques
- Wide Variety Of Situations
- Data Cleaning Is Very Necessary Part
- Cleaning Methods And Approaches Depend Upon  
The Type Of Data Comparison Of Data Cleaning Tools  
And Determines The Best Tool.

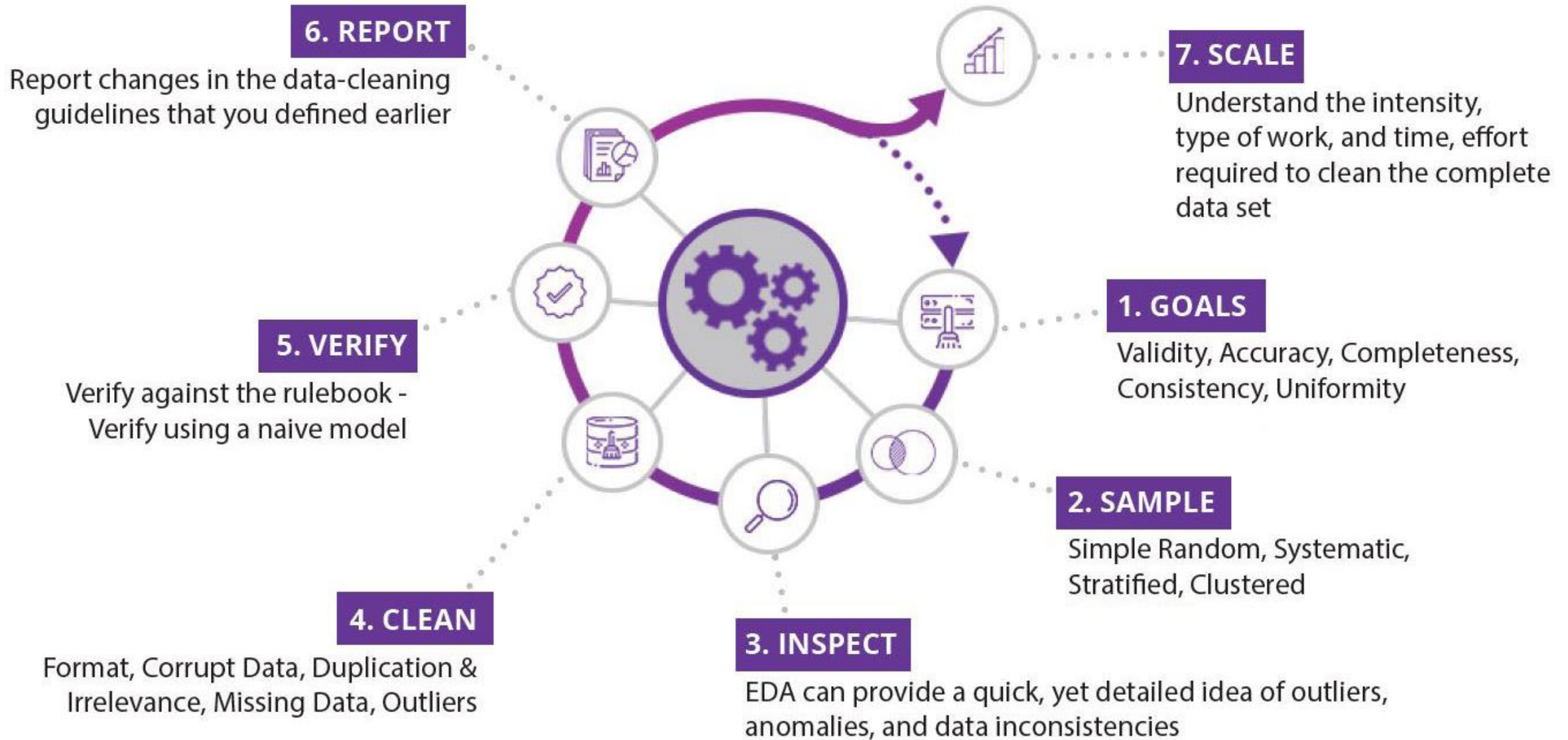


# APPROACH FOR DATA CLEANING





# DATA CLEANING ITERATIVE MODEL



## **PROBLEM STATEMENT:**

To analyse a data set containing uncleaned data and then apply various data cleaning methods on the data set according to the requirement.

## **OBJECTIVE:**

- To apply the general understanding of data cleaning.
- To analyze a dataset.
- To use the data cleaning tools.
- To apply various data cleaning methods on a data set.



