

# Big Data Platform for Health and Safety Accident Prediction

## Abstract

### Purpose

This paper highlights the use of the Big Data technologies for health and safety risks analytics in the power infrastructure domain with large datasets of health and safety risks, which are usually sparse and noisy.

### Design/methodology/approach

The study focuses on using Big Data frameworks for designing a robust architecture for handling and analysing (exploratory and predictive analytics) accidents in power infrastructure. The designed architecture is based on a well coherent health risk analytics lifecycle. A prototype of the architecture interfaced various technology artefacts was implemented in the Java language to predict the likelihoods of health hazards occurrence. A preliminary evaluation of the proposed architecture was carried out with a subset of an objective data, obtained from a leading UK power infrastructure company offering a broad range of power infrastructure services.

### Findings

The proposed architecture was able to identify relevant variables and improve preliminary prediction accuracies and explanatory capacities. It has also enabled conclusions to be drawn regarding the causes of health risks. The results represent a significant improvement in terms of managing information on construction accidents, particularly in power infrastructure domain.

### Originality/value

This study carries out a comprehensive literature review to advance the health and safety risk management in construction. It also highlights the inability of the conventional technologies in handling unstructured and incomplete dataset for real-time analytics processing. The study proposes a technique in Big Data technology for finding complex patterns and establishing the statistical cohesion of hidden patterns for optimal future decision-making.

Keywords: Big Data analytics, Machine learning, Health hazards analytics, Health and Safety

## 1. Introduction

Occupational accidents are things of worry in modern society, especially in construction sites where a high number of construction activities take place (Zhu et al. 2016). The power infrastructure delivery sector, for instance, has high incidences of nonfatal occupational injuries as workers using heavy machinery are confronted with health risks such as radiation, dust, temperature extremes, and chemicals amongst others (McDermott & Hayes 2016). According to the UK Health and Safety Executive, a total cost of £4.8 billion was expended in 2014/15 for workplace injury (HSE 2016). Similarly, repair costs of buried communication lines are significant when disrupted during excavations (McDermott & Hayes 2016).

Several machine-learning techniques have been used for health and safety risks prediction in construction. For instance, decision trees (Cheng et al. 2011), the generalised linear model (Esmaeili et al. 2015), and fuzzy-neural method (Debnath et al. 2016) have all been used to analyse incident data to reduce accident rates. Techniques such as the Bayesian network was

51 used to quantify occupational accident rates (Papazoglou et al. 2015), and fuzzy Bayesian  
52 networks for damaged equipment analysis (Zhang et al. 2016). Others are the bow tie  
53 representation for occupational risks assessment (Jacinto & Silva 2010), and Poisson models  
54 for occupational injury impacts modelling (Yorio et al. 2014).

55

56 However, a significant problem associated with these existing models is their limited ability to  
57 process large-scale raw data since considerable effort is needed to transform them into an  
58 appropriate internal form to achieve high prediction accuracy (Esmaeili et al. 2015).  
59 Construction accident data are typically large, heterogeneous and dynamic (Fenrick &  
60 Getachew 2012), nonlinear relationships among accident causation variables (Gholizadeh &  
61 Esmaeili 2016), imbalance data, and appreciable missing values (Bohle et al. 2015). Besides,  
62 these techniques simplify some key factors and pay little attention to analysing relationships  
63 between a safety phenomenon and the safety data (Landset et al. 2015).

64 Based on the preceding, the Big Data technology due to its parallel processing feature and  
65 ability to efficiently handle high dimensional, noisy data with nonlinear relationships, will be  
66 beneficial for health and safety risks analytics in the power infrastructure domain. Also, the  
67 technology will uncover potential factors contributing to accidents in this domain. The objectives  
68 of this study are, therefore, to chart lifecycle stages of occupational hazards analytics and  
69 develop a Big Data architecture for managing health and safety risks.

70

#### 71 *1.1. Big data for health and safety risk analytics*

72

73 Big Data is an emerging technology, which refers to data sets that are many orders of  
74 magnitude larger than the standard files transmitted via the Internet (Suthakar et al. 2016).  
75 There is tremendous interest in utilising information in Big Data for various analytics  
76 (exploratory, descriptive, predictive and prescriptive) to determine future occurrences. Most  
77 importantly, Big Data technologies support analytical techniques for occupational health and  
78 safety risk analytics; thus, a system being proposed in this study, named Big Data Accident  
79 Prediction Platform (B-DAPP) offers unparalleled opportunities to minimise occupational  
80 hazards at construction sites. The seamless combination of the following technologies: Big  
81 Data, Health and Safety, and Machine-learning is an outcome of a robust health and safety risk  
82 management tool to help stakeholders in making appropriate decisions to minimise  
83 occupational accidents in Power Infrastructure projects.

84 Health and safety risk analytics is dependent on a high-performance computation and large-  
85 scale data storage requiring a large number of diverse datasets of health and safety risks, and  
86 machine-learning knowledge to successfully provide the needed analytical responsibilities. The  
87 datasets, however, are unreliable, unstructured, incomplete, and imbalanced (Chen et al.  
88 2017). Hence, storing the datasets using conventional technologies and subjecting them to  
89 real-time processing for advanced analytics is highly challenging. A robust technique for finding  
90 complex patterns and establishing the statistical cohesion of hidden patterns in such datasets

91 for optimal future decision-making is inevitable. Thus, motivating the use of Big Data  
92 technologies to address these challenges.

93

#### 94 *1.2. Research Justification*

95 There exists an apparent technological gap in existing literature regarding health and safety  
96 risk management. In particular, there is limited research on the application of Big Data  
97 techniques for managing health and safety risk in Power Infrastructure. The development of a  
98 robust B-DAPP for health and safety risk is the objective of the ongoing R&D effort. The  
99 proposed tool will provide stakeholders with well-informed and data-driven insights to reduce  
100 accidents and incidents at construction sites. Therefore, a Big Data architecture is proposed  
101 for managing health and safety risks. Also, a presentation of components and relevant  
102 technologies of the proposed architecture necessary for storing and analysing health and safety  
103 risk datasets for real-time exploration and prediction is made. The term 'Architecture' as used  
104 in this text refers to high-level structures of a software system. Similarly in the context of this  
105 study, 'Accident' is an unplanned, unpremeditated event caused by unsafe acts or conditions  
106 resulting in injury while 'Incident' is an event causing actual damage to property (including plant  
107 or equipment) or other loss with potential to cause injury.

108

109 The remainder of the paper is structured as follows: Section 2 discusses on the research  
110 methodology, Big Data analytics, and Big Data ecosystem. Section 3 deliberates on the health  
111 hazards analytics lifecycle. Section 4 presents the proposed Big Data architecture for health  
112 and safety risk management while Section 5 presents the preliminary outcomes. Conclusions  
113 and future work are given in Section 6.

114

115

## 116 **2. Methodology**

117

118 In this section, a discussion on the methodology employed in this research is made. Foremost,  
119 a comprehensive literature review is performed to advance the health and safety risk  
120 management with respect to the system architecture and system analytics lifecycle. Then the  
121 proposed architecture and occupational hazard analytics lifecycle are validated in a preliminary  
122 analysis of the health and safety risk related data. To be able to offer a holistic Big Data  
123 architecture and occupational hazard analytics lifecycle, a careful review of existing literature  
124 on health and safety risk prediction models, Big Data, and machine learning have been carried  
125 out. In this regard, online databases such as Journal of Big Data, Big Data Research, Safety  
126 Science, Journal of construction engineering, Journal of Decision Systems, Journal of Safety  
127 Research, Journal of Construction Engineering and Management, Reliability Engineering and  
128 System Safety are searched for research articles between 2005 and 2017. Recent reviews of  
129 research and books on Big Data Analytics are also considered (Camann et al. 2011; Gandomi  
130 & Haider 2015; Guo et al. 2016).

131

132 Examples of search words used include: “managing health and safety risks”, “design strategies  
133 for occupational hazards in construction”, “Prediction models for occupational health risks”, “Big  
134 Data in Construction”, “Big Data based Application Architecture”, and “Big Data Analytics”. In  
135 general, 94 publications were selected even though literature search was in-exhaustive as a  
136 result of a vast amount of published articles. However, it is believed that the literature search  
137 has captured a representative balanced sample of the related research. Studies in which Big  
138 Data is used to develop enterprise applications were included, and those focusing on road  
139 traffic related hazards and health hazards in domains not related to construction (e.g., mining  
140 and fishing) were excluded. This elimination procedure further reduced the selected articles to  
141 66. These articles are furthermore scrutinised for relevancy by reading abstracts, introductions,  
142 and conclusions. Ultimately, the articles are reduced to 50. Table 1 depicts how these selected  
143 articles are relevant and contributing to the development of the proposed architecture, which is  
144 essentially based on three concepts, namely Big Data, Health and safety risk, and Machine  
145 learning. In this study, we introduce the proposed B-DAPP architecture and the occupational  
146 hazards analytics lifecycle stages for managing incidents and accidents.

