# Application of Bayesian Networks in Consumer Service Industry and Healthcare

Le Zhang[1], Yuan Gao[1], Balmatee Bidassie[3], and Vincent G. Duffy[1,2]

[1] School of Industrial Engineering
[2] School of Agriculture and Biological Engineering, Purdue University,
47907 West Lafayette, IN, USA
{zhan1255,gao186,duffy}@purdue.edu
[3] Department of Veteran Affairs, Center for Applied Systems Engineering,
Detroit, MI, USA
balmatee.bidassie@va.gov

**Abstract.** Bayesian networks are powerful in data mining and analyzing causal relationships of an uncertain-reasoning problem. The implementation of Bayesian networks in industry and healthcare diagnosis can facilitate the process of locating causations in complex issues. This study conducted two case studies by BayesiaLab in consumer service and healthcare domain. Case Study One used unsupervised learning and supervised learning on the individual data set of county road traffic volume in Indiana State and concluded that road type has the most significant impact on daily vehicle miles traveled. In Case Study Two, only supervised learning was used to observe the aggregated data of adverse mental health effect on civilians, deployed veterans and nondeployed veterans of different genders. Both types of veterans showed higher probability to have adverse mental health compared to civilians. In conclusion, Bayesian networks provided valid results to support prior research. Further research is needed to investigate the differences between using individual data and aggregated data, and to apply Bayesian networks in meta-analysis.

**Keywords:** Bayesian Networks, BayesiaLab, Traffic Volume, Mental Health.

## 1 Introduction

Data mining is the automated or convenient extraction of patterns representing knowledge implicitly stored or captured in large databases and other storage media. It assists human in data collection, instead of the traditional method which relies on manual analysis and interpretation [1, 2]. Among data mining techniques, Bayesian networks, based on Bayes' Theorem, is a powerful technique to work with probabilities rather than relying on factual observations, especially for complicated issues involving association and causal relationships yet to be discovered.

A Bayesian network is an annotated directed graph that encodes probabilistic relationships among distinctions of interest in an uncertain-reasoning problem [3]. Its graphic presentation provides an exceptional demonstration of the relationships (arcs)

among all the factors (nodes) in a complex problem. Since the model deals with dependencies among all variables, it can cope with incomplete and uncertain data, and especially uncertain rules of reasoning, strengthening the power of diagnosis and prediction [4]. Bayesian networks also support learning abilities, allowing automatic application of this methodology in complex problems [5]. Another important feature is "omnidirectional inference" - while traditional statistical models usually contain one dependent and many independent variables, all variables can be treated equivalently in a Bayesian network, which is interesting for exploratory research [6].

In the last decade, many fields, from the traditional industries to new areas, have seen applications of Bayesian networks. Bayesian networks have risen to prominence as the preferred technology for probabilistic reasoning in artificial intelligence, with a proliferation of techniques for fast and approximate updating and also for their automated discovery from data [7]. In the management of complex industrial systems, it is used for dependability, risk analysis and maintenance areas [8]. In finance and banking, it can be used for fraud detection and credit scoring. In marketing, it is used for consumer survey analysis, market segmentation and simulation. In healthcare, it's used in diagnostic systems. Other applications include quality management, operational safety simulation, etc [9].

Mental health disease diagnosis is a field with many factors involved and interacting with each other. Clinical diagnosis of mental disorders is more complicated and more difficult compared to other sectors in primary care, with more risk factors involved, including personal profile, residential and work environment, cultural and sociological settings, etc [10]. For veterans, the challenges are even bigger due to additional factors, such as their military service experience, higher chances of physical injuries and disabilities, and the need to get readjusted to life out of the armed forces [11]. Bayesian networks can be introduced as a holistic approach to monitor and evaluate mental health risks from patients' perspective, and further facilitate early diagnose and prevention of implicit mental conditions.

