
Mental Health in Tech Survey

Anmisha Reddy Ramidi **Sindhura Uppu**
Department of Computer Science Department of Computer Science
University at Buffalo University at Buffalo
anmishar@buffalo.edu suppu@buffalo.edu

Abstract

This project aims to measure and model the attitude towards mental health and frequency of mental health disorders in the tech workplace. The goal is to achieve this by generating a probabilistic graphical model or Bayesian Network to model the data answer queries. The dataset is obtained from [here](#). Exact and Approximate inference algorithms have been applied on this network using which the queries have been answered.

1 Problem Domain

The dataset that is being modeled in this project belongs to the Kaggle dataset. The data is actually collected from a survey that is conducted on people working in the Tech workplace. The domain Mental Health is a very important topic which is being taken seriously by many corporations. This dataset is from a 2014 survey that measures attitudes towards mental health and frequency of mental health disorders in the tech workplace. It is said that better data leads to better health services. Surveys such as these help in determining how issues related to mental health in technical corporations effect the majority of population. The major goal of this project is to analyze the attitude of people working at the Tech places and help in improving the conditions.

1.1 Data Set

The data set is the collection of the answers the workers have answered during the Survey. The data set contains around 1260 rows and has 25 variables. The variables help in describing the attitude and features of how mental illness is treated at the tech workplace. The data set is used to determine a Bayesian network, where links are generated based on the relationship of the variables with each other. The data set was cleaned for better understanding and for getting better results when modeled. The following are the descriptions of each variables and what they represent.

Timestamp: Timestamp at which the survey was taken.

Age: Age of the person taking the survey.

Gender: Gender of the person taking the survey.

Country: Country to which the person belongs.

State: If you live in the United States, which state or territory do you live in?

Self_employed: Are you self-employed?

Family_history: Do you have a family history of mental illness?

Treatment: Have you sought treatment for a mental health condition?

Work_interfere: If you have a mental health condition, do you feel that it interferes with your work?

44 No_employees: How many employees does your company or organization have?
45 Remote_work: Do you work remotely (outside of an office) at least 50% of the time?
46 Tech_company: Is your employer primarily a tech company/organization?
47 Benefits: Does your employer provide mental health benefits?
48 Care_options: Do you know the options for mental health care your employer provides?
49 Wellness_program: Has your employer ever discussed mental health as part of an employee
50 wellness program?
51 Seek_help: Does your employer provide resources to learn more about mental health issues
52 and how to seek help?
53 Anonymity: Is your anonymity protected if you choose to take advantage of mental health or
54 substance abuse treatment resources?
55 Leave: How easy is it for you to take medical leave for a mental health condition?
56 Mental_health_consequence: Do you think that discussing a mental health issue with your
57 employer would have negative consequences?
58 Phys_health_consequence: Do you think that discussing a physical health issue with your
59 employer would have negative consequences?
60 Coworkers: Would you be willing to discuss a mental health issue with your coworkers?
61 Supervisor: Would you be willing to discuss a mental health issue with your supervisor(s)?
62 Mental_health_interview: Would you bring up a mental health issue with a potential employer
63 in an interview?
64 Phys_health_interview: Would you bring up a physical health issue with a potential employer
65 in an interview?
66 Mental_vs_physical: Do you feel that your employer takes mental health as seriously as
67 physical health?
68 Obs_consequence: Have you heard of or observed negative consequences for coworkers with
69 mental health conditions in your workplace?
70 Comments: Any additional notes or comments

