
Term Deposit Subscription Prediction

Anmisha Reddy

Department of Computer Science Department of Computer Science
University at Buffalo University at Buffalo
anmishar@buffalo.edu suppu@buffalo.edu

Sindhura Uppu

Abstract

This project aims to predict if a client will subscribe to a term deposit in a bank. The data is related with direct marketing campaigns, phone calls, of a Portuguese banking institution. The goal is to achieve this prediction by generating a probabilistic graphical model or Markov Network to model the data and answer queries. The dataset is obtained from [here](#). Exact and Approximate inference algorithms have been applied on this network using which various queries have been answered that result in classifying which set of variables give accurate prediction.

1 Problem Domain

The dataset that is being modeled in this project is obtained from UCI Machine Learning Repository. The data is related with direct marketing campaigns of a Portuguese banking institution. A variety of variables are taken into consideration and their values play a role in the client's decision to either subscribe for a term deposit or not. This marketing campaign is carried out through phone calls. The major goal of this project is to see which set of variables give accurate prediction by observing the inference answers of the training data.

1.1 Data Set

The data set is the collection of the details of the clients of a Portuguese banking institution. The dataset contains around 4522 instances and 21 attributes. These attributes help in analyzing the state of a client by identifying whether he/she has taken a personal or home loan, their educational status and whether anyone from the bank has previously called them for asking about subscription and their decision whether to subscribe or not. The dataset contained many 'string' type variables, which are converted to categories, if found to be discrete. The following are the descriptions of the variables and what they each represent

Age: age of the client.

Job: type of job – admin, blue-collar, entrepreneur, housemaid, management, retired, self-employed, services, student, technician, unemployed, unknown.

Marital: marital status of the client – divorced, married, single, unknown.

Education: how much education does the client have – basic.4y, basic.6y, basic.9y, high.school, illiterate, professional.course, university.degree, unknown.

Default: Does the client have credit in default? – yes or no.

Housing: Did the client take any housing loan in the past? – yes or no.

Loan: Did the client take any personal loan in the past? – yes or no.

48 Contact: Contact communication type – cellular or telephone.
 49 Month: last contact month of the year – jan, feb, mar, apr, may, jun, jul, aug, sep, oct, nov,
 50 dec.
 51 Day_of_week: last contact day of the week – mon, tue, wed, thu, fri, sat.
 52 Duration: last contact duration, in seconds.
 53 Campaign: number of contacts performed during this campaign and for this client.
 54 Pdays: number of days that passed by after the client was last contacted from a previous
 55 campaign.
 56 Previous: number of contacts performed before this campaign and for this client.
 57 Poutcome: outcome of the previous marketing campaign – failure, nonexistent, success.
 58 Emp.var.rate: employment variation rate.
 59 Cons.price.idx: consumer price index.
 60 Cons.conf.idx: consumer confidence index.
 61 Euribor3m: euribor 3 month rate.
 62 Nr.employed: number of employees.
 63 Y: Has the client subscribed a term deposit? (Yes or No).
 64 Y is the output variable in this dataset that classifies whether the client has subscribed
 65 for a term deposit or not.

