

DATA ANALYTICS ASSESSMENT: ANALYSE A DATASET

By

Aaron Ward

Supervisor(s):

SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE OF
B.SC IN COMPUTING AND INFORMATION TECHNOLOGY
AT
INSTITUTE OF TECHNOLOGY BLANCHARDSTOWN
DUBLIN, IRELAND
2017

Declaration

I/We declare that this material, which I/We now submit for assessment, is entirely my/our own work and has not been taken from the work of others, except where otherwise stated. I/We have identified and included the source of all facts, ideas, opinions, and viewpoints of others in the assignment references. Direct quotations from books, journal articles, internet sources, module text, or any other source whatsoever are acknowledged and the source cited are identified in the assignment references.

I/We understand that plagiarism, collusion, and copying are grave and serious offences and accept the penalties that would be imposed should I/we engage in plagiarism, collusion or copying. I acknowledge that copying someone else's assignment, or part of it, is wrong, and that submitting identical work to others constitutes a form of plagiarism. I/We have read and understood the college's plagiarism policy 3AS08.

This material, or any part of it, has not been previously submitted for assessment for an academic purpose at this or any other academic institution. I have not allowed anyone to copy my work with the intention of passing it off as their own work.

Dated: 2017

Author:

Aaron Ward

1	Numeric Attribute Description.	4
2	Categorical Attribute Description.	5

Business Understanding and Data Understanding

Business Understanding

This dataset provides data on subjects with and without meningitis. It contains information such as age, gender, location, sum of health problems such as headaches, fevers and seizures. Additionally, it provides an attribute that indicates if the subject does or doesn't have meningitis, with a negative or positive value.

Business Objective

- Predict if someone is at risk of getting meningitis.

Data Mining objective

- The main objective is to create a model to predict the risk of a person getting meningitis
- This model will use the attributes provided in the dataset such as age, gender, seizure history etc to assess the prediction accuracy.
- The model shall test multiple data mining algorithms to obtain a prediction.

Data Understanding

Describing Data

In this section the allocated dataset is explained in terms of informational content, data quality and usability. The data set itself consists of attributes in relation to meningitis, a neurological infectious disease that can cause brain inflammation due to bacteria or viruses infecting that brain. As seen in tables below, the attributes have been segregated from numeric and categorical data. The numeric data given a description and a data type. Additionally, the mean, minimum, maximum and standard deviation values are given. Please refer to table 1 for information of the numeric data.

Numeric Attributes						
Name	Description	Data type	Mean	Min	Max	SD
AGE	List the age of each person	Numeric	37.6285	10.0	84.0	15.3853
COLD	Number of days since last cold	Numeric	2.6642	0.0	35.0	4.8273
HEADACHE	Days since last headache	Numeric	7.1857	0.0	63.0	9.1278
FEVER	Days since last fevers	Numeric	6.3428	0.0	63.0	8.0294
NAUSEA	Start of nausea	Numeric	2.4857	0.0	32.0	4.5856
LOC	When loss of consciousness occurs	Numeric	0.7428	0.0	26.0	2.6481

SEIZURE	When convulsions are observed	Numeric	0.1857	0.0	6.0	0.8780
BT	Body temperature	Numeric	37.625	35.5	40.2	1.3041
STIFF	Neck stiffness	Numeric	1.9571	0.0	5.0	1.4033
KERNIG	Kernig sign	Numeric	0.2142	0.0	1.0	0.4117
LASEGUE	Lasegue sign	Numeric	0.0785	0.0	1.0	0.2700
GCS	Glasgow coma scale	Numeric	14.7071	9.0	15.0	1.1536
WBC	White blood cell count	Numeric	8743.42	1070	90009	7795.80
CRP	C-Reactive protein	Numeric	1.6878	0.0	31.0	4.1317
ESR	Blood sedimentation test	Numeric	5.9285	0.0	60.0	11.880
CSF_CELL	Cell Count in Cerebulospinal Fluid	Numeric	1505.4	0.0	63350	5708.83
Cell_Poly	Polynuclear cell in CSF	Numeric	1025.85	0.0	61520	5402.38
Cell_Mono	Mononuclear cell in CSF	Numeric	465.08	0.0	7840	816.98
CSF_PRO	Protein in CSF	Numeric	99.414	0.0	474.0	96.307
CSF_GLU	Glucose in CSF	Numeric	56.578	0.0	520	44.3412
CSF_CELL3	Cell Count CSF 3 days after the treatment	Numeric	385.18	8	4860	1038.37

CSF_CELL7	Cell Count of CSF 7 days after treatment	Numeric	205.61	0.0	7840	816.98
-----------	--	---------	--------	-----	------	--------

Table 1: Numeric Attribute Description.