147

#### 148 *2.1. Big data analytics*

149 Big data consists of large and complex datasets often difficult to manipulate using the  
150 conventional processing methods. It has six defining attributes (Gandomi & Haider 2015), which  
151 are volume, variety, velocity, veracity, variability and complexity, and value. The term 'volume'  
152 represents the magnitude of the data (measured in terabytes, petabytes and beyond). 'Variety'  
153 is the structural heterogeneity in a dataset while the 'Velocity' is the rate of generating data.  
154 'Veracity' is the unreliability inherent in data sources while 'Variability' (complexity) represents  
155 the variation in data flow rates. Finally, 'Value' measures the information extracted from  
156 historical incident datasets for optimal control decision to mitigate incidents and reduce their  
157 impact.

158 These attributes are evident in a typical power infrastructure health and safety dataset, which  
159 is typically large, heterogeneous and dynamic (Ferrick & Getachew 2012). Big data analytics  
160 is a concept that inspects, cleans, transforms, and models the big data to discover useful  
161 information to support decision-making (Power 2014). The Big data analytics have rich  
162 intellectual traditions and borrow from a wide variety of related fields such as statistics, data  
163 mining, business analytics, knowledge discovery from data (KDD), and data science. The forms  
164 of big data analytics are descriptive (Schryver et al. 2012), predictive (Esmaeili et al. 2015),  
165 prescriptive (Delen & Demirkan 2013) and causal (Schryver et al. 2012).

166

#### 167 *2.2. Big Data for safety risk management*

168 A wide variety of technologies and heterogeneous architectures are available to implement Big  
169 data applications. Since this paper intends to develop a robust Big Data architecture for health  
170 hazards analytics, A brief discussion of tools and Big Data platforms to facilitate the creation of  
171 a compact architecture and increase the understanding of the concept is made. Primarily,  
172 focusing on the Hadoop ecosystem, a system designed for solving Big Data problems.

173  
174  
175

Table 1: Summary of articles reviewed

#	Article	Contribution to health and safety risk analytics architecture		
		Health and safety risk	Machine learning	Big data
1	Liu & Tsai (2012)	×	×	
2	Zhou et al. (2015)	×		
3	García-Herrero et al. (2012)	×	×	
4	Groves et al. (2007)	×		
5	Li et al. (2016)	×	×	
6	Soltanzadeh et al. (2016)		×	
7	Power (2014)			×
8	Yi et al. (2016)	×	×	
9	Cheng et al. (2011)	×	×	
10	Silva et al. (2016)	×		
11	Raviv et al. (2017)	×		
12	Liao & Perng (2008)		×	
13	Li & Bai (2008)			
14	Törner & Pousette (2009)	×		
15	Pinto et al. (2011)	×		
16	Tixier et al. (2016)	×	×	
17	Hallowell & Gambatese (2009)	×		
18	Pääkkönen & Pakkala (2015)			×
19	Venturini et al. (2017)			×
20	Suthakar et al. (2016)			×
21	Najafabadi et al. (2015)		×	×
22	Landset et al. (2015)			×
23	Tsai et al. (2015)			×
24	Zang et al. (2014)		×	×
25	Jin et al. (2015)			×
26	Rahman & Esmailpour (2016)			×
27	Al-Jarrah et al. (2015)			×
28	Zhang et al. (2016)	×	×	
29	Love & Teo (2017)	×	×	
30	Rivas et al. (2011)	×	×	
31	Guo et al. (2016)	×		×
32	Zou et al. (2007)	×		
33	Wu et al. (2010)	×		
34	Carbonari et al. (2011)	×		
35	Weng et al. (2013)	×	×	
36	Naderpour et al. (2016)	×	×	
37	Yoon et al. (2016)	×		
38	Favarò & Saleh (2016)	×	×	
39	Jocelyn et al. (2017)	×	×	
40	Papazoglou et al. (2017)	×	×	
41	Papazoglou et al. (2015)	×	×	
42	Fragiadakis et al. (2014)	×	×	
43	Ciarapica & Giacchetta (2009)	×	×	
44	Khakzad et al. (2015)	×	×	
45	Galizzi & Tempesti (2015)	×		
46	Gürçanlı & Mungena (2009)	×	×	
47	Debnath et al. (2016)	×	×	
48	Nanda et al. (2016)	×	×	
49	Zeng et al. (2008)	×		
50	Guo et al. (2016)	×	×	

176  
177

178 2.2.1. *Hadoop ecosystem*

179 Hadoop is a MapReduce processing engine with distributed file systems (White 2012).  
180 However, it has evolved into a vast web of projects (Hadoop ecosystem) related to every step  
181 of a Big Data workflow. The concept now is being referred to as the Hadoop ecosystem, which  
182 encompasses related projects and products developed to either complement or replace original  
183 components. Further examination of the two concepts for ease of understanding follows.

184

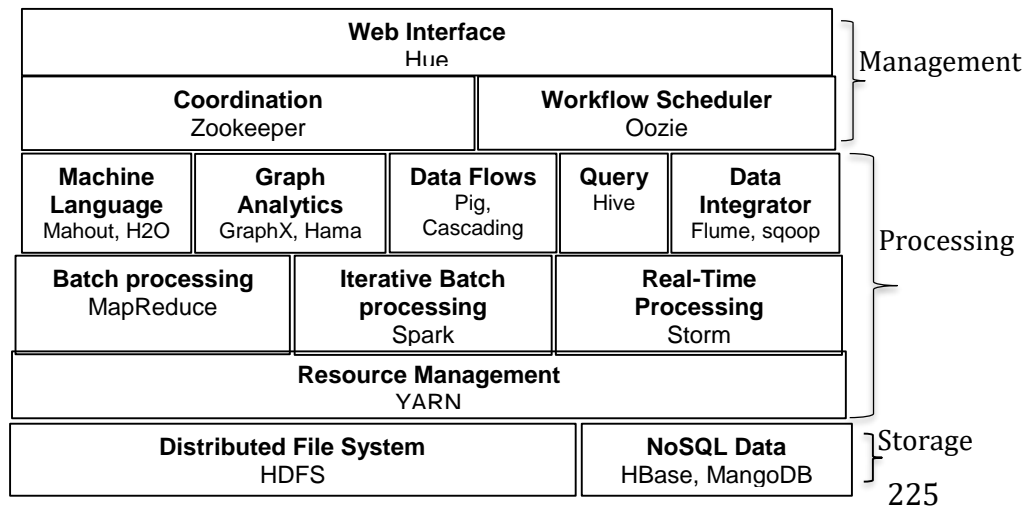
185 The Hadoop project consists of four modules (White 2012):

- 186 a) Hadoop distributed file system (HDFS) is a fault-tolerant file system designed to store  
187 massive data across multiple nodes of commodity hardware. It has a master-slave  
188 architecture that is made up of data nodes and name nodes. Data nodes store blocks  
189 of the data, retrieve data on request and report to the name node with inventory. The  
190 name node keeps records of the inventory and directs traffic to the data nodes upon  
191 client requests.
- 192 b) MapReduce Data processing engine. A MapReduce job consists of a map phase and  
193 a reduce phase. A map phase organises raw data into key/value pairs, while the reduce  
194 phase processes data in parallel.
- 195 c) YARN (“Yet Another Resource Negotiator”) is a resource manager of the Hadoop  
196 project introduced to address the limitations of the MapReduce. It separates  
197 infrastructures from program representations.
- 198 d) Common is a set of utilities required by the other Hadoop modules. These include  
199 compression codecs, I/O utilities, error detection, proxy users authorisation,  
200 authentication, and data confidentiality.

201

202 The Hadoop ecosystem consists of several tools built on top of the core Hadoop modules  
203 described above to support researchers and practitioners in all aspects of data analyses. The  
204 ecosystem structure has the following layers: storage, processing, and management. Figure 1.  
205 depicts examples of standard tools used in Big Data applications. The right selection requires  
206 in-depth knowledge of critical features of these platforms and the characteristics of the problem  
207 to be solved. In the case of health hazards analytics the platforms to adapt as a result of  
208 increased workload, outweighs the rest of the selection criteria. In the real sense, Hadoop  
209 ecosystem is made up of well over 100 projects, and readers are referred to (White 2012) or  
210 the Hadoop website for more information.

