

Project 1: Bayesian Structure Learning

Problem Statement

This project implements an algorithm to find a (somewhat) bayesian-optimal directed acyclic graph given some data.

Algorithm Description

I chose to use *opportunistic local search* for my graph search algorithm. This algorithm works by randomly choosing an operation - add, remove, or flip an edge - and then applying it to the current graph to generate a new graph. This new graph is kept and iterated upon if its Bayesian score is better than the original graph's score. The algorithm concludes when the overall best score has not been achieved in `maxAttempts` consecutive graphs.

This method of hill climbing is susceptible to getting stuck in local minima - to prevent this, there are a few methods used to introduce random resets into the algorithm. First, after a new graph is generated, even if its Bayesian score is worse than the original graph, I assign a 5% chance to keeping the new (worse) graph. Second, if the algorithm seems to get stuck and fails to find any improvement after `maxAttempts/20` graphs, I revert the most previous change made to the graph. Last, if there has not been an improvement on the overall best score `maxAttempts/10` new graphs, I reset the graph back to its state where the best overall score was achieved. For optimizing the graphs, I set `maxAttempts` with a goal of finishing optimization in approximately 10 to 20 minutes.

Pseudocode for the algorithm is shown in **Appendix 1**.

The runtime of these algorithms is in **Figure 1**. The runtime of these algorithms is summarized in **Table 1**.

Table 1. Summary of algorithm performance with various data dimensionality. Note that the 20 minute time limit was reached for the 50-variable dataset.

Number of variables	maxAttempts	Optimized bayesian score	Runtime (mins)	Graphs generated and scored per minute
8	1000	-3796.86	11.63	223
12	300	-42051.16	19.36	33
50	300	-483860.92	20	10

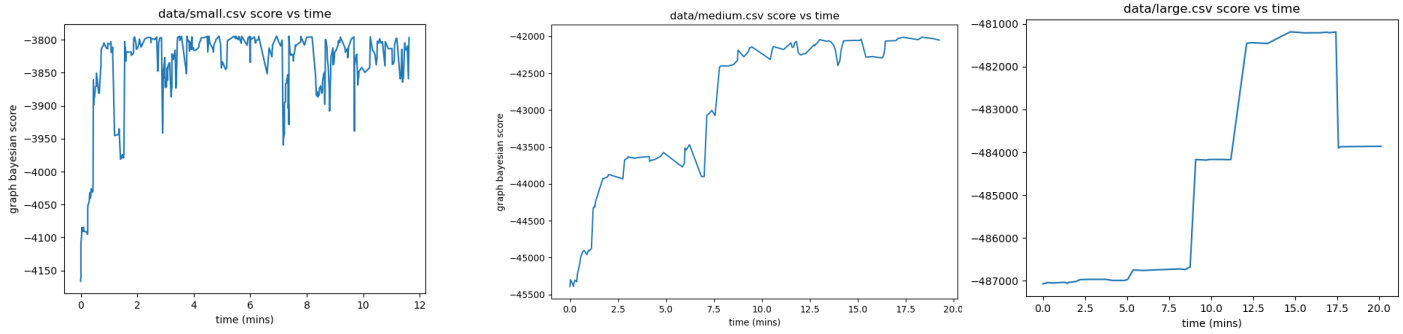


Figure 1. Plot of score versus runtime for the small (left), medium, and large (right) datasets. The algorithm struggled with the large dataset, as seen by the long duration in between data points in the rightmost plot.

Graph Plots

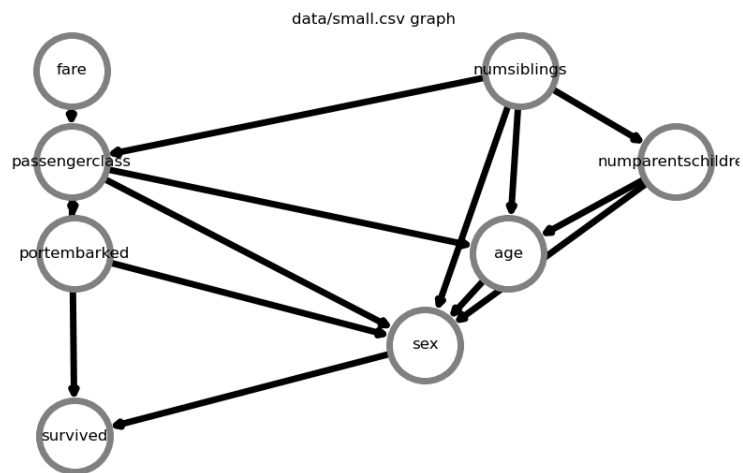


Figure 2. Graph from “small” dataset.