The categorical dataset is given a description to the attribute labels and given a data type. Most of the attributes consist of only 2 values, but does of whom that are multivalued are displayed with the highest and lowest values in the table. See table 2 for further insight to the dataset.

Categorical Attributes				
Name	Description	Data type	Value 1	Value 2
SEX	Gender of people	Nominal	M (82)	F (58)
Diag2	Diagnoses	Nominal	VIRUS (98)	BACTERIA (42)
ONSET	Inception	Nominal	CHRONIC (1)	ACUTE (130)
LOC_DAT	Loss of consciousness	Nominal	- (98)	+ (42)
FOCAL	Focal Sign	Nominal	- (105)	+ (35)
CT_FIND	CT findings	Nominal	normal (101)	abnormal (39)
EEG_WAVE	Electroencephalography Wave Findings	Nominal	abnormal (117)	normal (23)
EEG_FOCUS	Focal sign in EEG	Nominal	-(104)	+(36)
CULT_FIND	If bacteria or virus found	Nominal	F (107)	T(33)
CULTURE	Name of bacteria/virus found	Nominal	Tb (1)	- (107)
THERAPY2	Therapy	Nominal	PIPC+CTX (1)	no_therapy (58)
C_COURSE	Clinical course at discharge	Nominal	negative (117)	paralysis (1)
COURSE(Grouped)	Grouped attribute of C_COURSE	Nominal	n (117)	p (23)

RISK(Grouped)	Class label - at risk	Nominal	n (121)	p (19)
---------------	-----------------------	---------	---------	--------

Table 2: Categorical Attribute Description.

The data above could be divided into a number of sections. Attributes such as AGE and SEX can be categorized as *personal information*. COLD, HEADACHE, NAUSEA LOC etc. can be described as *subject history* as they provide some information on the commencement of the symptoms. BT, STIFF, KERNIG, GCS can be assigned to a category of *physical examination* as they attribute values obtained during investigation. Further more, *laboratory investigation* used to describe the attributes such as WBC, EEG_WAVE, CULT_FIND, ESR etc. These are values collected during further investigate of the bodily anomalies. Lastly, *postliminary treatment* if used to describe THERAPY2, CSF_CELL3 and CSF_CELL7 as they are attributes describing values after a subject has been treated for meningitis. The dimensionality of the dataset is a total of 36 attributes and there are 141 instances.

Data Exploration

The following section will utilize the Exploratory Data Analysis (EDA) techniques to further analyze the data. This will give insight on how to improve the quality of the data and reduce dimensionality among the dataset. It shall also provide a number of charts to visually represent the data in regards to attributes that may prove to be predictive, show indication of correlations, and data with low variability.

As seen in Figure1 a histogram is used to plot out the main symptoms of meningitis. This consists of the attributes COLD, HEADACHE, FEVER, NAUSEA. The values are attributed to the amount of days since these sensations have been felt. Therefore they are applicable to the mining objectives as they appear to be very predictive.