211  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224



226

Figure 1. Hadoop ecosystem

227

228

- a) Storage layer- This layer includes the HDFS described earlier and Non-relational databases (NoSQL). Non-relational databases are nested, semi-structured, and unstructured data that support machine-learning tasks. These databases use the following data representation models: Key-value stores (i.e. Redis), Document stores (i.e. MongoDB), Column-oriented Data (i.e. HBase), and Graph-based models (Neo4J). The graph model is regarded as more flexible than other models.

234

235

- b) Processing layer - This layer carries out the actual analysis using YARN, which allows one or more processing engines to run on a Hadoop cluster. Additionally, a layer has frameworks for data transfer, aggregation, and interaction. Examples include Flume, Sqoop, Hive, Spark, and Pig. Flume collects, aggregates, and moves data log in HDFS. Kafka is a distributed messaging system on HDFS, and Sqoop transports bulk data between the HDFS and relational databases. Hive is a query engine for querying data stored in the HDFS and NoSQL databases. Spark supports iterative computation, and it improves on speed and resource issues by utilising in-memory computation. Finally, Pig offers an execution framework and data flow language to support user-defined functions written in Python, Java, JavaScript, etc. Machine learning frameworks are used to perform machine-learning tasks in Hadoop. Examples are Mahout, H2O, etc. Mahout is one of the more well-known machine-learning tools. It is known for having a wide selection of robust algorithms, but with inefficient runtimes due to the slow MapReduce engine. H2O provides a parallel processing engine, analytics, math, and machine learning libraries for data pre-processing and evaluation.

250

251

- c) Management layer - This layer has tools for user interaction and high-level organisation. It carries out functions such as scheduling, monitoring, coordination, amongst others. Examples of tools available in this layer are Oozie, Zookeeper, and Hue. Oozie is a workflow scheduler, which manages jobs for many of the tools in the

252

253

254

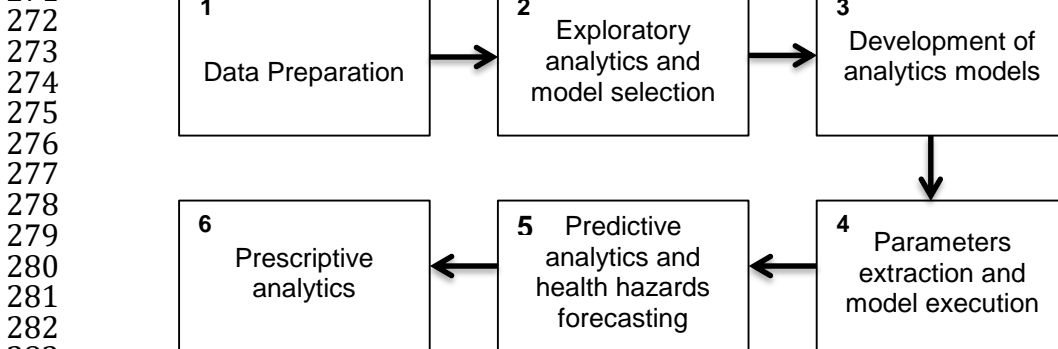
255 processing layer. Zookeeper provides tools to handle the coordination of data and  
256 protocols and can handle partial network failures. It includes APIs for Java and C and  
257 also has bindings for Python and REST clients. Hue is a web interface for Hadoop  
258 projects with support for widely used Hadoop ecosystem components.

259

### 260 3. Proposed health hazards analytics stages

261 Developing a health hazards analytics tool for health and safety risk data is a challenging task  
262 since the data are typically dynamic (Fenrick & Getachew 2012), and unbalanced with  
263 significant missing values (Bohle et al. 2015). Besides, the traditional accident-causing  
264 modelling may ignore or simplify some key factors as well as assume the same format for the  
265 input data. Thus, an efficient methodology to address these challenges requires a well-  
266 articulated process to break the task into smaller manageable stages to ensure adequate  
267 preparation of various analytical approaches. In this section, a discussion on the lifecycle of the  
268 proposed Big Data architecture for the health hazards analytics tool is made. The lifecycle has  
269 six stages (see Figure 2) that are iteratively executed to suit the requirements of the proposed  
270 tool.

271



283

284 Figure 2. Stages of the health hazards analytics

285

#### 286 3.1. Data preparation

287 Data preparation is a procedure to detect and repair errors in the dataset. For the health  
288 hazards analytics, sufficient data quality is necessary for high-quality analytics. Thus, data from  
289 various sources are obtained, transformed, and loaded into the centralised data store. Before  
290 this, outliers are inadvertently eliminated using techniques such as mean/mode imputation,  
291 transformation, and binning. Missing data issues should also be solved using appropriate  
292 technology. The k-nearest neighbour (kNN) imputation and mean/mode imputation are few  
293 examples to eliminate the missing data problem. Apparently, machine-learning techniques can  
294 also be applied to quickly filter through hundreds of thousands of narratives (texts) to accurately  
295 and consistently retrieve and track high-magnitude, high-risk and emerging causes of injury.  
296 The retrieved information is then utilised to guide the development of interventions to prevent  
297 future incidents.

298 In the event of having large data, methods for parallel data movement may be required, which  
299 may necessitate using the appropriate component of the Hadoop ecosystem. Data is often



300 analysed to get familiar with the health and safety risk as it pertains to the construction domain.  
301 For the sake of preliminary analysis presented here, the health and safety data are provided  
302 as .csv files that are stacked on the Hadoop cluster. The respective files are queried to retrieve  
303 specific details on health and safety hazards such as injured body parts, loss type injury, and  
304 damaged equipment amongst others. For this purpose, tools like Apache Flume are of  
305 immense relevance to capture current versions of datasets.

306

### 307 *3.2. Exploratory analytics and model selection*

308 For the health hazards management, the analysis starts with exploratory analytics and then to  
309 the predictive analytics. For each activity in the proposed tool, a clear objective is essential for  
310 the right selection of analytical approaches (prescription, exploratory, predictive, etc.) to  
311 execute. The data exploration of health and safety records is performed to understand the  
312 relationship between different explanatory variables. This exploratory data analysis informs the  
313 selection of relevant variables to build a robust health hazards prediction model. In this study,  
314 a visualisation technique is used for exploratory data analysis. At this phase, the purpose of the  
315 analysis is to capture essential predictors and independent variables while eliminating the least  
316 relevant ones for building the model. Variable selection methods include All Possible  
317 regression, Stepwise Forward regression, Best Subset regression, etc. These selection  
318 methods are often iterative and require a series of steps to identify the most useful variables  
319 for the given model. Tools such as R Studio could be exploited to build these models.

320

### 321 *3.3. Development of analytics models*

322 In this stage, analytics models are created for health and safety risk prediction using robust Big  
323 Data analytics techniques. The data are divided first into the training and test sets. The analytics  
324 models are then fitted to the training data and evaluated using the test data. Models with optimal  
325 accuracy or higher predictive power are selected. Often, this step may involve dealing with  
326 certain optimisation issues such as multicollinearity. The best model is selected and deployed  
327 to predict health and safety risk from a large volume of data. Many times the production  
328 environment may require adjusting and redeploying models to support more practical situations  
329 (Camann et al. 2011).

330

### 331 *3.4. Parameters extraction and model execution*

332 Here, vital parameters are extracted to execute the predictive models. Parameters such as  
333 task, equipment type, project complexity, etc. are extracted and the relationship between a  
334 safety phenomenon and safety data explored to uncover potential factors that contribute to the  
335 likelihood of accidents. These relationships bring those potential trends into the focus that could  
336 be utilised to predict the health and safety risk of an infrastructure project under execution. A  
337 series of transformations are applied to make the application user-friendly. Specifically, by  
338 standardising contents using the ifcOWL ontology (Chaudhuri & Dayal 1997). The data are then

339 stored as graph-annotated formats to support broader computations required from the  
340 proposed tool.

341

### 342 3.5. *Predictive analytics and health hazards forecasting*

343 Health hazards prediction provides the necessary foundation for understanding causes and  
344 types of health and safety risk arising from a construction project in execution. Thus, this stage  
345 employs predictive models generated through the big data analytics approaches to analyse  
346 health and safety risk database and give notice of a possible health hazard occurrence. Indeed,  
347 the critical thing about this evaluation is the accuracy of the health and safety risk prediction  
348 models that are employed.

349 The traditional accident-causing modelling has the following limitations: may ignore or simplify  
350 some key factors, uses qualitative analysis, and focuses on causality analysis and explanations  
351 of an accident (Landset et al. 2015). Hence, these methods pay little attention to the analysis  
352 of relationships between a safety phenomenon and safety data. They are also unable to  
353 uncover potential factors that contribute to the likelihood of accidents, such as frequency,  
354 relevance, locale, and timeliness.