In this study, we propose an application in consumer service to identify key factors and causal relationships, used as strategic guidelines in consumer satisfaction improvement and cost saving. Following the case studies of traffic volume analysis of the road system of the state of Indiana, an analysis of mental health risk factor among veterans and civilians will be showed to illustrate Bayesian networks implementation in industry and health care. This research demonstrates how to find out the hidden ties through BayesiaLab [12], and how it can benefit early diagnosis and prevention of mental health problems among veterans.

## 2     Methodology and Tool

The tool used in both case studies was BayesiaLab, a modeling software developed and supported by Bayesia, a designer of decision aid software packages in Bayesian networks for data mining [12]. BayesiaLab provides a complete laboratory for handling Bayesian networks to develop, communicate with and use readable illustrated decisional models that are strictly faithful to reality. Among its many features, the

most appreciable ones for this case study include: highly intuitive graphical network presentation, learning capability, and the Bayesian power of inference.

In Data mining, unsupervised learning is defined as when a learning human, animal or man-made system observes its surroundings and, based on observations adapts its behavior without being told to associate given observations to given desired responses, as opposed to supervised learning, which has a targeted classifier whose value is analyzed based on the influences of the other factors [13]. With its learning function, individual data can be translated into Bayesian networks. Both learning techniques are supported by BayesiaLab with various algorithms, among which Naive Bayesian algorithm is based on the assumption (known as class-conditional independence) that the effect of an attribute value on a given class is independent of the values of the other attributes [14]. It simplifies the computations involved and, in this sense, is considered "naive".

Bayesian networks can be developed based on individual data or aggregated data. Individual data presents the parameter of every factor on each subject. In contrast, aggregated data presents the percentages of each factor on every subject. Based on the learning capability of BayesiaLab, Case Study One used individual data to develop a model. Case Study Two showed how aggregated data can be used to conduct meta-analysis [15] on mental health.

## 3     Case Study One: Traffic Volume Analysis of the Road System of the State of Indiana

Daily Vehicle Miles Traveled (DVMT) is a measure of the traffic volume flowing along a roadway during an average 24 hours period [16]. It is calculated by multiplying the Average Annual Daily Traffic (AADT) by the length of the road. This data is an important indicator used to assess transportation needs, system performance and highway planning and program recommendations [17]. In this case study, we looked into the DVMT data of each county's highway system in the State of Indiana from 2006 to 2010, and use learning capability of Bayesian networks to build a simply model to explore the risk factors affecting DVMT.

### 3.1     Hypothesis

Based on the definition, DVMT is directly affected by the average daily traffic volume and the length of roads. In Indiana, traffic volume raw data are collected by permanent continuous count stations, as well as portable traffic counters, and then adjusted and seasonally factored [17]. Road traffic volume is affected by many factors, ranging from road conditions, gasoline prices and toll, weather, and social and demographic factors such as age of the population, household size, labor force participation and car ownership [18]. This case study was designed based on the hypothesis that a certain extent of relationship exists between county populations and DVMT. Compared to other factors like vehicle ownership and labor force, the total population of a county is not known to have an explicit causal relationship with traffic volume,

but it is logical to assume that the more people reside nearby, the higher traffic volume the area is likely to have. Therefore, it was chosen for the case study to test and demonstrate Bayesian networks' ability to explore unknown relationships.

## 3.2    Modeling and Analysis

DVMT data came from the Historic Vehicle Miles Traveled (VMT) by County & Systems data published by Indiana Department of Transportation (INDOT) [19]. For the 92 counties of Indiana, each of them includes some or all of 4 types of routes: City and County Roads, Interstate Roads (I), State Roads (SR) and U.S. Highways (US). County population is based on data published by STATS Indiana [20].

The source data were prepared so that each type of road in each county is listed as an individual data point. County population is duplicated for multiple road types. In the model, county names are omitted because name itself doesn't affect the traffic volume in any way. And a new parameter was calculated by dividing DVMT with the length of road. In theory, it should be equal to AADT. However, without the original AADT data from INDOT, or the exact method of calculating DVMT, it is named "Usage Rate" to indicate the extent of usage of the road. By doing this, the bias caused by the multiplying effect of road length in traffic volume is removed. Table 1 shows an example of the data of Allen County in 2006 (Note that the column "County" was removed during actual modeling).