# DESIGN OF SOLUTION

Color	Martin Scorsese	240	116866727	The Wolf of Wall Street	English	USA	100000000	2013	8.2
Color	Shane Black	195	408992272	Iron Man 3	English	USA	200000000	2013	7.2
color	Quentin Tarantino	187	54116191	The Hateful Eight	English	USA	44000000	2015	7.9
Color	Kenneth Lonergan	186	46495	Margaret	English	usa	14000000	2011	6.5
Color	Peter Jackson	186	258355354	The Hobbit: The Desolation of Smaug	English	USA	225000000	2013	7.9
	N/A	183	330249062	Batman v Superman: Dawn of Justice	English	USA	250000000	202	6.9
Color	Peter Jackson	-50	303001229	The Hobbit: An Unexpected Journey	English	USA	180000000	2012	7.9
Color	Edward Hall	180		Restless	English	UK		2012	7.2
Color	Joss Whedon	173	623279547	The Avengers	English	USA	220000000	2012	8.1
Color	Joss Whedon	173	623279547	The Avengers	English	USA	220000000	2012	8.1
	Tom Tykwer	172	27098580	Cloud Atlas	English	Germany	102000000	2012	-7.5
Color	Null	158	102515793	The Girl with the Dragon Tattoo	English	USA	90000000	2011	7.8
Color	Christopher Spencer	170	59696176	Son of God	English	USA	22000000	2014	5.6
Color	Peter Jackson	164	255108370	The Hobbit: The Battle of the Five Armies	English	New Zealand	250000000	2014	7.5
Color	Tom Hooper	158	148775460	Les Misérables	English	USA	61000000	2012	7.6
Color	Tom Hooper	158	148775460	Les Misérables	English	USA	61000000	2012	7.6

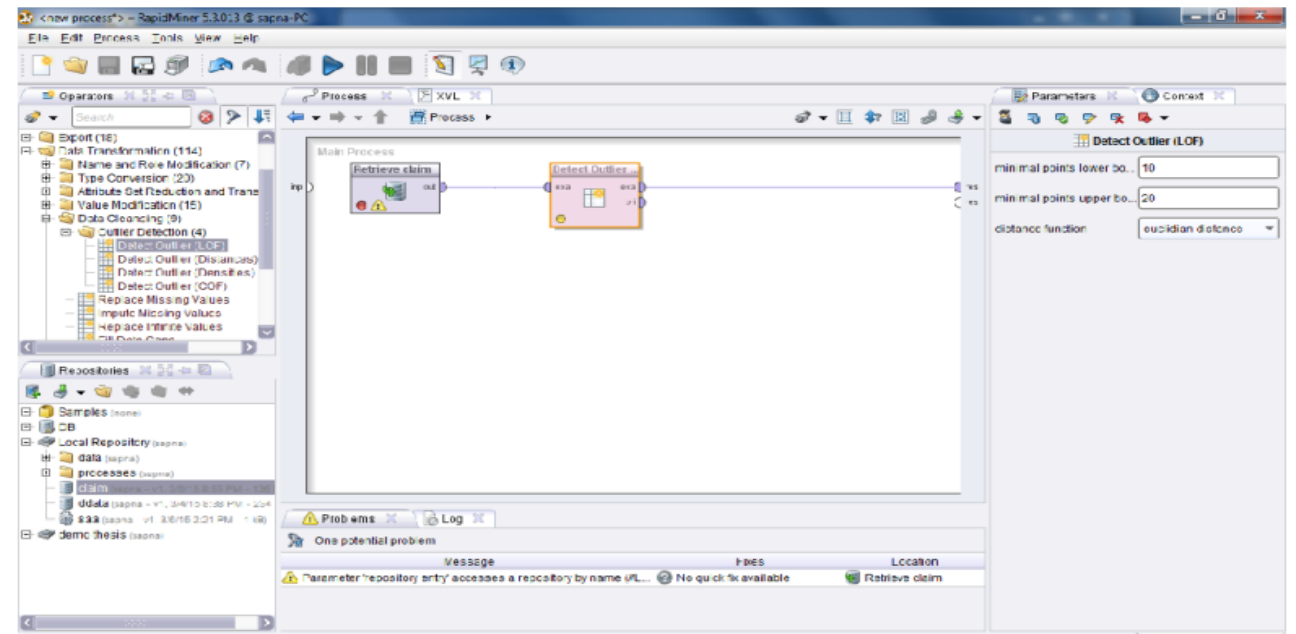
## Example of Unclean Data Set

	Name	Height	Roll	Department	Address
0	A	5.2	55	CSE	polashi
1	B	5.7	99	EEE	banani
2	C	5.6	15	BME	farmgate
3	D	5.5	88	CSE	mirpur
4	E	5.3	1	ME	dhanmondi
5	F	5.8	12	ME	ishwardi
6	G	5.6	47	CE	khulna
7	H	5.5	104	CSE	uttara