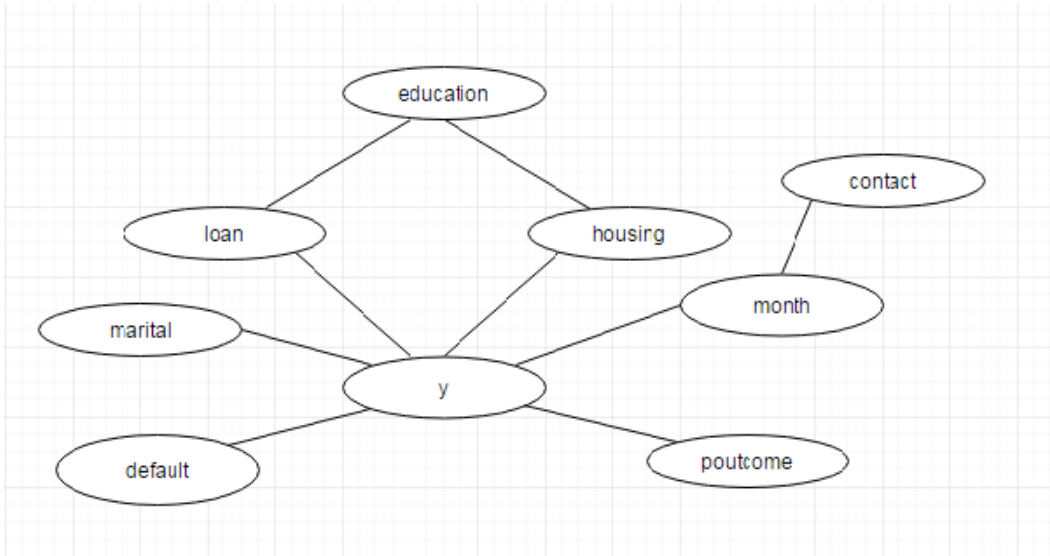
66 67 **1.2 Markov Network Model**

68 Markov models are probabilistic graphical models that are undirected. They represent the
 69 joint probability distribution over events which are represented by variables. The nodes
 70 in the graph represent variables and the edges correspond to a notion of direct probabilistic
 71 interaction between the neighboring variables – an interaction that is not mediated by any
 72 other variable in the network^[1]. Factors are associated with each complete subgraph in the
 73 network. They are usually non negative and do not necessarily represent probabilities. The
 74 joint distribution is the product of distribution of factors together. The factors can be
 75 normalized by taking the entire sum and dividing each factor with this total.

76 The variables of the data set were assessed and relations were drawn. The Markov
 77 Network consists of 9 variables with 9 edges between them. All these variables are
 78 considered as categories and the values are changed to numeric codes that are given to
 79 categories. These values are given in ascending order to the distinct values in each
 80 variable, based on their first letter. After performing this conversion, some of the values
 81 of the variables are as follows:

82 Default: 0 – no, 1 – yes.
 83 Marital: 0 – divorced, 1 – married, 2 – single.
 84 Contact: 0 – cellular, 1 – telephone, 3 – unknown.
 85 Housing: 0 – no, 1 – yes.
 86 Education: 0 – primary, 1 – secondary, 2 – tertiary, 3 – unknown.
 87 Loan: 0 – no, 1 – yes.
 88 Poutcome: 0 – failure, 1 – other, 2 – success, 3 – unknown.
 89 Month: 0 – apr, 1 – aug, 2 – dec, 3 – feb, 4 – jan, 5 – jul, 6 – jun, 7 – mar, 8 – may, 9 –
 90 nov, 10 – oct, 11 – sep.
 91 Y: 0 – no, 1 – yes.

92 The Markov network that we represented below was constructed manually based on our
 93 intuition.
 94



95
 96

97 The Markov network can be generated in Python as follows^[2]:

```

98 mark=
99 MarkovModel([("education","loan"),('education','housing'),('loan','y'),('housing','y'),(
100 'marital','y'),('default','y'),('contact','month'),('month','y'),('poutcome','y')])

```

101

102 Pgmpy library has inbuilt models including Bayesian Model and Markov Model. The
 103 nodes that have edges are sent as parameters using the node names as attributes in the
 104 model.

105

106 1.3 Factors

107 It is not possible to fit a Markov model with our existing data set. Hence, we can use
 108 factors to represent the joint probability distribution of the Markov model. Factors help
 109 parameterize the model. The factors are not probabilities or conditional probabilities
 110 between two variables, rather they seem to have values given by a user's intuition^[1]. This
 111 means that they are significantly harder to estimate from the data.

112 We aim to initialize the factors in such a way that the accuracy of inference algorithms is
 113 high. To this extent, we took the co-occurrence of each subset of variables in the Markov
 114 network and gave these values as discrete factors to the model. Thus, by doing this we
 115 obtained around 9 factors that are added to the Markov model declared as above. Our
 116 representation is complete as the graph structure is associated with a set of parameters.

117 Code snippet that shows how we obtained the factor values and declared them using
 118 DiscreteFactors in pgmpy library^[2].

```

119 data.groupby(["education","loan"]).size()

```

120 Results:

```

121 education  loan
122 0          0      584
123          1      94

```

```

124      1          0      1890
125          1      416
126      2          0      1176
127          1      174
128      3          0      180
129          1          7
130      dtype: int64
131
132      fl=
133      DiscreteFactor(["education","loan"],cardinality=[4,2],values=(584,94,1898,416,1176,17
134      4,180,7))
135
136      Groupby() function gives the count of co-occurrences of each subset of variables given,
137      which are education and loan. DiscreteFactor is a factor type present in
138      pgmpy.factors.discrete and it is used to declare discrete factors for each set of variables.
139      Cardinality is the size of unique elements in each variable and the values must be a list of
140      each co-occurrence in sorted order. In a similar fashion, the 9 discrete factors are declared
141      and these factors are added to the Markov network as shown below:
142
143      mark.add_factors(fl)
144
145      Local independencies are obtained from the model using get_local_independencies()
146      function.
147
148      mark.get_local_independencies()
149
150      Results:
151      (poutcome _|_ default, loan, month, contact, education, mar
152      ital, housing | y)
153      (default _|_ poutcome, month, loan, marital, contact, educa
154      tion, housing | y)
155      (loan _|_ poutcome, default, housing, month, contact, marit
156      al | y, education)
157      (marital _|_ poutcome, default, loan, month, contact, educa
158      tion, housing | y)
159      (contact _|_ poutcome, default, loan, marital, y, education
160      , housing | month)
161      (month _|_ poutcome, default, loan, marital, education, hou
162      sing | y, contact)
163      (y _|_ contact, education | poutcome, default, loan, month,
164      marital, housing)
165      (education _|_ poutcome, default, month, contact, y, marita
166      l | loan, housing)
167      (housing _|_ poutcome, default, loan, month, contact, marit
168      al | y, education)
169
170
171
172      The factors in a Markov model can be retrieved using the function get_factors() as
173      below:
174