**Table 1.** Data Example of Allen County in 2006

| Year | County | Road Type | Length (Miles) | DVMT | Usage Rate | Population |
|------|--------|-----------|----------------|------|------------|-----------|
| 2006 | Allen | City and County Roads | 2686.94 | 6061000 | **2255.73** | 345976 |
| 2006 | Allen | I | 98.8 | 2104000 | **21295.55** | 345976 |
| 2006 | Allen | SR | 87.24 | 936000 | **10729.02** | 345976 |
| 2006 | Allen | US | 65.87 | 1017000 | **15439.5** | 345976 |

During data import, KMeans was used as the discretization method for its benefits in removing the effect of outliers [21]. An initial analysis using BayesiaLab's unsupervised learning feature and distance mapping layout is shown in Fig. 1. In this model, BayesiaLab automatically connected nodes (factors) based on the mutual information and correlation between each two, and the amount of information each node can provide to reduce the uncertainties of other nodes, given the value of the specific node itself. The node "Year" is greyed out with no connection due to its insignificance. Using the feature Distance Mapping, the model is presented in a two-dimensional form where the length of arc between any two nodes indicates the Mutual Information as in information technology [22]. The closer two nodes are, the higher mutual information value they have in between. The color and figures of each arc reflect the relationship between two nodes in Pearson Correlation coefficient [23]: red color means negative correlation and blue means positive.
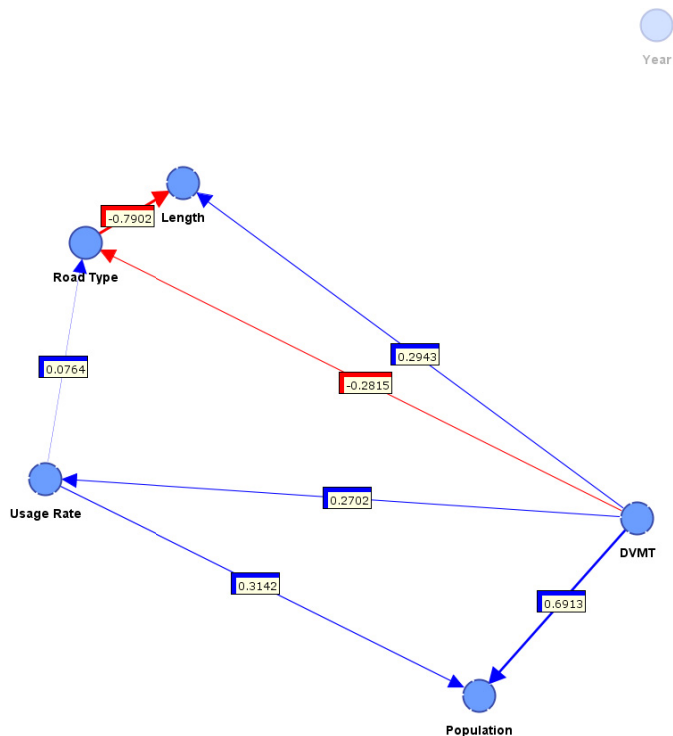
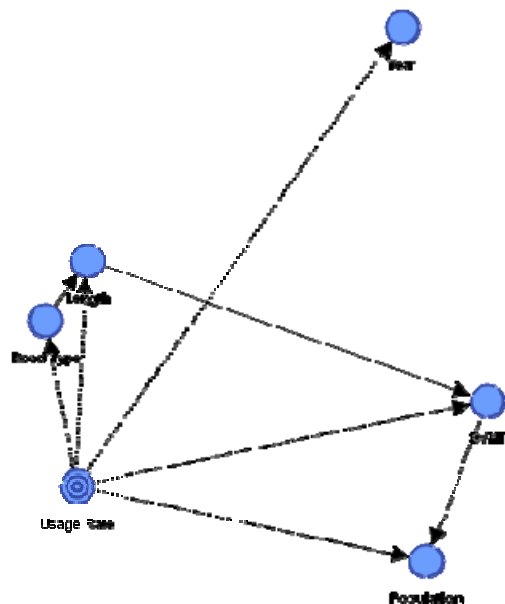**Fig. 1.** Unsupervised Model in Mutual Information Distance Mapping Layout