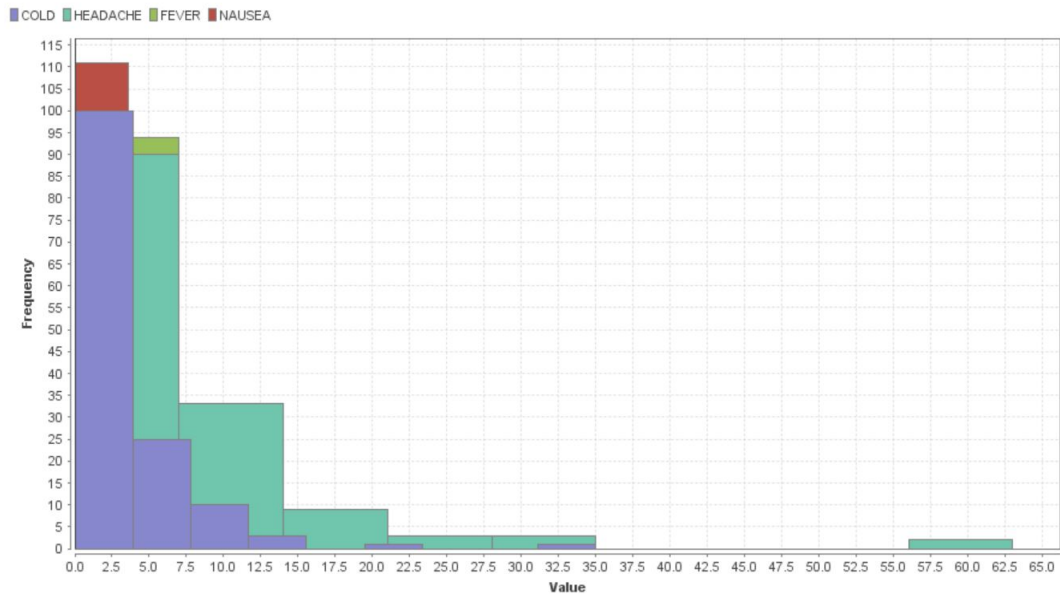


Figure 1: Main symptoms for meningitis

Furthermore, a high correlation between the attributes ONSET and CSF_CELL3 can be seen in Figure 2. Due to this high correlation one of these attributes can be removed from the dataset as they are redundant. Furthermore, They are scarce in values so they may not prove beneficial to the mining objectives.

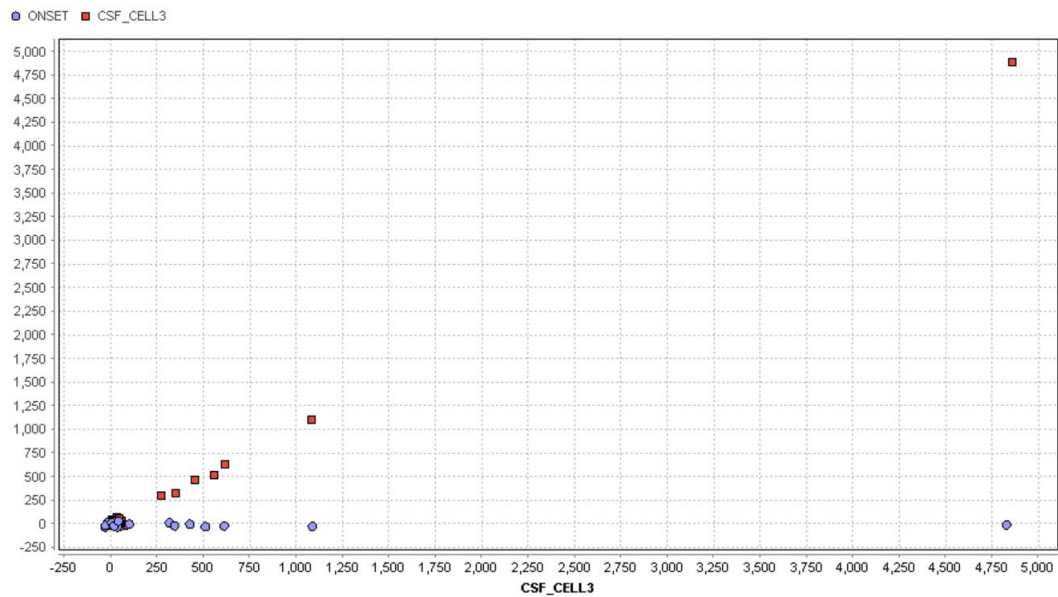


Figure 2: Scatter plot for ONSET and CSF_CELL3

The following histogram displays the values for positive (P) and negative (N) in the class label RISK(Grouped). This unbalance in inequitable group of values may prove to cause some difficulties as it may falsely predict everyone to be negative(n) for risk of having meningitis. See Figure 3 below.