```

```

175
176         for fact in mark.get_factors():
177             print(fact)
178
179 Results:
180 +-----+-----+-----+
181 | education | housing | phi(education,housing) |
182 |-----+-----+-----|
183 | education_0 | housing_0 | 295.0000 |
184 | education_0 | housing_1 | 383.0000 |
185 | education_1 | housing_0 | 876.0000 |
186 | education_1 | housing_1 | 1430.0000 |
187 | education_2 | housing_0 | 687.0000 |
188 | education_2 | housing_1 | 663.0000 |
189 | education_3 | housing_0 | 104.0000 |
190 | education_3 | housing_1 | 83.0000 |
191 +-----+-----+-----+
192 +-----+-----+-----+
193 | loan | y | phi(loan,y) |
194 |-----+-----+-----|
195 | loan_0 | y_0 | 3352.0000 |
196 | loan_0 | y_1 | 478.0000 |
197 | loan_1 | y_0 | 648.0000 |
198 | loan_1 | y_1 | 43.0000 |
199 +-----+-----+-----+

```

200

201

202 Inference Models

203

204 2.1 Belief Propagation

205 Computing the a posteriori belief of a variable in a general Markov Network is NP-hard.
 206 Belief Propagation is an Approximate Inference algorithm. It is an efficient way to solve
 207 inference problems by passing local messages. It is available in the pgmpy.inference
 208 library as a class for performing inference using BeliefPropagation model. It creates a
 209 junction tree or Clique tree for the input probabilistic graphical model and performs
 210 calibration of the junction tree so formed using belief propagation. Thus, we trained our
 211 model using belief propagation as follows^[3]:

```
212 belief_prop = BeliefPropagation(mark)
```

213

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

217

218 Query1:

```
219 bp1 = belief_prop.query(variables=['y'],evidence={'marital' : 0,'default' : 0})
```

```
220 print(bp1['y'])
```

221 We aim to find the variation of y, whether the client opts to subscribe or not, based on
 222 his/her marital status (divorced) and default value (no).

223 Results:

```
224 +-----+-----+
```

```

225 | y | phi (y) |
226 |-----+-----|
227 | y_0 | 1.0000 |
228 | y_1 | 0.0000 |
229 +-----+-----+

```

230

231 Query2:

```
232 bp2 = belief_prop.query(variables=['y'],evidence={'education' : 2,'housing' : 1})
```

```
233 print(bp2['y'])
```

234 We aim to find the variation of y, whether the client opts to subscribe or not, based
235 on his/her education (tertiary) and when the client took housing loan.

236 Results:

```

237 +-----+-----+
238 | y | phi (y) |
239 |-----+-----|
240 | y_0 | 1.0000 |
241 | y_1 | 0.0000 |
242 +-----+-----+

```

243

244 Query3:

```
245 bp3 = belief_prop.query(variables=['y'],evidence={'month' : 11,'poutcome' : 2})
```

```
246 print(bp3['y'])
```

247 We aim to find the variation of y, whether the client opts to subscribe or not, based
248 on the month they were last contacted (oct) and when the previous outcome is success.

249 Results:

```

250 +-----+-----+
251 | y | phi (y) |
252 |-----+-----|
253 | y_0 | 0.9997 |
254 | y_1 | 0.0003 |
255 +-----+-----+

```

256

257 Query4:

```
258 bp4 = belief_prop.query(variables=['y'],evidence={'poutcome' : 2, 'loan' : 1})
```

```
259 print(bp4['y'])
```

260 We aim to find the variation of y, whether the client opts to subscribe or not, based
261 on the case where the client has taken a personal loan and previous outcome for term
262 deposit subscription is a success.