**Fig. 2.** Supervised Model with Usage Rate as Target Node

DVMT and Road Type are the two most information-rich parameters in this diagram, and Road Type has strong mutual information ties with Length and Usage Rate, while DVMT is related to county Population. It's noteworthy that Road Type is treated as discrete data and ranked in alphabetic order by default: City and County Road first and U.S. last. This helps understand the positive/negative correlation.

Next, a targeted analysis using BayesiaLab's Supervised Learning feature was performed with Usage Rate set as target node. Using Augmented Naïve Bayes learning algorithm [14], the system is modeled as in Fig. 2. The predefined Naive Bayes structure is highlighted with the dotted arcs, while the augmented arcs (from the additional Unsupervised Learning) are shown with solid arcs.

### 3.3    Results

Fig. 3 illustrated the posterior probability analysis on the supervised model, showing that Road Type has a significant impact on DVMT, Usage Rate and other variables. Comparing the left and right sides conditioned on different road types, it's obvious that in Indiana, City and County Roads are most likely to have the longest lengths, but lowest average traffic volume, and they are most likely to exist in counties with smaller populations. While Interstate roads are the shortest in length, but have higher volume and are more likely to exist in counties with larger populations. We only present the comparison result of these two types of roads because they demonstrated the largest divergence. The other two road types, US and SR, are in the middle, with SR showing relatively higher probability of lower usage rate than US. Overall, from 2006 to 2010, time didn't play a big role in the amount of road traffic or usage rate.
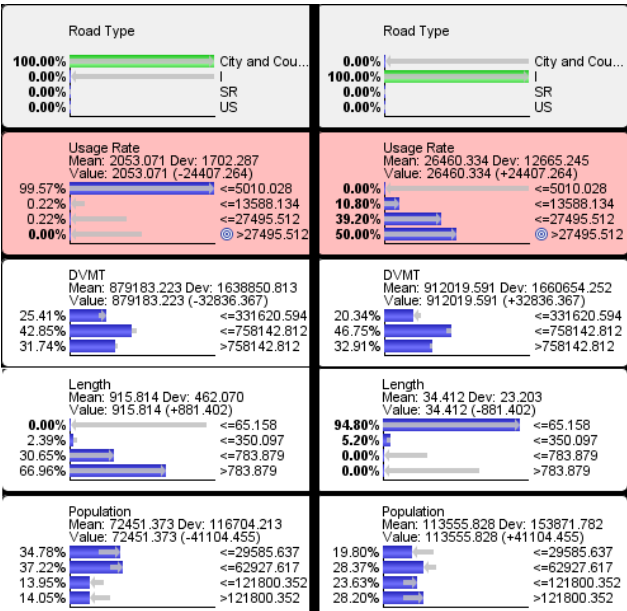


**Fig. 3.** Comparison of Impact of Road Type

# 4    Case Study Two: Analysis of Mental Health Risk Factor among Veterans and Civilians

Veterans' mental health problems are considered related to war experience for military service [24]. The numbers of mental health diagnosis are increasing and often coexisting with other medical problems, causing a "downward spiral" as stated by Secretary Eric K. Shinseki: "unseen" psychological wounds interplaying with biological and physical ailments, resulting in personal isolations and even more serious issues like homelessness and suicide [11]. Meanwhile, experience from military services could put veterans into long term struggles caused by brain injuries, disabilities and chronic diseases such as diabetes, obesity and hypertension, which are found to contribute to mental disorders, making their self-management more difficult and causing suicidal ideation [25]. Posttraumatic Stress Disorder (PTSD) is a psychiatric disorder that can occur after someone goes through a traumatic event like war, assault, or disaster [26]. To understand the adverse mental health effect within veterans, a large amount of studies focus on risk factors like gender, age, race, combat exposure, and war zone deployment, and PTSD [27, 28]. However, veterans who hadn't served in a war theater were also considered associated with war zone deployment and combat exposure in researches [29].