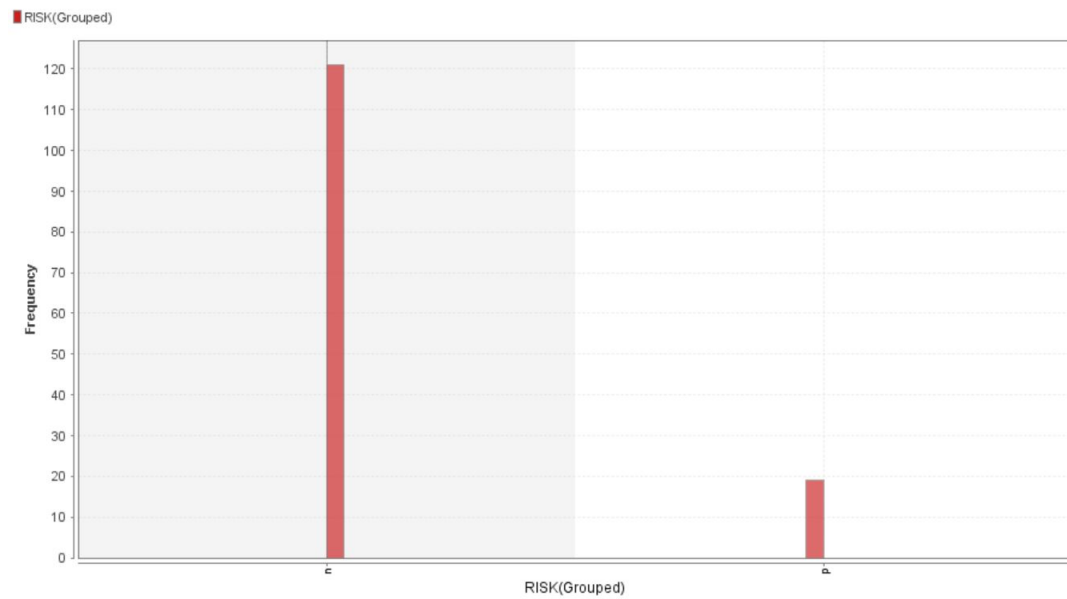


Figure 3: Histogram for class label RISK(Grouped)

Figure 4 shows a box plot of CSF_PRO the distribution of data. Although these values may prove beneficial in distinguishing between at risk and not at risk for meningitis, there are some outlier values seen which can skew the data.

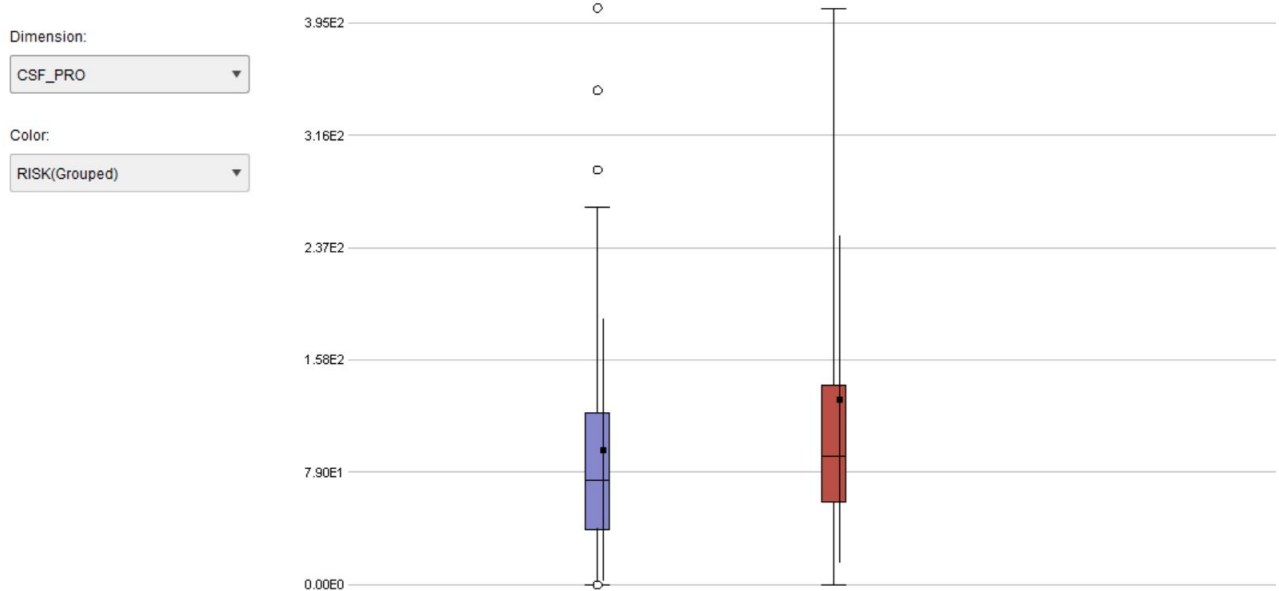


Figure 4: Protein in Cerebulospinal Fluid

Verifying Data Quality

Missing Values

Some missing values can be found in this data set. For example, row 122 contains no values at all. Additionally, CSF_CELL3 has a lot of missing data. There are a total of 119 out of the 141 instances of CSF_CELL3 that a missing values. This means that there are only 22 rows in the entire dataset that do not contain any missing data.