71

72 **1.2 Bayesian Network Model**

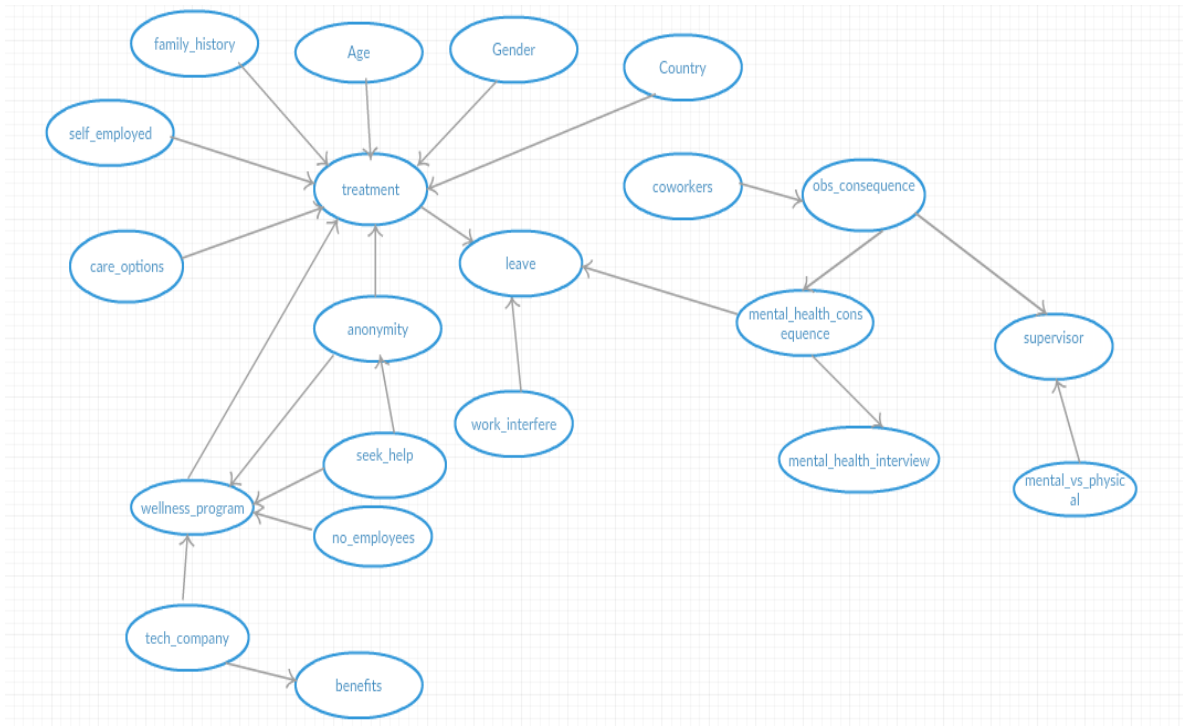
73

74 A Bayesian Network represents the causal probabilistic relationship among a set of random
75 variables, their conditional dependencies and it provides a compact representation of a joint
76 probability distribution^[1]. Edges represent conditional dependencies; nodes that are not connected
77 represent variables that are conditionally independent of each other. Each node is associated with a
78 probability function that take a particular set of values for the node's parent variables, and gives the
79 probability of the variable represented by the node. They handle uncertainty through the established
80 theory of probability.

81

82 The variables of the data set were assessed and relations were drawn. The Bayesian Network
83 consists of 21 variables which contain over 22 links among them. The links represent the causal
84 relationship between these 21 variables. The Bayesian Network that we represented was constructed
85 manually based on intuition.

86



The Bayesian Network can be generated in Python using pgmpy library as follows:

```

Mental_health_model=
BayesianModel([('Age','treatment'),('Gender','treatment'),('Country','treatment'),('family_hist
ory','treatment'),('self_employed','treatment'),('care_options','treatment'),('anonymity','treatm
ent'),('treatment','leave'),('work_interfere','leave'),('coworkers','obs_consequence'),('obs_cons
equence','supervisor'),('mental_vs_physical','supervisor'),('obs_consequence','mental_health_
consequence'),('mental_health_consequence','mental_health_interview'),('anonymity','wellne
ss_program'),('no_employees','wellness_program'),('seek_help','wellness_program'),('tech_co
mpany','wellness_program'),('tech_company','benefits'),('seek_help','anonymity'),('mental_he
alth_consequence','leave'),('wellness_program','treatment')])

```

Pgmpy library has inbuilt models including Bayesian Network model. The nodes that have edges are sent as parameters using the node names as attributes in the model. The attribute in the left represents the parent while the one in the right is the child node in the Bayesian network.

