

Building Belief Networks with pgmpy

In this notebook we will use the `pgmpy` [python package \(http://pgmpy.org\)](http://pgmpy.org) to build belief networks. See the website for detailed installation instructions, but basically:

```
$ pip3 install -r requirements.txt
$ pip3 install pgmpy
```

where the `requirements.txt` file (obtained from the github source [here](https://github.com/pgmpy/pgmpy/blob/dev/requirements.txt) (<https://github.com/pgmpy/pgmpy/blob/dev/requirements.txt>)) is a list of the `pgmpy` package dependencies that must also be installed.

See [the docs \(http://pgmpy.org\)](http://pgmpy.org) for reference or [these notebooks](https://github.com/pgmpy/pgmpy_notebook/tree/master/notebooks) (https://github.com/pgmpy/pgmpy_notebook/tree/master/notebooks) and [these examples](https://github.com/pgmpy/pgmpy/tree/dev/examples) (<https://github.com/pgmpy/pgmpy/tree/dev/examples>) from the `pgmpy` github site (although several aren't finished). We will cover the basics here.

Example 1: Inspector Clouseau

Here we re-do the Inspector Clouseau example from Barber, but this time by constructing a belief network.

Barber Example 1.3 (Inspector Clouseau)

Inspector Clouseau arrives at the scene of a crime. The victim lies dead in the room and the inspector quickly finds the murder weapon, a Knife (κ). The Butler (B) and Maid (M) are his main suspects. The inspector has a prior belief of 0.8 that the Butler is the murderer, and a prior belief of 0.2 that the Maid is the murderer. These probabilities are independent in the sense that $p(B, M) = p(B)p(M)$. (It is possible that both the Butler and the Maid murdered the victim or neither). The inspector's *prior* criminal knowledge can be formulated mathematically as follows.

(Here we use short-hand notation, e.g. $p(B)$ means the (*a priori*) probability that the butler is the murder, $p(K|B, \neg M)$ means the probability that the knife was used given that the butler was the murder and the maid was not.)

$$\begin{aligned}p(B) &= 0.6 \\p(M) &= 0.2 \\p(M, B) &= p(M)p(B) \\p(K|\neg B, \neg M) &= 0.3 \\p(K|\neg B, M) &= 0.2 \\p(K|B, \neg M) &= 0.6 \\p(K|B, M) &= 0.1\end{aligned}$$

Assuming the knife is the murder weapon, what is the probability that the Butler is the murderer?

We can solve this by creating a Bayesian network using the `pgmpy` package. First, we import it:

In [1]:

```
1 from pgmpy.models import BayesianModel
2 from pgmpy.factors.discrete import TabularCPD
```

Build Up Network

In `pgmpy` , a belief network is specified by variables and their conditional probabilities. To build the model, you describe the structure of network by defining nodes and links in a directed graph, and then assign the conditional probabilities (or priors) for each variable.

In [10]:

```
1 model = BayesianModel([( 'B', 'K'), ( 'M', 'K')])
2
3 # define p(B) and p(M)
4 # variable_card is cardinality = 2 for true/false
5 # values are defined in numeric order p(x_i = [false, true]), ie [0, 1]
6 priorB = TabularCPD(variable='B', variable_card=2, values=[[0.4, 0.6]])
7 priorM = TabularCPD(variable='M', variable_card=2, values=[[0.8, 0.2]])
8
9 # define p(K|B,M)
10 # Variables cycle in numerical order of evidence values,
11 # ie BM = 00, 01, 10, 11 for each value of K.
12 cpdK = TabularCPD(variable='K', variable_card=2,
13     evidence=[ 'B', 'M'], evidence_card=[2, 2],
14     values=[[0.7, 0.8, 0.4, 0.9],
15            [0.3, 0.2, 0.6, 0.1]])
16
17 # add probabilities to model
18 model.add_cpds(priorB, priorM, cpdK)
```

The `check_model` function check structure and values consistency. For example, it will give an error if the probabilities don't sum to one.