263 Results:

```

264 +-----+-----+
265 | y | phi (y) |
266 |-----+-----|
267 | y_0 | 1.0000 |
268 | y_1 | 0.0000 |
269 +-----+-----+

```

270

```

271 Query5:
272 Bp5 = belief_prop.query(variables=['y'],evidence={'marital' : 1, 'loan' : 1, 'contact' : 1,
273 'month' : 5})
274 print(bp5['y'])

```

275 We aim to find the variation of y, whether the client opts to subscribe or not, based
 276 on several factors such as the client is married, has taken loan, contact is through telephone
 277 and last contacted month is July.

```

278 Results:
279 +-----+-----+
280 | y      | phi (y) |
281 |-----+-----|
282 | y_0    | 1.0000  |
283 | y_1    | 0.0000  |
284 +-----+-----+
285

```

286 2.2 Variable Elimination

287
 288 Variable Elimination is an exact inference algorithm. It is available in the pgmpy.inference
 289 library as a class for performing inference^[3].

```

290
291 from pgmpy.inference import VariableElimination
292 infer = VariableElimination(mark)
293

```

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

```

297 Query1:
298 infer = VariableElimination(mark)
299 phi_query = infer.query(variables=['y'],evidence={'marital':0,'default':0})
300 print(bp1['y'])

```

```

301 Result:
302 +-----+-----+
303 | y      | phi (y) |
304 |-----+-----|
305 | y_0    | 1.0000  |
306 | y_1    | 0.0000  |
307 +-----+-----+
308

```

```

309 Query2:
310 infer = VariableElimination(mark)
311 phi_query = infer.query(variables=['y'],evidence={'education':2,'housing':1})
312 print(bp1['y'])

```

```

313 Result:
314 +-----+-----+
315 | y      | phi (y) |
316 |-----+-----|
317 | y_0    | 1.0000  |
318 | y_1    | 0.0000  |
319 +-----+-----+

```

```

320 Query3:
321 infer = VariableElimination(mark)
322 phi_query = infer.query(variables=['y'],evidence={'month':10,'poutcome':2})
323 print(bp1['y'])

```

```

324 Result:
325 +-----+-----+
326 | y      | phi(y) |
327 |-----+-----|
328 | y_0    | 1.0000 |
329 | y_1    | 0.0000 |
330 +-----+-----+

```

```

331
332 2.3 MPLP
333

```

MPLP is a class for performing approximate inference using Max-Product Linear Programming method^[3]. Markov model can be trained using this inference algorithm. The following code snippet shows how we have trained our model using MPLP inference algorithm.

```

338
339     from pgmpy.inference import Mplp
340     mplp = Mplp(mark)

```

```

341

```

```

342

```

3 Sampling

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

```

346

```

3.1 Gibbs Sampling

A distribution P_Φ is a Gibbs distribution parameterized by a set of factors $\Phi = \{\Phi_1(D_1), \dots, \Phi_k(D_k)\}$ if it is defined as follows:

$$P_\Phi(X_1, \dots, X_n) = 1 / Z (P_\Phi(X_1, \dots, X_n))$$

```

351

```

GibbsSampling is a class present in pgmpy.sampling. It generates samples from models from which variables are inherited and transition probabilities computed^[4]. The following code snippet shows how we obtained samples from Markov network using GibbsSampling.

```

356     from pgmpy.sampling import GibbsSampling
357     gibbs_chain = GibbsSampling(mark)
358     gibbs_chain.sample(size=5)

```

This code generates samples of size 5 from the model.

```

360

```

4 Evaluation Metrics

So far, we developed inference algorithms for determining the mean and entropy of each distribution.

```

364

```

4.1 Mean of a distribution

The mean for a distribution $P(x)$ is:

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

367

368 Using N samples, the mean can be computed as

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

369

370 Numpy package in python has predefined functions for computing mean.

371 Calculation of mean from the sample obtained:

372 `np.mean(df)`

373 Result:

```
374 marital      1.2
375 education    1.2
376 default      0.0
377 housing      0.8
378 loan         0.4
379 contact      0.8
380 month        6.4
381 poutcome     1.8
382 y            0.0
383 dtype: float64
```

384

385 4.2 Entropy

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

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

387

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

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

389

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

393 `scipy.stats.entropy(s1)`

394 Results:

```
395 array([ 1.56071041,  1.32966135,          -inf,  1.38629436,
396        0.69314718,
397        0.69314718,  1.37050239,  1.09861229,          -inf])
```

398

399
400

References

- 401 [1] Probabilistic Graphical Models Principles and Techniques, Daphne Koller and Nir
402 Friedman. The MIT Press, Cambridge, Massachusetts.
- 403 [2] <http://pgmpy.org/models.html>
- 404 [3] <http://pgmpy.org/inference.html>
- 405 [4] <http://pgmpy.org/sampling.html>