1.3 Conditional Probability Distributions

This model is now fitted with the Mental Health dataset using Maximum Likelihood Estimator. Likelihood of a dataset is the probability of obtaining that particular set of data, given the chosen probability distribution model. There are unknown model parameters in this expression whose values that maximize the sample likelihood are known as the Maximum Likelihood Estimates.

Snippet to fit using Maximum Likelihood Estimator:

```
Mental_health_model.fit(train, estimator = MaximumLikelihoodEstimator)
```

Train is the training data of 600 samples that have been retrieved from the dataset.

118 Conditional Probability Distributions for each node were also generated from the Maximum
 119 Likelihood Estimator using the edges in the Bayesian network. A conditional probability
 120 distribution over b via a conditional distribution with a means that the distribution over b
 121 depends on the value of a. The CPD's will be represented in a tabular format as follows:

```
122
123 +-----+-----+
124 | Age (20-40) | 0.5 |
125 +-----+-----+
126 | Age (40-60) | 0.5 |
127 +-----+-----+
```

128 Age is an independent (parent) node in the network.

```
129
130 +-----+-----+-----+
131 | seek_help          | seek_help(Don't know) | seek_help(No) |
132 +-----+-----+-----+
133 | anonymity(Don't know) | 0.3333333333333333 | 0.3333333333333333 |
134 +-----+-----+-----+
135 | anonymity(No)        | 0.3333333333333333 | 0.3333333333333333 |
136 +-----+-----+-----+
137 | anonymity(Yes)       | 0.3333333333333333 | 0.3333333333333333 |
138 +-----+-----+-----+
```

139 Seek_help is the parent node while anonymity is the child node and their CPD's based on their
 140 dependency is obtained as shown above.

141

142

143 **2 Inference Models**

144

145 **2.1 Belief Propagation**

146 Computing the a posteriori belief of a variable in a general Bayesian Network is NP-hard.
 147 Belief Propagation is an Approximate Inference algorithm. It is an efficient way to solve
 148 inference problems based on passing local messages^[2]. It is available in the pgmpy.inference
 149 library as a class for performing inference using Belief Propagation model. It creates a junction
 150 tree or Clique tree for the input probabilistic graphical model and performs calibration of the
 151 junction tree so formed using belief propagation.

152 Code snippet for Belief Propagation implementation^[3]:

```
153 belief_prop = BeliefPropagation(Mental_health_model)
```

154

155 Now that the inference is drawn from this model, queries can be run on the model to analyze
 156 the variation and independence of one or more variables over other such variables in the
 157 model. A few of the queries implemented are as follows:

158

159 Query1:

```
160 bp1 = belief_prop.query(variables=['leave','wellness_program'],evidence={'tech_company' :
161 0})
162 print(bp1['leave'])
163 print(bp1['wellness_program'])
```

164 Here are looking at how easy is it for the employee to take leave and whether they are
 165 aware of the wellness programs in their company provided that it is a tech company they are
 166 working for.

167 Tech_company : 0 maps to the 'Yes' value for the variable in the dataset.

168

169 Result:

```
170 +-----+-----+
171 | leave | phi(leave) |
172 |-----+-----|
173 | leave_0 | 0.2756 |
174 | leave_1 | 0.0977 |
175 | leave_2 | 0.2311 |
176 | leave_3 | 0.0977 |
177 | leave_4 | 0.2978 |
178 +-----+-----+
179 +-----+-----+
180 | wellness_program | phi(wellness_program) |
181 |-----+-----|
182 | wellness_program_0 | 0.2800 |
183 | wellness_program_1 | 0.4400 |
184 | wellness_program_2 | 0.2800 |
185 +-----+-----+
```

186

187 Query 2:

```
188 bp2 = belief_prop.query(variables=['treatment'],evidence={'Age' : 1, 'Gender' : 1,
189 'family_history' : 1})
190 print(bp2['treatment'])
```

191 The attribute 'Treatment' is conditionally dependent on the attributes 'Age', 'Gender'
192 and 'family_history'. This query calculates how many males within the age range 20-40 and
193 had a family history of mental illness have opted for mental health treatment.