355 The development of robust health hazards prediction models is the ultimate goal of this  
356 lifecycle, and using the prediction models, comprehensive accident and equipment damage  
357 forecasts are generated to organisations implement strategies and techniques to improve the  
358 safety of their construction sites.

359

### 360 3.6. *Prescriptive analytics*

361 This phase optimises various safety strategies based on myriad factors (the interaction  
362 between deficiencies in work teams, workplace, equipment and materials, weather, etc.) to  
363 recommend the best course of action for a given situation. It uses simulation and optimisation  
364 to offer the best strategy to employ for different health and safety risks. Consequently, a large  
365 number of alternative optimisation plans are generated and converted into user-friendly  
366 prescriptions for stakeholders to aid in data-driven decision-making for minimising accidents.

367

### 368 3.7. *Analysis and preliminary results*

369 The proposed architecture is further assured and validated with the objective data, obtained  
370 from a leading UK construction company, offering a broad range of power infrastructure  
371 services, including building and refurbishing overhead lines, substations, underground cabling,  
372 fibre optics, etc. The company uses a relational database to store the health and safety risks  
373 data, which consist of a large number of power infrastructure projects constructed over 13 years  
374 (2004 to 2016) across five UK regions. Each time an incident (or hazard) occurs, a digital  
375 record is created in the database. Details of some of the relevant explanatory variables in the  
376 database are shown in Table 2.

377

378 A subset of 5000 randomly selected projects from 20000 projects in total was used for a  
 379 preliminary evaluation and analysis presented in this study. The criteria for this selection include  
 380 project types (i.e. overhead lines, cabling, and substations) and construction mode (i.e. new  
 381 built, refurbishment). The distribution of data across the UK regions will help to generate  
 382 advanced visualisations such as geographic heat map. Data from the relational database is  
 383 accessed via the front-end application and exported to comma-separated files (.csv). Plainly,  
 384 occupational hazards data of 5000 projects will not be labelled as Big Data to justify the use of  
 385 data-intensive platforms for its analysis. However, the approach adopted in this study can be  
 386 used to analyse larger sets of health and safety risk data. Exploratory data analytics is applied  
 387 to understand the underlying trends in the data using geographical and chronological  
 388 dimensions. Thus, a variety of visualisations such as bar plot, box plot, and geographic heat  
 389 map are used for data investigation.

390

391 Table 2: Explanatory variables in the database

Variable	Meaning
Incident reference	Identification of a given incident
Project type	The specific project (overhead line, cabling, offshore, etc.)
Project contract	The nature construction project being built (i.e. new built, maintenance, refurbishment)
Region	The specific region of the construction site (Scotland, North, South East, Midlands, etc.)
Sub region	The sub-region where the site is located i.e. Yorkshire East, Midlands North, East England, Tyrone, etc.
City	UK cities where the construction site is located.
Location	A specific area or location of the site
Client	An organisation using the services of the power infrastructure company.
Equipment type	Specifies the machinery (e.g. drill, hammer, haulage, etc.) used for a task.
Age	The age of the victim at the time of the accident.
Year	The year when the health hazard occurred.
Season	External factor such as the weather
Month	The month (1-12) when the incident occurred
Time	The period incident happened (0-6- early morning, 6-12- morning, 12-18 afternoon, 18-23 -evening).
Day of the week	Day (1-31) when the accident occurred.
Weekday	The weekday i.e. Monday, Tuesday, Wednesday, etc.
Task	Specific task or operation to be carried out ( excavating, lifting, cutting, etc.
Accident type	The type of accident, for instance, fall, trip, struck by, Inhalation, Caught in/between, etc.
Injury type	The physical consequence for a victim, i.e. first aid, fatal, no injury, etc.
Severity cost	Financial cost incurred as a result of the accident
Hazard type	Forms of health hazards, for example, illness, injury, loss or damage, etc.,
Injured body part	The part of the body that is injured, i.e. Fingers, shoulder, head, back, etc.
Total cost	The cost of the project
Equipment	Part of the equipment damaged during operation.

392

393

394

395

#### 396 4. Proposed big data architecture for health hazards analytics

397 This section discusses the proposed Big Data architecture for health hazards analytics (see  
398 Figure 3). Components of the architecture are the Application layer, Analytics and Functional  
399 Model layer, Semantic layer, and Data Storage layer which are discussed in subsequent  
400 subsections.

401

##### 402 4.1. *Data storage*

403 This layer is the data source (finance and health and safety risks), which are needed for efficient  
404 functioning of B-DAPP and analytics models (predictive and prescriptive) development. The  
405 finance data includes information such as project cost, margin, labour cost, material cost, etc.  
406 The health and safety data contains historical occupational risk data while multimedia data  
407 consists of images and videos depicting accidents scenes.

408 As a result of the diverse nature of data to be stored in this layer, a NoSQL database (i.e.  
409 MongoDB, Neo4J, Oracle NoSQL) is used for the implementation due to its robust storage  
410 mechanisms and efficient handling of structured, semi-structured and unstructured data (Leavitt  
411 2010).

412

##### 413 4.2. *Semantic layer*

414 This layer provides the data exchange formatting and data provisioning to the application layer.  
415 The data exchange formatting allows the sharing of a common data format in the entire system.  
416 The DDAXML is used to share data among different modules in the system since it is an  
417 industrially supported schema for sharing information. The data provisioning functionality  
418 provides the application layer of the architecture with seamless access to databases through  
419 the Representation State Transfer (REST) web service. This database access approach is  
420 considered the most appropriate due to the different nature of health and safety risk data.