### Example of Cleaned Data Set

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	Name	salary	totcomp	tenure	age	sales	profits	assets	Email										
2	robert jhon	3030	8138	7	61	161315.0	2356.0	257389.0	robert.jhon@gmail.com										
3	mary	6050	14538	8	51	144416.0	22071.0	237545.0	mary101@yahoo.com										
4	ajay choudhary	3571	7433	11	63	133208.0	4430.0	49271.0	ajay.choudhary@gmail.com										
5	kiran thakur	3300	13454	6	60	102697.0	6370.0	92630.0	ajay.choudhary@gmail.com										
6	kirti thakur	10030	68235	18	63	103469.0	9296.0	365935.0	kirti.thakur@gmail.com										
7	sehi	9375	42391	6	57	81667.0	6326.0	86100.0	sehi101@gmail.com										
8	suman baa	9525	21155	15	60	75431.0	5307.0	668641.0	suman.baa@gmail.com										
9	aran sharma	5030	24424	5	61	57813.0	5372.0	59920.0	aran.sharma@gmail.com										
10	varun kashyap	999	2916	3	57	55154.0	1120.0	36672.0	varun.kashyap@gmail.com										
11	robert jhon	3300	7457	2	60	53588.0	6396.0	59550.0	robert.jhon11@gmail.com										
12	sanjay sharma	3530	3677	16	63	53777.0	5165.0	617679.0	sanjay.sharma@gmail.com										
13	sanjay sharma	2493	6728	6	61	47678.0	1704.0	42764.0	sanjay.sharma@yahoo.com										
14	sehi thakur	1911	4727	7	58	47061.0	2345.0	33673.0	sehi.thakur@gmail.com										
15	sehi thakur	2130	2383	4	59	41322.0	1046.0	37675.0	sehi102@gmail.com										
16	sanjay sharma	1135	4358	8	56	37154.0	3780.0	30966.0	sanjay.sharma@yahoo.com										
17	suresh kumar	5236	20592	2	60	35853.0	1255.0	299904.0	suresh.kumar@gmail.com										
18	suresh sharma	1930	14936	4	60	33674.0	568.0	14166.0	suresh.sharma@gmail.com										
19	sushel sharma	6030	31238	32	74	33296.0	3765.6	194398.0	sushel.sharma@gmail.com										
20	kritika thakur	6229	1621	5	63	32379.0	3762.0	365875.0	kritika.thakur@gmail.com										
21	kirti thakur	1523	4354	3	56	31707.0	576.0	28570.0	kirti.thakur@gmail.com										
22	kirti thakur	2050	6854	4	52	31565.0	2365.3	55143.0	kirti.thakur10@gmail.com										
23	sehi thakur	4417	8623	15	57	31260.0	703.0	29350.0	sehi.thakur@gmail.com										
24	sehi sharma	8838	48952	17	55	31131.0	3276.0	317590.0	sehi.sharma@gmail.com										
25	robert jhon	3343	1931	5	56	30951.0	935.0	15666.0	robert.jhon111@gmail.com										
26	aran thakur	1852	1821	4	57	33678.0	594.0	23638.0	aran.thakur12@gmail.com										
27	sandeep rood	2930	3116	2	56	33219.0	1514.0	13465.0	sandeep.rood@gmail.com										
28	sushel sharma	2952	8536	2	52	33147.0	970.0	26720.0	sushel.sharma10@gmail.com										
29	sunita sharma	1830	2327	2	49	29398.0	-362.0	28728.0	Sunita.sharma105@gmail.com										
30	sushel choudhary	6153	10112	9	57	23777.0	4323.0	45066.0	sushel.choudhary@gmail.com										
31	suresh kumar	1136	3327	9	59	29303.0	410.8	8780.0	suresh.kumar102@gmail.com										
32	suman sharma	2550	2557	6	58	25898.0	6718.0	31853.0	suman.sharma106@gmail.com										

## MS EXCEL WITH DATA CLEANER



## RAPIDMINOR

WinPure Clean and Match (Trial Version) - Project

File Home Settings Help DATA CLEAN MATCH

Select Refresh No Filter View Move Wizard

Table1 'ceo' imported from : ceo.xls

F1	F2	F3	F4	F5	F6	F7	F8	F9
Name	salary	totcomp	tenure	age	sales	profits	assets	Email
robert.jhon	3030	8138	7	61	161315.0	2956.0	257389.0	robert.jhon@gmail...
mary	6050	14530	0	51	144416.0	22071.0	237545.0	mary101@yahoo...
ajay choudhary	3571	7433	11	63	139208.0	4430.0	48271.0	ajay.choudhary@...
kiran thakur	3300	13464	6	60	100697.0	6370.0	92630.0	ajay.choudhary@...
koti thakur	10000	68285	18	63	100469.0	9296.0	355935.0	koti.thakur@gmail...
sehw	9375	42381	6	57	81667.0	6328.0	86100.0	sehw101@gmail.c...
suman bala	9525	12165	15	60	76431.0	5807.0	668641.0	suman.bala@gm...
arun sharma	5000	24424	5	61	57813.0	5372.0	59920.0	arun.sharma@gm...