194 Result:

```
195 +-----+-----+
196 | treatment | phi(treatment) |
197 |-----+-----|
198 | treatment_0 | 0.4432 |
199 | treatment_1 | 0.5568 |
200 +-----+-----+
```

201

202 Query 3:

```
203 bp3 = belief_prop.query(variables=['benefits','treatment'],evidence={'tech_company' : 1})
204 print(bp3['benefits'])
205 print(bp3['treatment'])
```

206 The above query shows how being in a tech company relates to whether their
207 employer provides any benefits for mental health illness and how many have taken treatment
208 from the organization. It should be noted that tech_treatment has an indirect dependence on
209 tech_company while benefits has a direct dependence on tech_company in our Bayesian
210 network.

211 Result:

```
212 +-----+-----+
213 | benefits | phi(benefits) |
214 |-----+-----|
215 | benefits_0 | 0.0000 |
216 | benefits_1 | 0.0000 |
217 | benefits_2 | 1.0000 |
218 +-----+-----+
219
```

```

220 +-----+-----+
221 | treatment | phi(treatment) |
222 |-----+-----|
223 | treatment_0 | 0.4441 |
224 | treatment_1 | 0.5559 |
225 +-----+-----+

```

226

227 Query 4:

```

228 bp7 = belief_prop.query(variables=['treatment','leave'],evidence={'seek_help' :
229 2,'care_options' : 1})
230 print(bp7['treatment'])
231 print(bp7['leave'])

```

232 This query identifies the relation between treatment, leave attributes independently
233 while being conditionally independent on seek_help and care_options. A value 2 in
234 seek_options means that the employer provides enough resources to learn more about health
235 issues and a value 1 in care_options means that the employee knows the options for mental
236 health that the employer provides.

237 Result:

```

238 +-----+-----+
239 | treatment | phi(treatment) |
240 |-----+-----|
241 | treatment_0 | 0.3933 |
242 | treatment_1 | 0.6067 |
243 +-----+-----+
244 +-----+-----+
245 | leave | phi(leave) |
246 |-----+-----|
247 | leave_0 | 0.2825 |
248 | leave_1 | 0.0884 |
249 | leave_2 | 0.2340 |
250 | leave_3 | 0.0884 |
251 | leave_4 | 0.3068 |
252 +-----+-----+

```

253

254 Query 5:

```

255 bp5 =
256 belief_prop.query(variables=['mental_health_interview','supervisor'],evidence={'obs_consequence' : 1})
257
258 print(bp5['mental_health_interview'])
259 print(bp5['supervisor'])

```

260 This query evaluates how many people are willing to disclose their mental health state
261 to prospective employers in an interview and also to their supervisor in the present
262 organization based on the consequences they observed for coworkers with mental health
263 conditions in their workplace. A value 1 in 'obs_consequence' maps to 'Yes' in the dataset.

264 Result:

```

265 +-----+-----+
266 | mental_health_interview | phi(mental_health_interview) |
267 |-----+-----|
268 | mental_health_interview_0 | 0.2000 |
269 | mental_health_interview_1 | 0.8000 |
270 | mental_health_interview_2 | 0.0000 |
271 +-----+-----+

```

```

272 +-----+-----+
273 | supervisor | phi(supervisor) |
274 +-----+-----+
275 | supervisor_0 | 0.0000 |
276 | supervisor_1 | 0.2000 |
277 | supervisor_2 | 0.8000 |
278 +-----+-----+

```

279

280 **3 Sampling**

281 We draw samples from the Bayesian network so that we can better understand the data and
282 make statistical inferences on them.