In [11]:

```
1 model.check_model()
```

Out[11]:

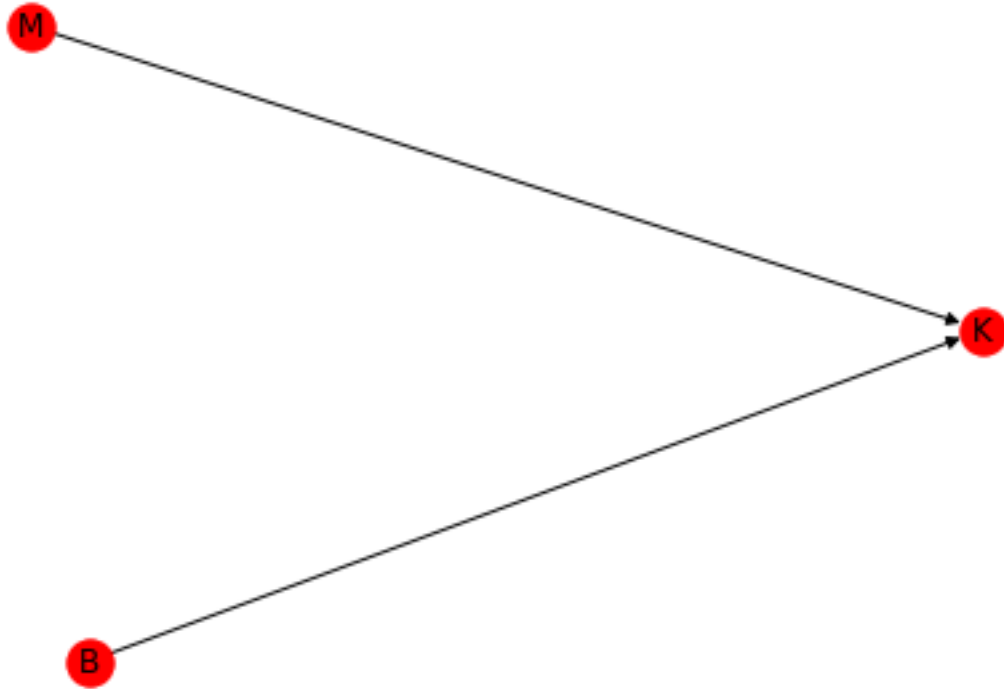
True

Two methods for plotting

1. Using `networkx` (<https://networkx.github.io/documentation/networkx-1.9/index.html>) and `pyplot` (https://matplotlib.org/api/pyplot_api.html) packages. This plots directly from model.

In [13]:

```
1 # TODO: Improve the layout and style
2 import networkx as nx
3 import matplotlib.pyplot as plt
4 nx.draw(model, with_labels=True)
5 plt.show()
```

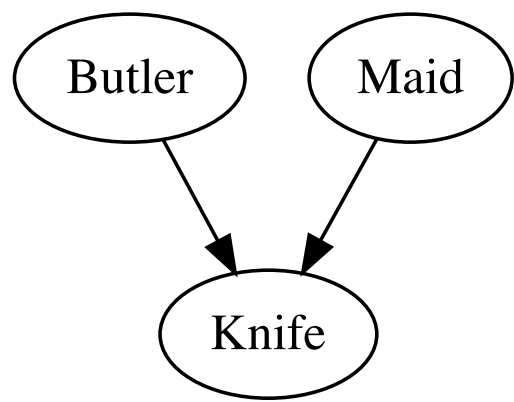


2. Use the [graphviz \(https://github.com/xflr6/graphviz\)](https://github.com/xflr6/graphviz) package, which produces a nice layout by default, but you have to define the graph separately.

In [14]:

```
1  # TODO: Can we plot with graphviz using model directly?
2
3  from graphviz import Digraph
4
5  dot = Digraph()
6  dot.node('B', 'Butler') # variable name, label
7  dot.node('M', 'Maid')
8  dot.node('K', 'Knife')
9  dot.edges(['BK', 'MK'])
10 # render inline
11 dot
12
13 # Or print to a pdf:
14 #print(dot.source)
15 #dot.render("Clouseau", view=True)
```

Out[14]:



Displaying the conditional probability tables

We can also display the conditional probability tables, which we can verify match the definitions above. Note that in `pgmpy`, the rows are the states of the variable and and columns are the states of evidence variables. $P(K|B, M)$ is shown below. Note that the variables of the CDF are in the rows and the conditioned variables are in the columns.

In [15]:

```
1  print(model.get_cpds('K'))
```

	B	B_0	B_0	B_1	B_1
	M	M_0	M_1	M_0	M_1
K_0	0.7	0.8	0.4	0.9	
K_1	0.3	0.2	0.6	0.1	

We can also print the probabilities for the other variables. Here there are no conditioned variables.