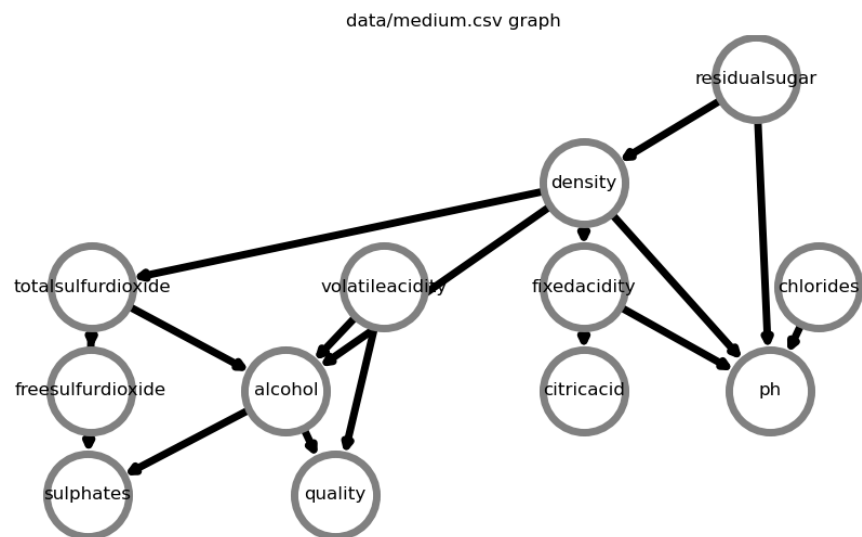


Figure 3. Graph from “medium” dataset. Note that “totalsulfurdioxide” and “freesulfurdioxide” are both parents of “sulphates” (this seems to have gotten overlapped in the plot).

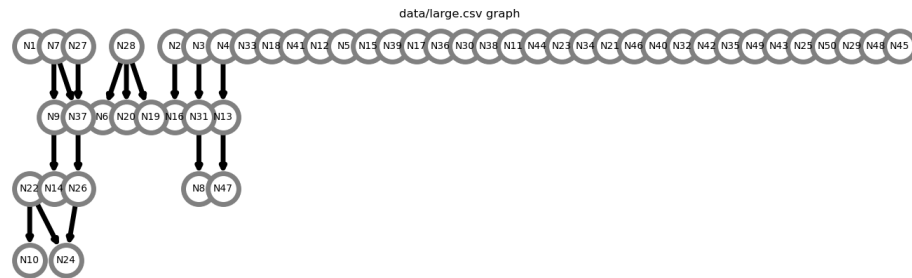


Figure 4. Graph from the “large” dataset. It seems like, even with 20 minutes of optimization, more time would have made this graph more representative of the data.