Hoglund [29] summarized in his study that adverse mental health effect was associated with veterans who had military service in a combat and war zone (deployed veterans) and female veterans who didn't have military service in a combat or war zone (nondepoyed veterans). Calculated the odds ratio of risk factors, he compared mental health conditions of two groups of veterans with civilians by exploring gender differences. This result here is worth a meta-analysis with Bayesian networks to explore probability relationships between those complex factors [15]. Based on his results, this case study developed a supervised learning model in BaysiaLab [12] to investigate the relationship between mental health effect and deployment status.

## 4.1    Hypothesis

This study takes deployment status of veterans into consideration, as well as social demographic characteristics to locate the risk factors and causal relationships with adverse mental health. The hypothesis in this study is either deployed veterans or nondeployed veterans are more associated with adverse mental health than civilians.

## 4.2    Modeling and Analysis

Instead of using odds ratio, this study developed a supervised learning model with BayesiaLab [12] by following the distribution of men and women (stratified by civilian, developed veteran, and nondeveloped veteran) with respect to race/ethnicity, marital status, education, employment status, general health, and mental health.  This study used aggregated data of Behavior Risk Factor Surveillance survey which was summarized by Hoglund [29]. Table 2 showed the data set for deployed veterans.

The data set included the race/ethnicity, marital status, education, employment status, general health, and a self-report of the number of days when mental health was not good during the previous 30 days (14 days or more indicated adverse mental health).

**Table 2.** Aggregated Data from Behavior Risk Factor Surveillance Survey [29]

|  |  | Deployed Veterans (N=978) | |
|  |  | Male (N=846) | Female (N=132) |
| --- | --- | --- | --- |
| Race/Ethnicity | Nonwhite and/or Hispanic | 0.20 | 0.33 |
|  | White Non-Hispanic | 0.80 | 0.67 |
| Marital Status | No Spouse or Partner | 0.25 | 0.41 |
|  | Spouse or Partner | 0.75 | 0.59 |
| Education | High School or Less | 0.27 | 0.08 |
|  | Some College or More | 0.73 | 0.92 |
| Employment | Not Employed | 0.12 | 0.13 |
|  | Employed | 0.88 | 0.87 |
| General Health | Fair-to-Poor | 0.12 | 0.15 |
|  | Good-to-Excellent | 0.88 | 0.85 |
| Mental Health | 14+Poor Mental Health Days | 0.09 | 0.17 |
|  | 13 or Fewer Poor Mental Health Days | 0.91 | 0.83 |

Eight factors (nodes) were created and mental health was set as the target factor in supervised learning. The probability for the effects of each factor were given by aggregated data from Behavior Risk Factor Surveillance survey [29]. As showed in Table 2, deployment and gender had relation to other factors. Thus, we manually connected arcs between deployment and gender with other factors. The model was presented with radial layout as Fig. 4. In probability mode, the probability of deployment and gender can be changed to 100% or 0. In this way, the probability of adverse mental health effect for each combination of deployment and gender can be observed. The probabilities were presented in histogram as Fig. 5.
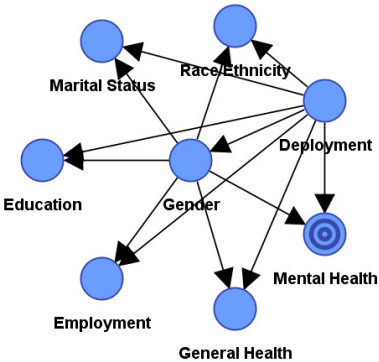