In [16]:

```
1 print(model.get_cpds('B'))
```

```
+-----+-----+
| B_0 | 0.4 |
+-----+-----+
| B_1 | 0.6 |
+-----+-----+
```

In [17]:

```
1 print(model.get_cpds('M'))
```

```
+-----+-----+
| M_0 | 0.8 |
+-----+-----+
| M_1 | 0.2 |
+-----+-----+
```

Inference in Belief Networks

From Barber, we can solve for a query pdf such as $p(B|K)$ using variable elimination as follows:

$$\begin{aligned} p(B|K) &= \sum_m p(B, m|K) \\ &= \sum_m \frac{p(B, m, K)}{p(K)} \\ &= \frac{p(B) \sum_m p(K|B, m)p(m)}{\sum_b p(b) \sum_m p(K|b, m)p(m)} \end{aligned}$$

Note how the summation variables are lowercase versions of the random variables to keep them distinct. The capital variables are set from $p(B|K)$, but the summation variables are b and m , which sum over the domains of B and M (which here are the same). Filling in the values above, we get

$$\begin{aligned} &= \frac{0.6 \times (0.2 \times 0.1 + 0.8 \times 0.6)}{0.6 \times (0.2 \times 0.1 + 0.8 \times 0.6) + 0.4 \times (0.2 \times 0.2 + 0.8 \times 0.3)} \\ &\approx 0.73 \end{aligned}$$

To solve for $p(B|K = \text{true})$ using variable elimination with `pgmpy` :

In [18]:

```
1 from pgmpy.inference import VariableElimination
2
3 inference = VariableElimination(model)
4 print(inference.query(['B'], evidence={'K' : 1}) ['B'])
```

```
+-----+-----+
| B      | phi(B) |
+=====+=====+
| B_0    | 0.2718  |
+-----+-----+
| B_1    | 0.7282  |
+-----+-----+
```

/anaconda3/lib/python3.6/site-packages/pgmpy/factors/discrete/DiscreteFactor.py:598: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
    phil.values = phil.values[slice_]
```

We can see that the result is return as a dictionary with a single key 'B':

In [19]:

```
1 inference.query(['B'], evidence={'K' : 1})
```

Out[19]:

```
{'B': <DiscreteFactor representing phi(B:2) at 0x1a26908e48>}
```

We can easily compute the probability for the other suspect $p(M|K)$

In [20]:

```
1 print(inference.query(['M'], evidence={'K' : 1}) ['M'])
```

```
+-----+-----+
| M      | phi(M) |
+=====+=====+
| M_0    | 0.9320  |
+-----+-----+
| M_1    | 0.0680  |
+-----+-----+
```

We can also add evidence. For example if we know that the maid could not have been present, then $p(B|K = \text{true}, M = \text{false})$

In [21]:

```
1 print(inference.query(['B'], evidence={'K' : 1, 'M' : 0})['B'])
```

```
+-----+-----+
| B      | phi(B) |
+=====+=====+
| B_0    | 0.2500 |
+-----+-----+
| B_1    | 0.7500 |
+-----+-----+
```

Which seems odd. How can the probability that the butler is the murder go *down* when we learn that the maid could not have been present? We can check the equations.

$$\begin{aligned} p(B|K, M) &= \frac{p(B, M, K)}{p(K, M)} \\ &= \frac{p(K|B, M)p(B)P(M)}{p(M) \sum_b p(K|b, M)p(b)} \end{aligned}$$

First we need to know the basics of defining and manipulating probability distributions in `pgmpy`.

In [22]:

```
1 from pgmpy.factors.discrete import DiscreteFactor
2 # DiscreteFactor(variables, cardinality, values)
3 phi1 = DiscreteFactor(['x1'], [2], [0.4, 0.6])
4 print(phi1)
```

```
+-----+-----+
| x1     | phi(x1) |
+=====+=====+
| x1_0   | 0.4000 |
+-----+-----+
| x1_1   | 0.6000 |
+-----+-----+
```

In [23]:

```
1 phi2 = DiscreteFactor(['x2'], [2], [0.3, 0.7])
2 print(phi2)
```

```
+-----+-----+
| x2     | phi(x2) |
+=====+=====+
| x2_0   | 0.3000 |
+-----+-----+
| x2_1   | 0.7000 |
+-----+-----+
```


In [24]:

```
1 phi3=phi1.product(phi2, inplace=False) # inplace means make copy, don't overwri
2 print(phi3)
```

x1	x2	phi(x1,x2)
x1_0	x2_0	0.1200
x1_0	x2_1	0.2800
x1_1	x2_0	0.1800
x1_1	x2_1	0.4200

/anaconda3/lib/python3.6/site-packages/pgmpy/factors/discrete/DiscreteFactor.py:586: FutureWarning: Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

```
phi.values = phi.values[slice_]
```

In [25]:

```
1 from pgmpy.factors.discrete import TabularCPD
2 px1 = TabularCPD('x1', 2, [[0.4, 0.6]])
3 #Or with keywords:
4 #px1 = TabularCPD(variable='x1', variable_card=2, values=[[0.4, 0.6]])
5 print(px1)
```

x1_0	0.4
x1_1	0.6

In [26]:

```
1 px2 = TabularCPD('x2', 2, [[0.3, 0.7]])
2 print(px2)
```

x2_0	0.3
x2_1	0.7

In []:

```
1 print(px1.product(px2, inplace=False))
```

So we can compute the product of the distributions in the numerator, to get the unnormalized values of $p(B|K, M)$

In [27]:

```
1 print(model.get_cpds('K'))
2 print(model.get_cpds('B'))
3 print(model.get_cpds('M'))
```

+-----+-----+-----+-----+-----+							
B	B_0	B_0	B_1	B_1			
+-----+-----+-----+-----+-----+							
M	M_0	M_1	M_0	M_1			
+-----+-----+-----+-----+-----+							
K_0	0.7	0.8	0.4	0.9			
+-----+-----+-----+-----+-----+							
K_1	0.3	0.2	0.6	0.1			
+-----+-----+-----+-----+-----+							
+-----+-----+							
B_0	0.4						
+-----+-----+							
B_1	0.6						
+-----+-----+							
+-----+-----+							
M_0	0.8						
+-----+-----+							
M_1	0.2						
+-----+-----+							

In [28]:

```
1 pB = model.get_cpds('B')
2 pM = model.get_cpds('M')
3 print(pB.product(pM, inplace=False))
```

+-----+-----+-----+-----+-----+					
M	M_0			M_1	
+-----+-----+-----+-----+-----+					
B_0	0.3200000000000000006			0.0800000000000000002	
+-----+-----+-----+-----+-----+					
B_1	0.48			0.12	
+-----+-----+-----+-----+-----+					

```
1 pKgBM = model.get_cpds('K')
2 print(pKgBM.product(pB, inplace=False))
```

In [30]:

```
1 print(pKgBM.product(pB, inplace=False).product(pM, inplace=False))
```

The probability that neither of the suspects is the murderer can be computed by computing $p(B, M|K)$. First we derive the theoretical solution using variable elimination.

So, all we have to do is compute the numerator, which we can do directly from the model, and normalize it. What we want to do get the product of the three distributions in the numerator, and then normalize the result.

In []:

```
1 model.get_cpds('K').values
```

In []:

```
1 print(model.get_cpds('K'))
```

In []:

```
1 pKgBM = model.get_cpds('K').values
2 for b in [1, 2]:
3     for m in [1, 2]:
4         for k in [1, 2]:
5             pBMgK[b][m][k] =
```

In [31]:

```
1 inference.query(['B', 'M'], evidence={'K' : 1})
```

Out[31]:

```
{'B': <DiscreteFactor representing phi(B:2) at 0x1a2773a2e8>,
'M': <DiscreteFactor representing phi(M:2) at 0x1a2773a748>}
```

In [32]:

```
1 # This does not do what we want, which is to
2 # compute the joint probability  $p(B, M \mid K)$ 
3 inference.query(['B', 'M'], evidence={'K' : 1})['B']
```

Out[32]:

```
<DiscreteFactor representing phi(B:2) at 0x1a2773a710>
```

In [33]:

```
1 # TODO: How do we print a probability table from a dictionary?
2 print(inference.query(['B', 'M'], evidence={'K' : 1})['B'])
3 print(inference.query(['B', 'M'], evidence={'K' : 1})['M'])
```

+-----+-----+	
B	phi(B)
+=====+	
B_0	0.2718
+-----+-----+	
B_1	0.7282
+-----+-----+	
+-----+-----+	
M	phi(M)
+=====+	
M_0	0.9320
+-----+-----+	
M_1	0.0680
+-----+-----+	

Explaining away

Here we illustrate the concept of *explaining away* with the wet grass example.

In []:

```
1
```