## WINPURE CLEAN AND MATCH



### Comparison of Data Cleaning tools

<b>Tools Problems</b>	<b>MS Excel with data cleaner</b>	<b>RapidMinor</b>	<b>Winpure Clean &amp; Match</b>
<b>Missing Values</b>	<b>No</b>	<b>Yes</b>	<b>Yes</b>
<b>Availability</b>	<b>Desktop</b>	<b>Desktop</b>	<b>Desktop</b>
<b>Duplication</b>	<b>Yes</b>	<b>Yes</b>	<b>Yes uses the matching</b>
<b>Illegal Values Elimination</b>	<b>No</b>	<b>No</b>	<b>Yes</b>
<b>Misspelling</b>	<b>No</b>	<b>No</b>	<b>No</b>
<b>Merge</b>	<b>No</b>	<b>Yes</b>	<b>Yes</b>
<b>File Format</b>	<b>Excel</b>	<b>CSV, Database, Excel, Access, binary, XML</b>	<b>Text files, Excel , commercial DBMS,</b>
<b>Ease of use</b>	<b>Moderate</b>	<b>Moderate</b>	<b>High</b>

## **Research Done In Past**

- Data Cleaning for Misspelled Proper Nouns (Border Detection Algorithm)
- Robust and Efficient Fuzzy Match for Online Data Cleaning(Fuzzy Match similarity Algorithm)
- Data Cleaning by Clustering and Association Methods (Data Mining Algorithms)

## DIFFERENT APPROACHES AND COMPARATIVE ANALYSIS

	Border Detection Data Algorithm	Data Mining Algorithm-Attribute Correction Algorithm	Fuzzy Match Similarity Function Algorithm
Features	Simple, effective to compute clusters in the validated against reference to match the reference data then string data.	The given attributes are validated against reference to match the reference data to provide cleansing solution fuzzy match similarity (fms) that explicitly considers IDF token weights and input errors while comparing tuples.	data mining techniques in the area of attribute correction are: context-independent attribute correction implemented using clustering and robust results . If one techniques and context-dependent
Significance / performance	It produces good cleansing results for string data with large distances between centers of clusters and small distances within the clusters	Quality of fms is better than ed (edit distance) using two Datasets.	Algorithm shows better results for longer strings

- PROJECT STATUS
- PLANNING
- FUTURESSCOPE





[illegible]

# REFERENCES

- [1] STEPS TO DO DATA CLEANING: <https://www.tableau.com/learn/articles/what-is-data-cleaning>
- [2] DATA CLEANING: CURRENT APPROACHES AND ISSUES : Vaishali Chandrakant Wangikar and Ratnadeep R. Deshmukh  
[https://www.researchgate.net/publication/278301609 Data Cleaning Current Approaches and Issues](https://www.researchgate.net/publication/278301609_Data_Cleaning_Current_Approaches_and_Issues)
- [3] STUDY OF DATA CLEANING & COMPARISON OF DATA CLEANING TOOLS : Sapna Devi, Dr Arvind Kalia  
<https://www.ijcsmc.com/docs/papers/March2015/V4I3201599a30.pdf>
- [4] AN OVERVIEW STUDY ON DATA CLEANING, ITS TYPES AND ITS METHODS FOR DATA MINING : S.Lakshmi AND Dr S.V  
<https://acadpubl.eu/hub/2018-119-12/articles/6/1564.pdf>
- [5] ITERATIVE PROCESS FOR DATA CLEANING  
<https://innotescus.io/data-cleaning/complete-guide-iterative-process-for-data-cleaning/>
- [6] FUZZY MATCHING ALGORITHM  
[https://nanonets.com/blog/fuzzy-matching-fuzzy-logic/#:~:text=Fuzzy%20Matching%20\(also%20called%20Approximate,Priceline%20in%20the%20graphic%20below](https://nanonets.com/blog/fuzzy-matching-fuzzy-logic/#:~:text=Fuzzy%20Matching%20(also%20called%20Approximate,Priceline%20in%20the%20graphic%20below)
- [7] BORDER DETECTION ALGORITHM  
[https://www.researchgate.net/publication/220803114 Cleansing Databases of Misspelled Proper Nouns](https://www.researchgate.net/publication/220803114_Cleansing_Databases_of_Misspelled_Proper_Nouns)
- [8] A review on data cleaning methods for Big Data : Fakhitha Ridzuan  
[https://www.researchgate.net/publication/338348131 A Review on Data Cleansing Methods for Big Data](https://www.researchgate.net/publication/338348131_A_Review_on_Data_Cleansing_Methods_for_Big_Data)