**Fig. 4.** Bayesian model of deployment data

## 4.3    Results

The histogram from Fig. 5 showed a higher probability of adverse mental health effect on veterans (9% and 11% for male deployed veterans and nondeployed veterans, 17% and 19% for female deployed veterans and nondeployed veterans) than on civilians (8% for male civilians, 13% for female civilians), which evidently support the hypothesis. Hoglund's research [29] indicated deployed veterans were associated with adverse mental health for men and possibly women; nondeployed veterans were associated with adverse mental health for women, but not for men. This study provided similar results on male deployed veterans, female deployed veterans, and female nondeployed veterans. However, male nondeployed veterans also showed a high probability on having adverse mental health.
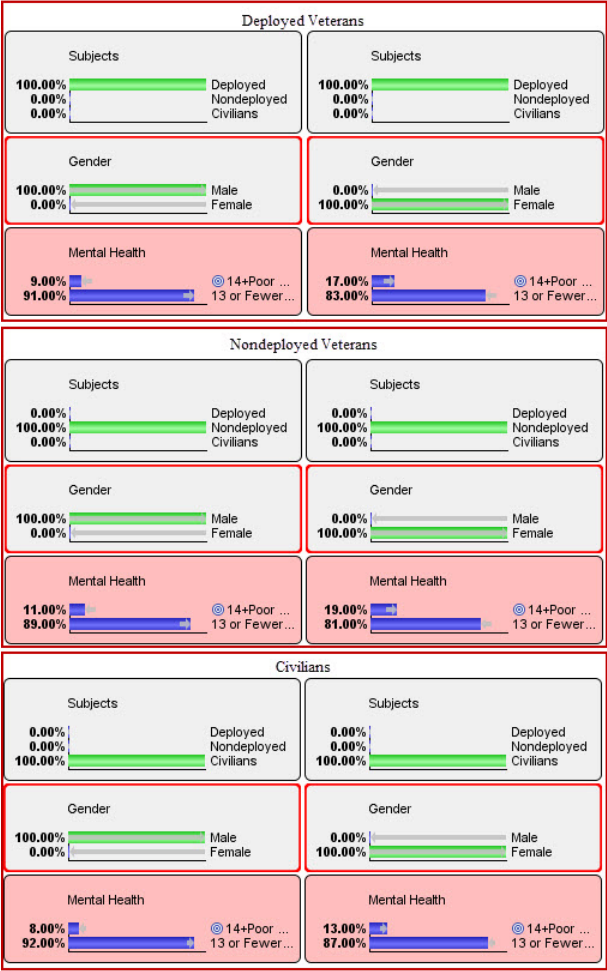


**Fig. 5.** Comparison of mental health between developed veterans, nondeveloped veterans, civilians by gender

# 5     Discussions

The findings of Case Study One are in line with existing literature. The interstate highway system has a much higher density of use than other components of U.S. surface transportation system. A 1996 report [30] shows interstate highways take up 23% of the market share of America's surface transport, and is 26 times as many person miles per route mile as all other roads (including low usage rural roads). Another more recent report in 2008 [18] shows that the rate of growth in VMT (Vehicle Miles Traveled) has fallen sharply since 2000, with a miniscule 0.6% in 2006. And this leveling off in VMT growth may be a long-term trend due to various socio-economic factors. This report also points out that regional population has a minor influence on VMT, though specific age/occupational groups have a more significant role in the growth of VMT, such as car owners, female labor force and working age population.

The Bayesian networks model of Case Study Two was developed from aggregated data, which had been translated into the probability of effect on each factor. Aggregated data also defines the relation between each factor. In such condition, only supervised learning can be applied on this data. Similar results from Hoglund's research [29] were provided. However, this model showed male nondeployed veterans (11%) presented a higher probability on having adverse mental health than male deployed veterans (9%), which differs from Hoglund's conclusion that male nondeployed veterans has insignificant association with adverse mental health effect. Since Bayesian networks provide probability relationships of each event rather than significant influence, even if adverse mental health effects present insignificant influence on nondeployed veterans, it still has probability to happen on nondeployed veterans and this probability is higher than deployed veterans.