Appendix 1: Algorithm Pseudocode

```
Generate graph G
Set allTimeBest_graph to G
Set allTimeBest_score and score to bayesScore(G)

While True:
    Randomly add an edge, remove an edge, or flip an edge to generate G_new
    Determine score_new of G_new
    If score_new > allTimeBest_score:
        Update allTimeBest_score to score_new
        Update allTimeBest_graph to G_new
        Set countWithoutAllTimeImprovement to 0
    Else:
        Increment countWithoutAllTimeImprovement
    If score_new > score or randint(0,20) == 0:
        Update score to score_new
        Update G to G_new
        Set countWithoutImprovement to 0
    Else:
        Increment countWithoutImprovement

    If countWithoutImprovement > maxAttempts/10:
        Undo last change to G
    If countWithoutAllTimeImprovement > maxAttempts/5:
        Reset G to allTimeBest_graph
    If countWithoutAllTimeImprovement > maxAttempts:
        Break

Save allTimeBest_G and allTimeBest_score
Plot
```

```

1  import sys
2  import numpy as np
3  import networkx as nx
4  import pandas
5  import matplotlib.pyplot as plt
6  import scipy.special
7  import copy
8  from datetime import datetime
9  from networkx.drawing.nx_pydot import graphviz_layout
10
11 def write_gph(dag, idx2names, filename):
12     with open(filename, 'w') as f:
13         for edge in dag.edges():
14             f.write("{}\n".format(idx2names[edge[0]], idx2names[edge[1]]))
15
16 def lg(x):
17     return scipy.special.loggamma(x)
18
19 def bayes_score_component(M, alpha):
20
21     p = np.sum(lg(alpha + M))
22     p = p - np.sum(lg(alpha))
23     p = p + np.sum(lg(np.sum(alpha, axis=1)))
24     p = p - np.sum(lg(np.sum(alpha, axis=1) + np.sum(M, axis=1)))
25
26     return p
27
28 def getIdxFromParentState(r, idxOfParents, stateOfParents):
29     # r = num instantiations of each variable
30     returnIdx = 0
31
32     nStatesForEachParent = [r[p] for p in idxOfParents]
33     q = np.prod(nStatesForEachParent) # total number of parent states
34     idxOfAllParentStates = range(q)
35     idxOfAllParentStates = np.reshape(idxOfAllParentStates, (nStatesForEachParent))
36     return idxOfAllParentStates[tuple(stateOfParents)]
37
38 def counts(vars, G, D, graphInfo):
39     # D = data of shape(n,m) where n is num variables and m is num datapts
40     # G = bayes struct
41     # returns array M of shape(n,)
42
43     # n = num variables
44     n = np.shape(D)[0]
45     m = np.shape(D)[1]
46
47     # r = num instantiations of var[i]
48     r = graphInfo['r']
49     # q = num instantiations of var[i]'s parents
50     q = graphInfo['q']
51     # m = list of size n, each element containing a shape(q[i], r[i]) containing counts
52     # of each state
53     M = [np.zeros((q[i], r[i])) for i in range(n)]
54
55     for datapoint in range(m):
56         for datapoint_item in range(n):
57             state_of_self = D[datapoint_item, datapoint]
58             parents = list(G.predecessors(datapoint_item))
59             state_of_parents_idx = 0
60             if len(parents) != 0:
61                 state_of_each_parent = [D[parent, datapoint] for parent in parents]
62                 state_of_parents_idx = getIdxFromParentState(r, parents,
63                                                             state_of_each_parent)
64
65             M[datapoint_item][state_of_parents_idx, state_of_self] += 1
66     return M

```

```

66
67 def prior(vars, graphInfo):
68     n = len(vars)
69     r = graphInfo['r']
70     q = graphInfo['q']
71
72     # m = list of size n, each element containing a shape(q[i], r[i]) containing counts
73     # of each state
74     M = [np.ones((q[i], r[i])) for i in range(n)]
75     return M
76
77
78 def bayes_score(vars, G, D, graphInfo):
79     n = np.shape(vars)[0]
80     M = counts(vars, G, D, graphInfo)
81     alpha = prior(vars, graphInfo)
82     return np.sum([bayes_score_component(M[i], alpha[i]) for i in range(n)])
83
84 def getGraphInfo(n, G, D):
85     # r = num instantiations of var[i]
86     r = [np.max(D[i])+1 for i in range(n)]
87     # q = num instantiations of var[i]'s parents
88     q = [int(np.prod([r[j] for j in list(G.predecessors(i))])) for i in range(n)]
89
90     returnVal = {'r': r, 'q': q}
91     return returnVal
92
93 def compute(infile, outfile):
94     startt=datetime.now()
95
96     # (1) read in data
97     D = pandas.read_csv(infile)
98     vars = D.columns
99     D = np.array(D).T - 1 # subtract 1 for 0-indexing
100    n = len(vars)
101
102    # (2) get a graph, set nodes to 0, 1, ..., n
103    G = nx.DiGraph()
104    G.add_nodes_from(range(n))
105    graphInfo = getGraphInfo(n, G, D)
106
107    # (3) score the graph
108    score = bayes_score(vars, G, D, graphInfo)
109    allTimeBest = {'G': G, 'score': score}
110
111    # (4) optimize the graph using local directed graph search
112    maxTime = 20 # mins
113    maxAttempts = 1000
114    undoRate = np.floor(maxAttempts/10)
115    resetRate = np.floor(maxAttempts/5)
116    consecutiveNewGraphsWithoutImprovement = 0
117    consecutiveNewGraphsWithoutAllTimeImprovement = 0
118    totalGraphsGenerated = 0
119    scores = np.array([[0, score]])
120    mostRecentChange = {'type': None, 'nodes': None}
121    while True:
122
123        runtime_mins = np.round((datetime.now() - startt).total_seconds()/60, 2)
124        G_new = copy.deepcopy(G)
125        existingEdges = list(G_new.edges)
126        # (4.1) add an edge
127        decision = np.random.randint(0,3)
128        if decision == 0 or len(existingEdges) == 0:
129
130            for i in range(maxAttempts):
131                # pick start and end node, make sure they're not the same and don't