Noise, Bias and Outlier Values

Many outlier values can be seen. For instance, the FEVER attribute has a mean value of 6.343. But this value may be skewed as there is a high outlier value, 63, at row 8. Furthermore, the value 90009.0 in the white blood cell count (WBC). This may be an erroneous value as the next highest value in the column is only 19500. This outlier value may disrupt the accuracy of prediction in the model.

Sufficient Attributes

There are some redundant attributes found in the dataset. One example of this is the SEIZURE

attribute. There are only 8 individual and 4 distinct instances therefore it would be insufficient for representation of a class. Secondly, the CULT_FIND and CULTURE represent the same data. CULTURE would be more sufficient for prediction as it contain more varied values of bacteria/viruses found such as influenza, tb, herpes, strepto etc. rather than the true or false values of CULT_FIND. Therefore CULT_FIND can be removed from the dataset.

Data Preparation

Firstly, the model was tested using a decision tree without any data preprocessing. Firstly, the attribute RISK(Grouped) was assigned as the class label. As many rows contain missing values, a filter was applied to ignore this row of data using the *no_missing_attribute* condition class. The split validation operator was then used for splitting the dataset in a training and testing data set with a split ratio of 0.7. Initially, linear sampling was used for testing the model. This gave an accuracy result of 57.14%. Shuffled sampling was used, which randomly selects 70% of the rows in the dataset as training data. This returned a result of 85.71% percent. Lastly, stratified sampling was used to select 70% of the rows from each class and this gave a result of 83.33%. Secondly, the k-Nearest Neighbor (k-NN) algorithm was used. Linear sampling and shuffled sampling with k-NN both provided an accuracy of 57.14% and stratified sampling gave an accuracy of 66.67%.

Although these accuracies are relatively high, they are only trained on 22 rows of values. This provides a false accuracy as the dataset is not large enough. Therefore, further preprocessing measures are needed to give a fairer accuracy result.

Select Data

There is not a high number rows need to be removed from the dataset as there should be at least 20 times more rows than columns in a dataset. However, some of the rows or columns may be transformed due to missing values, which shall be addressed in the following section.

Clean Data

Due to missing data in a number of cells in the dataset a number of transformational steps must be taken to ensure an increase in the overall accuracy of the final model. As row 122 is the only row in the dataset that has a large amount of missing data, it shall be removed from the dataset as it makes up less than 5% of the number of rows, therefore it will not significantly affect the final result. This was implemented using the *Filter Example Range*, which was used to filter out row 122 from the dataset. Using stratified sampling with a decision tree, this produced an accuracy of 73.81%. Please refer to Figure 5

accuracy: 73.81%

	true n	true p	class precision
pred. n	27	2	93.10%
pred. p	9	4	30.77%
class recall	75.00%	66.67%	

Figure 5: Stratified Sampling Using a Decision Tree

The attribute CSF_CELL3 lacks a lot of values. There are a total of 119 rows out of 141 that contain no value at all which means there is an 84.4% value absence. This column shall be removed as it is missing more 40% of its values. This is applied using the *Select Attributes* Filter and exclude CSF_CELL3 from the dataset. There are other attribute that should also be removed. For example, attributes such as LASUGUE, and GCS should be removed because to provide too little variance in values to the dataset. Following the application of these steps, the accuracy of the model has risen from 73.81% to 78.57%, as seen in Figure 6.

accuracy: 78.57%

	true n	true p	class precision
pred. n	31	4	88.57%
pred. p	5	2	28.57%
class recall	86.11%	33.33%	

Figure 6: Results After Filtering Useless Attributes

More ever, WBC possibly has a very high value outlier value of 90009. To accommodate for this *outlier detection* is used to

Construct Data

Modeling / Data Mining

- Use at least two mining algorithms on the dataset. -

Modeling technique

Test Design

- Explain how you will evaluate the tests

Build and Assess the Model

Evaluation

- Discuss the overall accuracy of your final model - non-technical terms, what information you have learned from the dataset. - Also discuss any limitations of the dataset that may have effected model accuracy