This paper demonstrated the strength of Bayesian networks in discovering relationships in systems involving multiple risk factors, and the validity of this method through verification by existing literature and comparative study. The second case study showcased a new application area of this methodology in diagnosis of mental health issues among veterans. Compared to traditional statistics analysis methods, Bayesian networks enable modeling in intuitive graphic views for analysts to identify the key factors and relationships at a glance. Modern software such as BayesiaLab have incorporated advanced visualization features and learning algorithms to further utilize the probabilistic predictive and diagnostic power of Bayes' theorem. However, it must be pointed out that the direction of the arc should not be interpreted as causal relationship, but statistical association. The existence of causation should be verified by further study and literature.

The difference in data processing between the two cases reflect different scenario in real world research. Sometimes, data of each individual subject is available. However, in some cases, such data is not directly accessible, or will take enormous time and efforts to collect, especially in the health care field where the subjects are usually patients and clinical experiments are expensive. The concept of meta-analysis advocates leveraging existing researches on similar topics in a systematic view [15]. To further verify the validity of modeling with aggregated data, a comparative study was performed based on data set of Case Study One to sum up individual data into

percentages and obtained similar results: road type affects DVMT and usage rate most. But the results are not exactly the same with minor variance (e.g. the probability for DVMT to be greater than 657,000 with aggregate data is higher than individual date by 1.47 percentage point). This may be due to the integration effect during aggregation.  It is also noteworthy that different discretization methods will affect the presentation of results. Therefore, using aggregated data should be based on a good knowledge of the type and distribution of the data set, and a good understanding of the research subject, content are and context. While this study shows the potential of using Bayesian networks in the context of meta-analysis, further study should be conducted to investigate the detailed application in different areas.

# References

 1. Fayyad, U., Piatetsky-shapiro, G., Smyth, P.: From Data Mining to Knowledge Discovery in Databases. 17, 37–54 (1996)
 2. Han, J., Kamber, M., Pei, J.: Data Mining Concepts and Techniques, 3rd edn. (2011)
 3. Heckerman, D., Mamdani, A., Wellman, M.P.: Real-world applications of Bayesian networks. Commun. ACM. 38, 24–26 (1995)
 4. Stassopoulou, A., Petrou, M.: Obtaining the correspondence between Bayesian and Neural Networks. Int. J. pattern Recognit. Artif. Intell. 12, 901–920 (1998)
 5. Santander Meteorology Group: Data mining and artificial intelligence: Bayesian and Neural networks, http://www.meteo.unican.es/research/datamining
 6. Conrady, S., Jouffe, L.: Vehicle Size, Weight, and Injury Risk High-Dimensional Modeling and Causal Inference with Bayesian Networks Table of Contents (2013)
 7. Korb, K.B., Nicholson, A.E.: The Causal Interpretation of Bayesian Networks. In: Holmes, D.E., Jain, L.C. (eds.) Innovations in Bayesian Networks. SCI, vol. 156, pp. 83–116. Springer, Heidelberg (2008)
 8. Weber, P., Medina-Oliva, G., Simon, C., Iung, B.: Overview on Bayesian networks applications for dependability, risk analysis and maintenance areas. Eng. Appl. Artif. Intell. 25, 671–682 (2012)
 9. Neapolitan, R.E.: Learning Bayesian Networks. Pearson Prentice Hall Upper Saddle River (2004)
10. Laufer, N., Zilber, N., Jecsmien, P., Maoz, B., Grupper, D., Hermesh, H., Gilad, R., Weizman, A., Munitz, H.: Mental disorders in primary care in Israel: Prevalence and risk factors. Soc. Psychiatry Psychiatr. Epidemiol. 48, 1539–1554 (2013)
11. Department of Veteran Affairs: Strategic Plan Refresh, FY 2011-2015., Washington, DC 20420 (2011)
12. Bayesia: Bayesia, http://www.bayesia.com/en/about-us/index.php
13. Huang, T.-M., Kecman, V., Kopriva, I.: Unsupervised Learning by Principal and Independent Component Analysis. In: Huang, T.-M., Kecman, V., Kopriva, I.: Kernel Based Algorithms for Mining Huge Data Sets. SCI, vol. 17, pp. 175–208. Springer, Berlin (2006)