```

```

132         already exist
133     for j in range(maxAttempts):
134         startNode = np.random.randint(0, n)
135         endNode = np.random.randint(0, n)
136         if startNode != endNode and (startNode, endNode) not in list(G.edges)
137         and (endNode, startNode) not in list(G.edges):
138             break
139
140     # add to graph. If we still have a DAG, we can go get the score
141     G_new.add_edge(startNode, endNode)
142
143     if nx.is_directed_acyclic_graph(G_new):
144         break
145
146 # (4.2) alternatively, remove or flip and existing edge
147 else:
148     for i in range(maxAttempts):
149         existingEdges = list(G_new.edges)
150         edgeToChange = existingEdges[np.random.randint(0, len(existingEdges))]
151         startNode = edgeToChange[0]
152         endNode = edgeToChange[1]
153         if decision == 1: # remove
154             G_new.remove_edge(startNode, endNode)
155         else: # flip
156             G_new.remove_edge(startNode, endNode)
157             G_new.add_edge(endNode, startNode)
158         if nx.is_directed_acyclic_graph(G_new):
159             break
160
161 # (4.3) get the score of G_new, and update G if the score is an improvement
162 if nx.is_directed_acyclic_graph(G_new):
163     graphInfo_new = getGraphInfo(n, G_new, D)
164     score_new = bayes_score(vars, G_new, D, graphInfo_new)
165     totalGraphsGenerated += 1
166
167     if score_new > score or np.random.randint(0, 20) == 0:
168         if score_new > score:
169             consecutiveNewGraphsWithoutImprovement = 0
170         if score_new > allTimeBest['score']:
171             improvement = round(score_new - allTimeBest['score'], 2)
172             print(f'New all time best score! {round(score_new, 2)} (change: {
173                 improvement})')
174             allTimeBest['score'] = score_new
175             allTimeBest['G'] = G_new
176             consecutiveNewGraphsWithoutAllTimeImprovement = 0
177
178     G = copy.deepcopy(G_new)
179     score = copy.deepcopy(score_new)
180     scores = np.append(scores, np.array([[runtime_mins, score]]), axis=0)
181     mostRecentChange['type'] = decision
182     mostRecentChange['nodes'] = [startNode, endNode]
183     print(f'Improved score {round(score, 2)}. Runtime = {runtime_mins}')
184 else:
185     consecutiveNewGraphsWithoutImprovement += 1
186     consecutiveNewGraphsWithoutAllTimeImprovement += 1
187
188 else:
189     consecutiveNewGraphsWithoutImprovement += 1
190     consecutiveNewGraphsWithoutAllTimeImprovement += 1
191
192 # (4.4) print status
193 if np.mod(consecutiveNewGraphsWithoutAllTimeImprovement, 100) == 0 and
194 consecutiveNewGraphsWithoutAllTimeImprovement > 99:
195     print(f'consecutive graphs without all-time score improvement = {
196         consecutiveNewGraphsWithoutAllTimeImprovement}' +
197         f' runtime = {runtime_mins}')