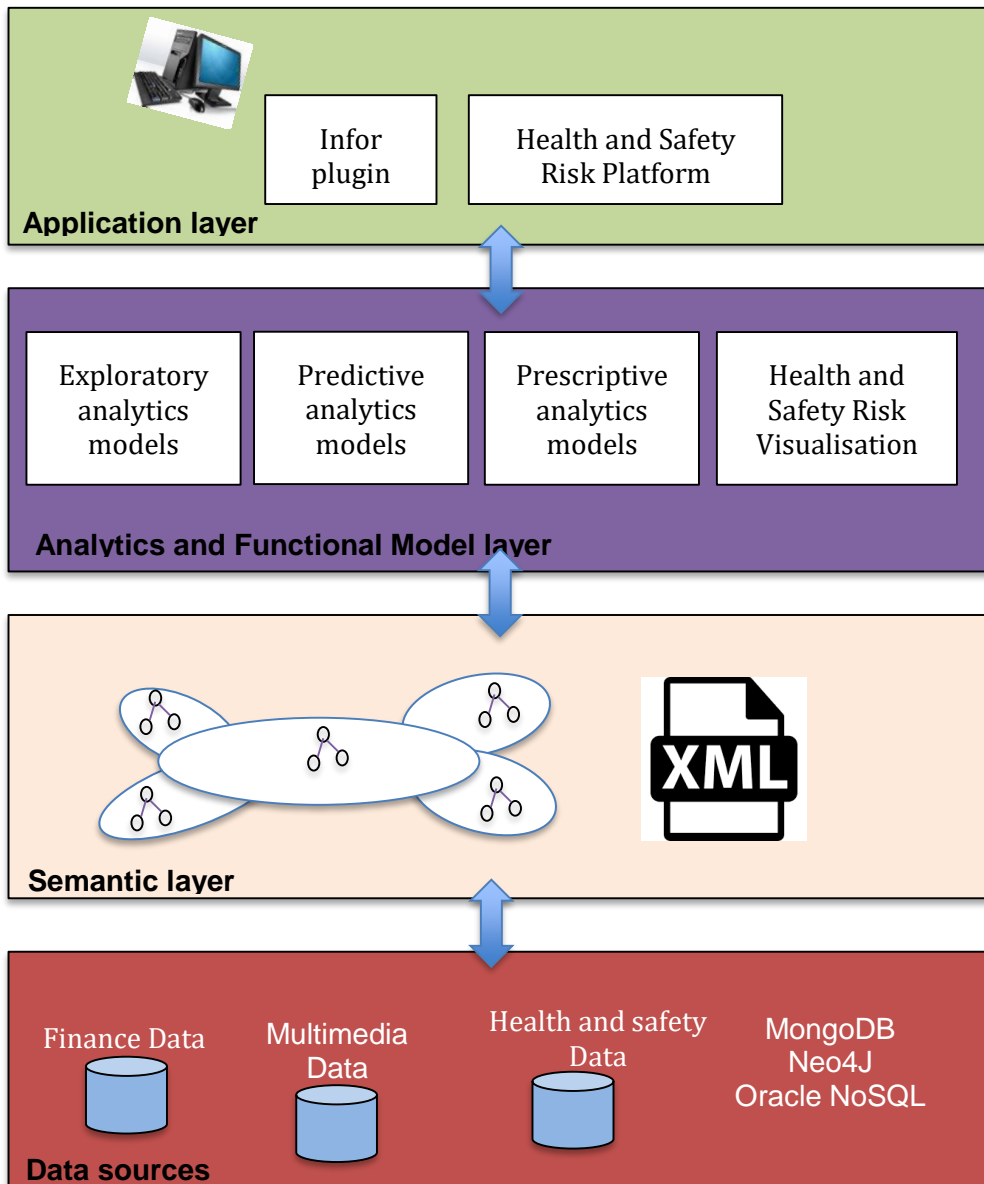
421

##### 422 4.3. *Analytics and functional model layer*

423 The significance of health and safety risk management tool lies in its ability to analyse and  
424 promptly act upon complex and high volume data. The layer has one functional model (Health  
425 and Safety visualisation) and three analytics models (discussed earlier), which are exploratory  
426 analytics, predictive analytics, and prescriptive analytics. As discussed earlier, predicting and  
427 managing health hazards is data-driven and highly intensive. Consequently, the Apache Spark  
428 engine was chosen over the MapReduce to build the analytics (predictive and prescriptive),  
429 due to its efficient in-memory storage and computation (Ryza et al. 2015). The analytical  
430 pipelines for health hazards management are actualised using SparkR, H2O, and GraphX.

431

432



433  
434  
435  
436

Figure 3. B-DAPP architecture

437 During each iteration in the analytical pipeline, different predictive models for health hazards  
438 are explored and optimised for optimum accuracy.

439 The H2O framework is selected because of its rich graphical user interface (GUI) and numerous  
440 tools for developing deep neural networks models. Additionally, it offers a comprehensive open  
441 source machine learning toolkit that is suitable for big data (Landset et al. 2015). It also provides  
442 tools for varied machine learning tasks, optimisation tools, data preprocessing and deep neural  
443 networks. Additionally, it offers coherent integration with Java, Python, R and R Studio, as well  
444 as Sparkling Water for integration with Spark and MLlib. Prior to or during an infrastructure  
445 project construction, health hazards are predicted and disseminated to stakeholders to help in  
446 mitigating the impact of hazards.

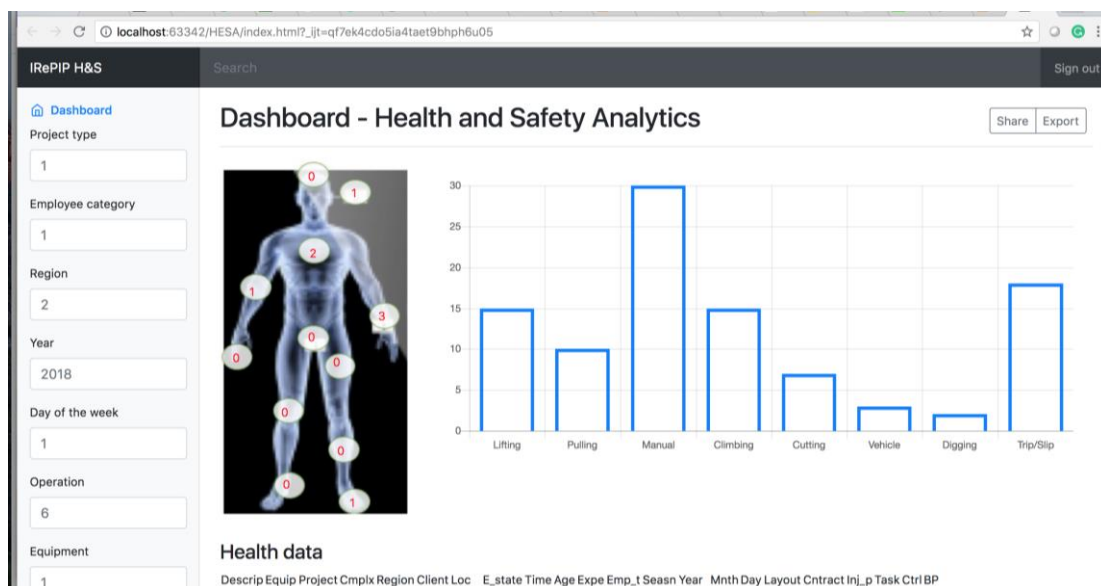
447 4.4. Application layer  
448 This layer is built by exploiting its powerful API programs. The end users of the tool are  
449 stakeholders (Engineers, Health and Safety officers, Site managers, Top level directors, etc.).  
450 The explanatory variables for infrastructure projects under B-DAPP are captured through  
451 appropriate the user interface and loaded to the HDFS and then to the Triplestore. Spark  
452 Streaming triggers the analytics pipeline to predict health hazards and suggests actionable  
453 insights to minimise health hazards. The predictions and prescriptions are communicated as  
454 the Predictive Model Markup Language (PMML). Stakeholders are provided with information to  
455 manage health hazards effectively.

456  
457  
458 **5. Results and discussions**

459 The prototype of the B-DAPP architecture is implemented by considering and interfacing the  
460 various technology artefacts. A sample screenshot produced by simulating the B-DAPP system  
461 is as shown in Figure 4, where the system predicts probable and number of injuries to body  
462 parts after the specification of input parameters (i.e. "Project type", "Region", "Operation", etc).  
463 It informs stakeholders of probable risks and allowing them adequate attention to risk factors  
464 when managing occupational hazards to achieve a safer environment.

465 The B-DAPP architecture is evaluated using exploratory data analysis and some preliminary  
466 results are provided. The purpose of this evaluation is to test the appropriateness of the B-  
467 DAPP architectural components and present some of these initial results. Interestingly, results  
468 obtained support findings in the literature. The future goal is to conduct a more rigorous  
469 evaluation through predictive analytics, by exploiting the preliminary analysis results presented  
470 in this paper.

471



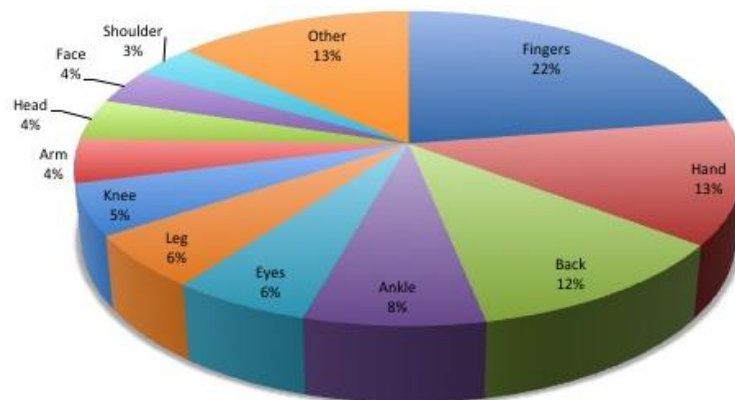
472  
473  
474  
475  
476

Figure 4. Screenshot of sub-module

477 5.1. Injury distribution by body parts

478 Since, Health and Safety dataset include the operation type variable, which describes the type  
479 of operation (lifting, pulling, cutting, etc.) with the specific tool (equipment) for the given task.  
480 Understanding the distribution of injury by body parts can highlight the top-k operations, for  
481 instance, that result in accidents to body parts. A graphical statistical tool (Pie chart) to explore  
482 this information is as depicted in Figure 5, where it is observed that certain body parts are prone  
483 to injuries during the power infrastructure project construction. The injury distribution of the top-  
484 5 body parts as specified in the database is as follows: Fingers (23%), Hand (13%),  
485 Back/Buttocks (12%), and Ankle (8%). The top five operations resulting in these injuries are  
486 pulling (stringing), lifting, loading/offloading, manual handling, and cutting because these parts  
487 are essential for carrying out these operations (Chi & Han 2013). The observation from this is  
488 probably that most of the accidents are as a result of carelessness, distractions, and disregard  
489 for safety procedures. The exploratory analysis results are in agreement with Fan et al. (2014).  
490 This fine-grained knowledge is not only integral to the development of robust construction  
491 health and safety risk management but also critical for stakeholders to enforce best safety  
492 practices to minimise accidents.