14. Han, J., Kamber, M., Pei, J.: Classification: Basic Concepts. In: Data Mining Concepts and Techniques. pp. 327–391 (2012)
15. Borenstein, M., Hedges, L.V., Higgins, J.P., Rothstein, H.R.: Introduction to meta-analysis. Wiley.com (2011)
16. Division of Planning, Office of Technical Services, O.D. of T.: Daily Vehicle Miles Traveled Report (DVMT),
    `http://www.dot.state.oh.us/Divisions/Planning/TechServ/traff`
    `ic/Pages/DVMT.aspx`
17. Autumn Young: Latest INDOT Traffic Adjustment Factors., Indianapolis (2013)
18. East-West Gateway: Trends in Regional Traffic Volumes Signs of Change (2008),
    `http://www.ewgateway.org/pdffiles/library/trans/trafficvolum`
    `es/vmtrpt.pdf`
19. Indiana Department of Transportation: Historic VMT by County & Systems (2006-2011),
    `http://www.in.gov/indot/files/TrafficStastics_HistoricIndian`
    `aVMTByCounty2006-2011.pdf`
20. STATS Indiana: Population,
    `https://www.stats.indiana.edu/topic/population.asp`
21. Hautamäki, V., Cherednichenko, S., Kärkkäinen, I., Kinnunen, T., Fränti, P.: Improving K-means by outlier removal. In: Kalviainen, H., Parkkinen, J., Kaarna, A. (eds.) SCIA 2005. LNCS, vol. 3540, pp. 978–987. Springer, Heidelberg (2005)
22. Cellucci, C.J., Albano, A.M., Rapp, P.E.: Statistical validation of mutual information calculations: Comparison of alternative numerical algorithms. Phys. Rev. E. 71, 66208 (2005)
23. Lee Rodgers, J., Nicewander, W.A.: Thirteen Ways to Look at the Correlation Coefficient. Am. Stat. 42, 59–66 (1988)
24. Jakupcak, M., Hoerster, K.D., Blais, R.K., Malte, C.A., Hunt, S., Seal, K.: Readiness for Change Predicts VA Mental Healthcare Utilization Among Iraq and Afghanistan War Veterans. 165–168 (2013)
25. Bossarte, R.M., Blosnich, J.R., Piegari, R.I., Hill, L.L., Kane, V.: Housing instability and mental distress among US veterans. Am. J. Public Health. 103(suppl.), S213–6 (2013)
26. Buckley, T., Tofler, G., Prigerson, H.G.: Posttraumatic Stress Disorder as a Risk Factor for Cardiovascular Disease: A Literature Review and Proposed Mechanisms. Curr. Cardiovasc. Risk Rep. 7, 506–513 (2013)
27. Maguen, S., Ren, L., Bosch, J.O., Marmar, C.R., Seal, K.H.: Gender differences in mental health diagnoses among Iraq and Afghanistan veterans enrolled in veterans affairs health care. Am. J. Public Health. 100, 2450–2456 (2010)
28. Metraux, S., Clegg, L.X., Daigh, J.D., Culhane, D.P., Kane, V.: Risk factors for becoming homeless among a cohort of veterans who served in the era of the Iraq and Afghanistan conflicts. Am. J. Public Health. 103(suppl.), S255–61 (2013)
29. Hoglund, M.W., Schwartz, R.M.: Mental health in deployed and nondeployed veteran men and women in comparison with their civilian counterparts. Mil. Med. 179, 19–25 (2014)
30. Cox, W., Love, J.: 40 Years of the US Interstate Highway System: An Analysis, The Best Investment A Nation Ever Made (1996)