```

```

194     # (4.5) if no changes seem to help, undo the last change
195     if np.mod(consecutiveNewGraphsWithoutImprovement, undoRate) == undoRate-1:
196         print(f'No changes are helping - removing last change')
197         if mostRecentChange['type'] == 0: # undo an add
198             try:
199                 G.remove_edge(mostRecentChange['nodes'][0], mostRecentChange['nodes']
200                               ][1])
201             except:
202                 print(f'error, cannot remove nodes from graph')
203         elif mostRecentChange['type'] == 1: # undo a remove
204             try:
205                 G.add_edge(mostRecentChange['nodes'][0], mostRecentChange['nodes'][1
206                               ])
207             except:
208                 print(f'error, cannot add nodes to graph')
209         elif mostRecentChange['type'] == 2: # undo a flip
210             try:
211                 G.remove_edge(mostRecentChange['nodes'][1], mostRecentChange['nodes']
212                               )[0])
213                 G.add_edge(mostRecentChange['nodes'][0], mostRecentChange['nodes'][1
214                               ])
215             except:
216                 print(f'error, cannot flip nodes')
217
218     if np.mod(consecutiveNewGraphsWithoutAllTimeImprovement, resetRate) == resetRate-
219     1:
220         print(f'{resetRate} consecutive graphs w/o all-time improvement - resetting
221         back to all-time G')
222         G = copy.deepcopy(allTimeBest['G'])
223
224     # (4.6) if we've had too many consecutive duds, we're done
225     if consecutiveNewGraphsWithoutImprovement > maxAttempts/10 or
226     consecutiveNewGraphsWithoutAllTimeImprovement > maxAttempts or runtime_mins >
227     maxTime:
228         break
229
230     # (5) return the graph when done
231     network_str = ''
232     G = copy.deepcopy(allTimeBest['G'])
233
234     for i in range(len(vars)):
235         parentName = vars[i]
236         kidsIdxs = list(G.successors(i))
237         if len(kidsIdxs) != 0:
238             for k in kidsIdxs:
239                 strToAdd = f'{parentName}, {vars[k]} \n'
240                 network_str += strToAdd
241         else:
242             network_str += f'{parentName}, \n'
243     graphsPerMin = round(totalGraphsGenerated/runtime_mins, 2)
244     print(f'generated {totalGraphsGenerated} graphs in {runtime_mins} mins ({graphsPerMin
245     } graphs/min), best score = {score}')
246     print(network_str)
247
248     if 1:
249         plt.figure()
250         plt.plot(scores[:,0], scores[:,1])
251         plt.xlabel('time (mins)')
252         plt.ylabel('graph bayesian score')
253         plt.title(f'{infile} score vs time')
254
255         nodeMapping = dict(zip(list(range(len(vars))), list(vars)))
256         G = nx.relabel_nodes(G, nodeMapping)
257         options = {
258             "font_size": 10,

```



```

252         "node_size": 1000,
253         "node_color": "white",
254         "edgecolors": "gray",
255         "linewidths": 5,
256         "width": 5,
257     }
258     plt.figure()
259     # pos = nx.shell_layout(G)
260     pos = graphviz_layout(G, prog='dot')
261     nx.draw_networkx(G, pos=pos, arrows=True, arrowstyle='->', **options)
262     plt.title(f'{infile} graph')
263     plt.axis("off")
264     plt.show()
265     print('done')
266
267 def main():
268
269     if len(sys.argv) != 3:
270         raise Exception("usage: python project1.py <infile>.csv <outfile>.gph")
271
272     inputfilename = sys.argv[1]
273     outputfilename = sys.argv[2]
274     compute(inputfilename, outputfilename)
275
276
277 if __name__ == '__main__':
278     if 0:
279         main()
280     else:
281         inputfilename = "data/large.csv"
282         outputfilename = "graphs/large.gph"
283         compute(inputfilename, outputfilename)
284

```