493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507



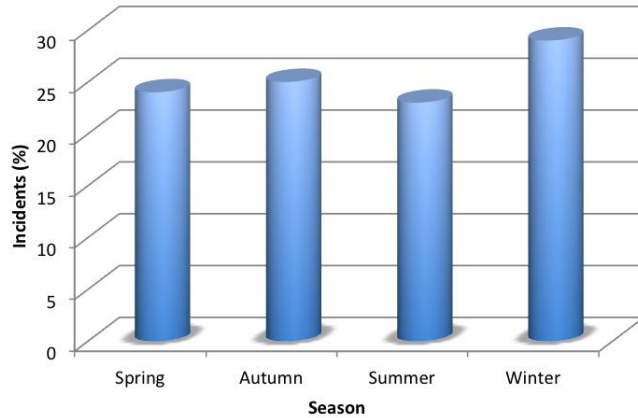
508 Figure 5. Injury distribution by body parts

509  
510 5.2. Incident distribution by season

511 Constructing power infrastructure (i.e. overhead lines) is mostly an outdoor activity, and certain  
512 types of accidents are more likely due to the changing seasonal conditions (summer, winter,  
513 autumn, and spring). Figure 6 shows that winter has the highest percentage of incidents (29%),  
514 followed by autumn (25%), spring (24%) and summer (23%). Scotland has a temperate and  
515 oceanic climate that is very cold in winter, due to frequent and heavy hail and snow showers.  
516 Wales likewise, has a temperate climate and tends to be wetter than England.  
517 Trips, slips, and falls are among the most common incidents in these regions due to the reduced  
518 visibility. Temperatures near or below freezing and strong winds can also result in severe illness  
519 and injury. Additionally, vehicle accidents occur due to the effects of ice and snow on muddy  
520 roads.

521 The use of Big Data analytics for automatic extraction and dissemination of climatic conditions  
522 of a region in real-time will go a long way at mitigating injuries that are synonymous to that  
523 region (location).

524



525

526 Figure 6. Incidents distribution by season

527

### 528 5.3. Accident distribution by spatial analysis

529 Often, the top management of a construction company may be interested in regions with high  
530 incident rates. Offering this service will equip managers with adequate information to  
531 proactively react to health and safety challenges in such regions. Thus, spatial analysis is of  
532 immense importance in such situations in that it enables the analysis of incidents over the  
533 topological and geographical spread. In the health and safety dataset, the location information  
534 is captured in the 'site' column. For the spatial analysis, the dataset is pre-processed to extract  
535 the UK postcode of each incident record and linked with the corresponding latitude and  
536 longitude data from Doogal (<http://www.doogal.co.uk/UKPostcodes.php>). The geographical  
537 heat map is employed to visualise the resulting data. Figure 7 shows the summary of this  
538 distribution, where the size of spheres represents the proportion of accidents (computed as  
539 percentages) in each region. Scotland has the highest (30%), followed by Wales and South  
540 West (25%), North (16%), South East (14%), and Midlands (2%). The frequency of severe  
541 weather is observed to be the leading cause of accidents in Scotland as well as Wales and  
542 South West regions. Strong wind, for instance, may lead to shattering of vehicle windscreens  
543 and a collapse of a fence or unit. Icy weather may result in trips and slips. Also, heavy-duty  
544 machinery operation (i.e. excavation and road cutting) is often the cause of utility service  
545 damage (i.e. gas pipelines, water supply). Even though geological conditions in different cities  
546 are complex, existing health and safety risk management approaches do not consider making  
547 this information available for proper health and safety risk prevention. To efficiently bring health  
548 and safety risk in the site under control, incorporating a module to automatically compute the  
549 geology and hydrology condition of construction sites in real-time will improve the optimal  
550 control of occupational hazards.

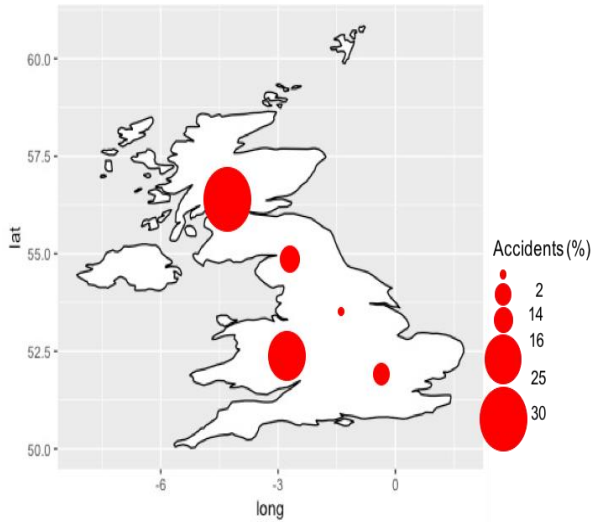
551

552

553

554





556

557

Figure 7. Spatial analysis of accidents

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

Additionally, the result of viewing the regions with respect to incident (or accident) rate can further be narrowed to cities and a specific location. The impact of location on incidents is worth further exploration. This investigation is the focus of future research on the proposed architecture.

#### 5.4. Modelling the relationship between variables

Tremendous R&D efforts have been carried out to reduce the impacts of occupational health hazards. One such attempt is in modelling and analysing several variables (i.e. determining the relationships between the predictors (independent variables) and the dependent variable. Robust and efficient machine learning techniques such as deep learning, gradient boosting machines, and linear multivariate regression are employed in modelling relationships among variables. In this paper, a demonstration of the linear regression technique is made due to its simplicity.

Linear multivariate regression, in this regard, advocates methods for analysing health hazards with respect to the project cost. This concept not only enables the exploratory analysis of injury but also allows predictive accident analytics. The principle of the linear multivariate regression is to predict  $Y$  as a linear combination of the input variables  $(x_1, x_2, \dots, x_p)$  plus an error term  $\epsilon_i$ .

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i, i \in [1, n]$$

$n$  is the number of sample data,  $p$  the number of variables and  $\beta_0$  a bias. This model can conveniently be written as  $y = X\beta + \epsilon$ , where

$$y = (y_1, \dots, y_n)^T, \epsilon = (\epsilon_1, \dots, \epsilon_n)^T, \beta = (\beta_1, \dots, \beta_n)^T, \text{ and } X = \begin{pmatrix} 1 & x_{11} & \vdots & x_{1p} \\ 1 & x_{21} & \vdots & x_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & \vdots & x_{np} \end{pmatrix}$$

The predicted or fitted value is thus,  $\hat{y} = X\hat{\beta}$ , where  $\hat{\beta}$  is the least squares estimate of  $\beta$ .

584 The model can be used for example to predict the body part injured given a set of inputs such  
 585 as the type of operation (task), equipment being used, kind of power infrastructure project, the  
 586 project complexity, project contract type, etc. A practical but straightforward illustration is to  
 587 determine the relationship between the project cost and occupational hazards (linear  
 588 regression with one predictor) is depicted using a line plot (Figure 8). The x-axis of the plot  
 589 represents the project cost while the y-axis represents the health hazards risk (incidents and  
 590 accidents). The line plot shows a significant increase in the number of health hazards (accident  
 591 and incidents) as the project cost increases. Consequently, the number of occupational health  
 592 risk is proportional to the project cost. This result is expected since the project cost is a crucial  
 593 factor in determining the complexity of a project. Thus, the more complex a project is, the more  
 594 are incidents associated with it.

595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622

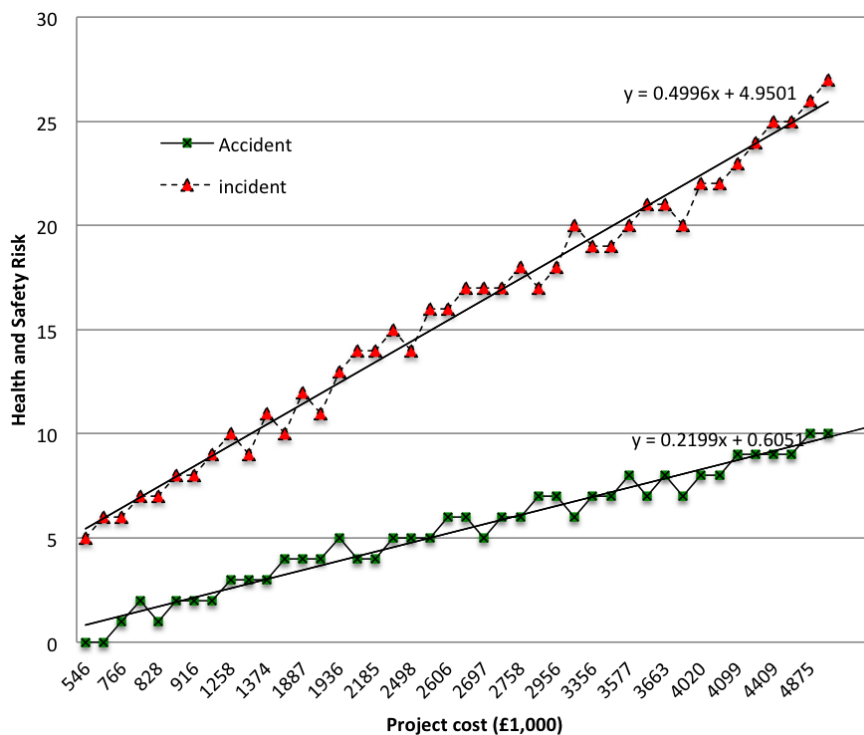


Figure 8. Relationship among variables

## 6. Conclusions

623 Construction safety risk analyses are currently limited because existing techniques overlook  
 624 the complex and dynamic nature of construction sites. Besides, they ignore or simplify some  
 625 key factors and pay little attention to analysing the relationship between a safety phenomenon  
 626 and safety data. Today, large and dynamic data with various data types are to be analysed. In  
 627 implementing the health hazards management tool, the Big Data architecture that is based on  
 628 a well coherent health risk analytics lifecycle is proposed. The Big Data technology was  
 629 selected due to its support for massive, high dimensional, heterogeneous, complex,  
 630 unstructured, incomplete, and noisy data.