283

284 **3.1 Bayesian Model Sampling**

285 Bayesian Model Sampling is available in pgmpy.sampling library as a class for sampling
286 methods specific to Bayesian Models. The model to be given as an input to the Sampling
287 function should be an instance of the Bayesian Model. Forward_sample function generates
288 samples from joint distribution of the Bayesian network, which we are using here as shown in
289 the below code snippet^[4]:

```
290 infer1 = BayesianModelSampling(Mental_health_model)
```

```
291 evidence1 = [State('treatment',1)]
```

```
292 sample1 = infer1.forward_sample(evidence1,5)
```

293 sample1 now contains a sample of the size 5 (rows) that correspond to the value 1 in
294 'treatment'. This generates a data frame from the dataset randomly based on the given
295 conditions.

296 Likelihood_weighted_sample generates weighted samples from joint distribution of the
297 Bayesian network that comply with the given evidence.

298 Sample code snippet:

```
299 infer1 = BayesianModelSampling(Mental_health_model)
```

```
300 evidence2 = [State('treatment',1)]
```

```
301 sample2 = infer1.likelihood_weighted_sample(evidence2,5)
```

302 As in the forward sampling, the data from the Mental health survey dataset is selected
303 randomly based on the evidence, although here a new column will be added to the data frame
304 sample2, namely weights.

305

306 **4 Evaluation Metrics**

307 We developed inference algorithms so far for determining the mean and entropy of each
308 distribution.

309

310 **4.1 Mean of a distribution**

311

312 The Mean for a distribution $p(x)$ is:

$$E[p(\mathbf{x})] = \sum_{\mathbf{x}} \mathbf{x}p(\mathbf{x})$$

313

314 Using N samples, the mean can be computed as:

$$\hat{E}[p(\mathbf{x})] = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$$

315
316 Numpy package in Python has predefined functions for computing mean.
317 Calculation of mean from the forward sample obtained above:

318
319 `np.mean(sample1)`

320

321 **Result:**

```
322 self_employed          0.0
323 coworkers              1.2
324 obs_consequence        0.0
325 Country                12.0
326 care_options           1.8
327 work_interfere         2.4
328 mental_vs_physical     1.0
329 Age                    1.2
330 Gender                 0.8
331 mental_health_consequence 1.0
332 mental_health_interview 1.0
333 seek_help              1.2
334 anonymity              1.2
335 supervisor             1.8
336 family_history         1.0
337 no_employees           2.6
338 tech_company           0.4
339 wellness_program       0.4
340 treatment              0.0
341 leave                  0.4
342 benefits               2.0
343 dtype: float64
```

344

345 **4.2 Entropy**

346 The entropy of a distribution $p(\mathbf{x})$ is:

$$H[p(\mathbf{x})] = - \sum_{\mathbf{x}} p(\mathbf{x}) \ln p(\mathbf{x})$$

347

348 When using N samples, the entropy can be calculated as:

$$\hat{H}[p(\mathbf{x})] = - \frac{1}{N} \sum_{k=1}^N \ln p(\mathbf{x}_k)$$

349

350 In order to calculate the entropy, the data frame containing the sample is converted into an
351 array in Python and the probabilities are computed for each cell column-wise. The entropy is
352 calculated using the entropy function in ‘Scipy’ package as follows:

353 `scipy.stats.entropy(s1)`

354

355


```
356 Result:
357 array([      -inf,  1.56071041,      -inf,  1.60943791,  1.58109375,
358         1.58902692,  1.05492017,  1.56071041,  1.38629436,  1.60943791,
359         1.60943791,  1.09861229,  1.09861229,  1.58109375,  1.60943791,
360         1.51938266,  0.69314718,  0.69314718,      -inf,  0.          ,
361         1.60943791])
```

362

363 References

- 364 [1] Dimitris Margaritis (2003) Learning Bayesian Network Model Structure from Data. Carnegie Mellon
365 University, PA.
- 366 [2] Jonathan S. Yedidia, William T. Freeman, Yair Weiss (2001) Understanding Belief Propagation and
367 its Generalizations. *Mitsubishi Electric Research Laboratories TR2001-22*
- 368 [3] <http://pgmpy.org/inference.html>
- 369 [4] <http://pgmpy.org/sampling.html>