631

632 The preliminary results obtained in this study using the various Big Data frameworks have  
633 enabled us to design a robust architecture to handle and analyse power infrastructure accident  
634 data. The proposed architecture can identify relevant variables and improve preliminary  
635 prediction accuracies and explanatory capacities. It has also enabled conclusions to be drawn  
636 regarding the causes of health hazards. The results obtained in this study represent a  
637 significant improvement in terms of managing information on construction accidents,  
638 particularly for power infrastructure companies. The satisfactory results of the B-DAPP tool  
639 have indicated the reliability and appropriateness of the selected Big Data components for  
640 studies of construction health risks and their causes.

641

642 Future research is aimed at rigorously evaluating accuracies of both the prediction and  
643 prescription of the software deployed in real-time. Additionally, other researchers should look  
644 in the area of designing and planning a more ambitious, larger scale models to gain a deeper  
645 understanding of accident causes in various industrial sectors.

646

647

#### 648 References

- 649 Al-Jarrah, O.Y. et al., 2015. Efficient Machine Learning for Big Data: A Review. *Big Data*  
650 *Research*, 2(3), pp.87–93.
- 651 Bohle, P. et al., 2015. Health and well-being of older workers: Comparing their associations  
652 with effort–reward imbalance and pressure, disorganisation and regulatory failure. *Work*  
653 *& Stress*, 8373, pp.1–14.
- 654 Camann, D.E. et al., 2011. *Data Science and Big Data Analytics*, New York: EMC.
- 655 Carbonari, A., Giretti, A. & Naticchia, B., 2011. A proactive system for real-time safety  
656 management in construction sites. *Automation in Construction*, 20(6), pp.686–698.  
657 Available at: <http://dx.doi.org/10.1016/j.autcon.2011.04.019>.
- 658 Chaudhuri, S. & Dayal, U., 1997. An overview of data warehousing and OLAP technology.  
659 *ACM SIGMOD Rec*, 26, pp.65–74.
- 660 Chen, J., Qiu, J. & Ahn, C., 2017. Construction worker's awkward posture recognition through  
661 supervised motion tensor decomposition. *Automation in Construction*, 77, pp.67–81.
- 662 Cheng, C. et al., 2011. Applying data mining techniques to explore factors contributing to  
663 occupational injuries in Taiwan's construction industry. *Accident Analysis and*  
664 *Prevention*, 48, pp.214–222.
- 665 Chi, S. & Han, S., 2013. Analyses of systems theory for construction accident prevention with  
666 specific reference to OSHA accident reports. *International Journal of Project*  
667 *Management*, 31(7), pp.1027–1041.
- 668 Ciarapica, F.E. & Giacchetta, G., 2009. Classification and prediction of occupational injury  
669 risk using soft computing techniques : An Italian study. *Safety Science*, 47(1), pp.36–49.  
670 Available at: <http://dx.doi.org/10.1016/j.ssci.2008.01.006>.
- 671 Debnath, J. et al., 2016. Fuzzy inference model for assessing occupational risks in  
672 construction sites. *International Journal of Industrial Ergonomics*, 55, pp.114–128.
- 673 Delen, D. & Demirkan, H., 2013. Data, information and analytics as services. *Decision*  
674 *Support Systems*, 55(1), pp.359–363.
- 675 Esmaili, B., Hallowell, M.R. & Rajagopalan, B., 2015. Attribute-based safety risk  
676 assessment. II: Predicting safety outcomes using generalized linear models. *Journal of*  
677 *Construction Engineering and Management*, 141(8), pp.1–11.
- 678 Fan, Z.J. et al., 2014. The association between combination of hand force and forearm  
679 posture and incidence of lateral epicondylitis in a working population. *Human factors*,  
680 56, pp.151–165.
- 681 Favaro, F.M. & Saleh, J.H., 2016. Toward risk assessment 2.0: Safety supervisory control  
682 and model-based hazard monitoring for risk-informed safety interventions. *Reliability*  
683 *Engineering and System Safety*, 152, pp.316–330. Available at:  
684 <http://dx.doi.org/10.1016/j.ress.2016.03.022>.

685 Fenrick, L. & Getachew, S., 2012. Cost and reliability comparisons of underground and  
686 overhead power lines. *Utilities Policy*, 20(1), pp.31–37.

687 Fragiadakis, N., Tsoukalas, V. & Papazoglou, V., 2014. An adaptive neuro-fuzzy inference  
688 system (anfis) model for assessing occupational risk in the shipbuilding industry. *Safety*  
689 *Science*, 63, pp.226–235.

690 Galizzi, M. & Tempesti, T., 2015. Workers' risk tolerance and occupational injuries. *Risk*  
691 *Analysis*, 35(10), pp.1858–1875.

692 Gandomi, M. & Haider, A., 2015. Beyond the hype: Big data concepts, methods, and  
693 analytics. *International Journal of Information Management*, 35(2), pp.137–144.

694 García-Herrero, S. et al., 2012. Working conditions, psychological/physical symptoms and  
695 occupational accidents Bayesian network models. *Safety Science*, 50(9), pp.1760–  
696 1774.

697 Gholizadeh, P. & Esmaeili, B., 2016. Applying classification trees to analyze electrical  
698 contractors' accidents. In *Construction Research Congress*. San Juan, Puerto Rico, pp.  
699 2699–2708.

700 Groves, W., Kecejovic, V. & Komljenovic, D., 2007. Analysis of fatalities and injuries involving  
701 mining equipment. *Journal of Safety Research*, 38(4), pp.461–470.

702 Guo, B., Yiu, T. & González, V., 2016. Predicting safety behavior in the construction industry:  
703 Development and test of an integrative model. *Safety Science*, 84, pp.1–11. Available  
704 at: <http://dx.doi.org/10.1016/j.ssci.2015.11.020>.

705 Guo, S. et al., 2016. A Big-Data-based platform of workers' behavior: Observations from the  
706 field. *Accident Analysis and Prevention*, 93, pp.299–309. Available at:  
707 <http://dx.doi.org/10.1016/j.aap.2015.09.024>.

708 Gürçanlı, G. & Müngena, U., 2009. An occupational safety risk analysis method at  
709 construction sites using fuzzy sets. *International Journal of Industrial Ergonomics*, 39(2),  
710 pp.371–387.

711 Hallowell, M.R. & Gambatese, J.A., 2009. Construction Safety Risk Mitigation. *Journal of*  
712 *Construction Engineering and Management*, 135(12), pp.1316–1323. Available at:  
713 <http://ascelibrary.org/doi/10.1061/%28ASCE%29CO.1943-7862.0000107>.

714 Haslam, R.A. et al., 2005. Contributing factors in construction accidents. *Applied Ergonomics*,  
715 36(4), p.401–415.

716 HSE, 2016. Health and safety at work Summary statistics for Great Britain. Available at:  
717 <http://www.hse.gov.uk/statistics/overall/hssh1516.pdf?pdf=hssh1516> [Accessed April  
718 14, 2017].

719 Jacinto, C. & Silva, C., 2010. A semi-quantitative assessment of occupational risks using  
720 bow-tie representation. *Safety Science*, 48, pp.973–979.

721 Jin, X. et al., 2015. Significance and Challenges of Big Data Research. *Big Data Research*,  
722 2(2), pp.59–64.

723 Jocelyn, S. et al., 2017. Application of logical analysis of data to machinery-related accident  
724 prevention based on scarce data. *Reliability Engineering and System Safety*, 159(May  
725 2016), pp.223–236. Available at: <http://dx.doi.org/10.1016/j.ress.2016.11.015>.

726 Khakzad, N., Khan, F. & Amyotte, P., 2015. Major Accidents (Gray Swans) Likelihood  
727 Modeling Using Accident Precursors and Approximate Reasoning. *Risk Analysis*, 35(7),  
728 pp.1336–1347.

729 Landset, S. et al., 2015. A survey of open source tools for machine learning with big data in  
730 the Hadoop ecosystem. *Journal of Big Data*, 2(1), pp.1–36.

731 Leavitt, N., 2010. Will NoSQL Databases Live Up to Their Promise? *IEE Computer Journal*,  
732 43(2), pp.12–14.

733 Li, H. et al., 2016. Stochastic state sequence model to predict construction site safety states  
734 through real-time location systems. *Safety Science*, 84, pp.78–87.

735 Li, Y. & Bai, Y., 2008. Comparison of characteristics between fatal and injury accidents in the  
736 highway construction zones. *Safety Science*, 46(4), pp.646–660.

737 Liao, C.-W. & Perng, Y.-H., 2008. Data mining for occupational injuries in the Taiwan  
738 construction industry. *Safety Science*, 46(7), pp.1091–1102.

739 Liu, H. & Tsai, Y., 2012. A fuzzy risk assessment approach for occupational hazards in the  
740 construction industry. *Safety Science*, 50(4), pp.1067–1078. Available at:  
741 <http://dx.doi.org/10.1016/j.ssci.2011.11.021>.

742 Love, P.E.D. & Teo, P., 2017. Statistical Analysis of Injury and Nonconformance Frequencies  
743 in Construction : Negative Binomial Regression Model. *Journal of Construction*  
744 *Engineering and Management*, 143(8), pp.1–9.

745 McDermott, V. & Hayes, J., 2016. "We're still hitting things": The effectiveness of third party  
746 processes for pipeline strike prevention. In *Proceedings of the eleventh international*  
747 *pipeline conference (IPC 2016)*. Calgary, Alberta, Canada, pp. 1–10.

748 Naderpour, M., Lu, J. & Zhang, G., 2016. A safety-critical decision support system evaluation  
749 using situation awareness and workload measures. *Reliability Engineering and System*  
750 *Safety*, 150, pp.147–159. Available at: <http://dx.doi.org/10.1016/j.ress.2016.01.024>.

751 Najafabadi, M.M. et al., 2015. Deep learning applications and challenges in big data analytics.  
752 *Journal of Big Data*, 2(1), pp.1–21.

753 Nanda, G. et al., 2016. Bayesian decision support for coding occupational injury data. *Journal*  
754 *of Safety Research*, 57, pp.71–82. Available at:  
755 <http://dx.doi.org/10.1016/j.jsr.2016.03.001>.

756 Pääkkönen, P. & Pakkala, D., 2015. Reference Architecture and Classification of  
757 Technologies, Products and Services for Big Data Systems. *Big Data Research*, 2(4),  
758 pp.166–186. Available at: <http://dx.doi.org/10.1016/j.bdr.2015.01.001>.

759 Papazoglou, I. et al., 2017. Quantitative occupational risk model: Single hazard. *Reliability*  
760 *Engineering & System Safety*, 160, pp.162–173.

761 Papazoglou, I. et al., 2015. Uncertainty Assessment in the Quantification of Risk Rates of  
762 Occupational Accidents. *Risk Analysis*, 35(8), pp.1536–1561.

763 Pinto, A., Nunes, I. & Ribeiro, R., 2011. Occupational risk assessment in construction industry  
764 – Overview and reflection. *Safety Science*, 49, pp.616–624.

765 Power, D., 2014. Using "Big Data" for analytics and decision support. *Journal of Decision*  
766 *Systems*, 23(2), pp.222–228.

767 Rahman, M.N. & Esmailpour, A., 2016. A Hybrid Data Center Architecture for Big Data. *Big*  
768 *Data Research*, 3, pp.29–40.

769 Raviv, G., Shapira, A. & Fishbain, B., 2017. AHP-based analysis of the risk potential of safety  
770 incidents: Case study of cranes in the construction industry. *Safety Science*, 91,  
771 pp.298–309. Available at: <http://dx.doi.org/10.1016/j.ssci.2016.08.027>.

772 Rivas, T. et al., 2011. Explaining and predicting workplace accidents using data-mining  
773 techniques. *Reliability Engineering & System Safety*, 96(7), pp.739–747.

774 Ryza, OJ S. et al., 2015. *Advanced Analytics with Spark*, Cambridge: O'Reilly,.

775 Schryver, J., Shankar, M. & Xu, S., 2012. Moving from descriptive to causal analytics: Case  
776 study of discovering knowledge from US health indicators warehouse. In *ACM SIGKDD*  
777 *Workshop on Health Informatics*. Beijing, China, pp. 1–8.

778 Silva, S.A. et al., 2016. Organizational practices for learning with work accidents throughout  
779 their information cycle. *Safety Science*, In Press.

780 Soltanzadeh, A. et al., 2016. Analysis of occupational accidents induced human injuries: A  
781 case study in construction industries and sites. *Journal of Civil Engineering and*  
782 *Construction Technology*, 7(1), pp.1–7. Available at:  
783 <http://academicjournals.org/journal/JCECT/article-abstract/15EEFC357741>.

784 Suthakar, U. et al., 2016. An efficient strategy for the collection and storage of large volumes  
785 of data for computation. *Journal of Big Data*, 3(1), pp.1–17.

786 Tixier, A.. et al., 2016. Application of machine learning to construction injury prediction.  
787 *Automation in Construction*, 69, pp.102–114.

788 Törner, M. & Pousette, A., 2009. Safety in construction - a comprehensive description of the  
789 characteristics of high safety standards in construction work, from the combined  
790 perspective of supervisors and experienced workers. *Journal of Safety Research*, 40(6),  
791 pp.399–409.

792 Tsai, C.W. et al., 2015. Big data analytics: a survey. *Journal of Big Data*, 2(1), pp.1–32.

793 Venturini, L., Baralis, E. & Garza, P., 2017. Scaling associative classification for very large  
794 datasets. *Journal of Big Data*, 4(1). Available at: [https://doi.org/10.1186/s40537-017-](https://doi.org/10.1186/s40537-017-0107-2)  
795 0107-2.

796 Weng, J., Meng, Q. & Wang, D.Z.W., 2013. Tree-based logistic regression approach for work  
797 zone casualty risk assessment. *Risk Analysis*, 33(3), pp.493–504.

798 White, T., 2012. *Hadoop: The Definitive Guide*, Sebastopol, CA: O'Reilly Media, Inc.

799 Wu, W. et al., 2010. Towards an autonomous real-time tracking system of near-miss  
800 accidents on construction sites. *Automation in Construction*, 19(2), pp.134–141.  
801 Available at: <http://dx.doi.org/10.1016/j.autcon.2009.11.017>.

802 Yi, W. et al., 2016. Development of an early-warning system for site work in hot and humid  
803 environments: A case study. *Automation in Construction*, 62, pp.101–113. Available at:  
804 <http://dx.doi.org/10.1016/j.autcon.2015.11.003>.

805 Yoon, Y.S., Ham, D.H. & Yoon, W.C., 2016. Application of activity theory to analysis of  
806 human-related accidents: Method and case studies. *Reliability Engineering and System*  
807 *Safety*, 150, pp.22–34. Available at: <http://dx.doi.org/10.1016/j.res.2016.01.013>.  
808 Yorio, P.L., Willmer, D.R. & Haight, J.M., 2014. Interpreting MSHA citations through the lens  
809 of occupational health and safety management systems: Investigating their impact on  
810 mine injuries and illnesses 2003-2010. *Risk Analysis*, 34(8), pp.1538–1553.  
811 Zang, W. et al., 2014. Comparative study between incremental and ensemble learning on  
812 data streams: Case study. *Journal Of Big Data*, pp.1–16. Available at:  
813 <http://www.journalofbigdata.com/content/1/1/5/abstract>.  
814 Zeng, S.X., Tam, V.W.Y. & Tam, C.M., 2008. Towards occupational health and safety  
815 systems in the construction industry of China. *Safety Science*, 46, pp.1155–1168.  
816 Zhang, L. et al., 2016. Towards a Fuzzy Bayesian Network Based Approach for Safety Risk  
817 Analysis of Tunnel-Induced Pipeline Damage. *Risk Analysis*, 36(2), pp.278–301.  
818 Zhou, Z., Goh, Y. & Li, Q., 2015. Overview and analysis of safety management studies in the  
819 construction industry. *Safety Science*, 72, pp.337–350.  
820 Zhu, Z. et al., 2016. Predicting movements of onsite workers and mobile equipment for  
821 enhancing construction site safety. *Automation in Construction*, 68, pp.95–101.  
822 Zou, P.X.W., Zhang, G. & Wang, J., 2007. Understanding the key risks in construction  
823 projects in China. *International Journal of Project Management*, 25(6), pp.601–614.  
824