

贝叶斯网络python实战（以泰坦尼克号数据集为例，pgmpy库）

原创 leida_wt 发布于2019-03-24 23:05:36 阅读数 3091 ☆ 收藏



文章目录

- 贝叶斯网络简介
 - 贝叶斯推断思路
 - 贝叶斯网络
- 贝叶斯网络的实现
 - 应用步骤
 - 泰坦尼克数据集背景介绍
 - 模型结构搭建
 - 模型参数构建
 - 贝叶斯估计器
 - 推理
 - 自动设计网络结构->使用结构学习方法
 - 模型保存
- 先验
- 练手数据集
 - Binary Classification
 - Multiclass Classification
- 下载

本文的相关数据集，代码见文末百度云

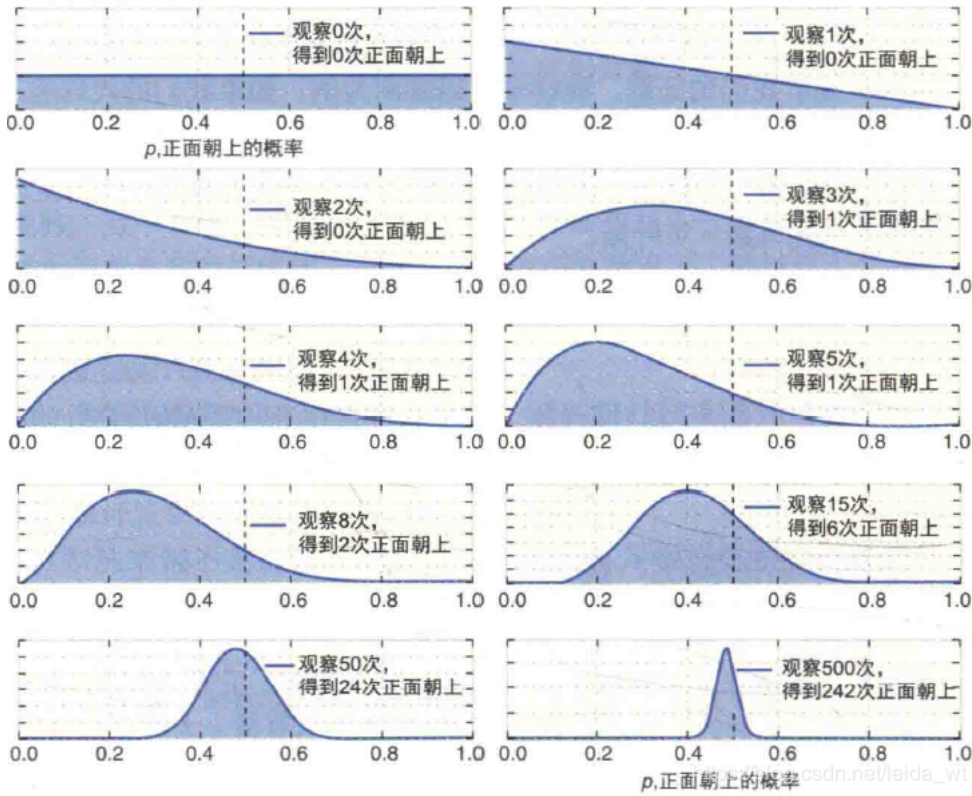
贝叶斯网络简介

贝叶斯网络是一种置信网络，一个生成模型。（判别模型，生成模型的区别可以这样：回答 $p(\text{label}|x)$ 即样本 x 属于某一类别的可能的，就是判别模型， $p(x, \text{label})$ 和 $p(x|\text{label})$ 的，即回答在给定的类别中找样本 x 及样本分布情况的，即为生成模型。生成模型给出的联合分布比判别网络能给出更多的信息，分布即可得 $p(\text{label}|x) p(x|\text{label})$)同时贝叶斯网络还是一个简单的白盒网络，提供了高可解释性的可能。相比于大热的几乎无所不能的深度神经网络，仍有他的优势和应用场景。比如在故障分析，疾病诊断里，我们不仅需要回答是不是，更重要的是回答为什么，并给出依据。这样的场景下，以贝叶斯一些可解释好的白盒网络更加有优势。

贝叶斯推断思路

与频率派直接从数据统计分析构建模型不同，贝叶斯派引入一个先验概率，表达对事件的已有了解，然后利用观测数据对先验知识进行修正，如通常把概率认为是0.5，这是个很朴素的先验知识，若是实验结果抛出了500-500的结果，那么证明先验知识是可靠的，合适的，若是出现100-900结果，那么

被逐渐修改 (越来越相信这是个作弊硬币), 当实验数据足够多的时候, 先验知识就几乎不再体现, 这时候得到与频率派几乎相同的结果。如图



具体例子推导可见[here](#)

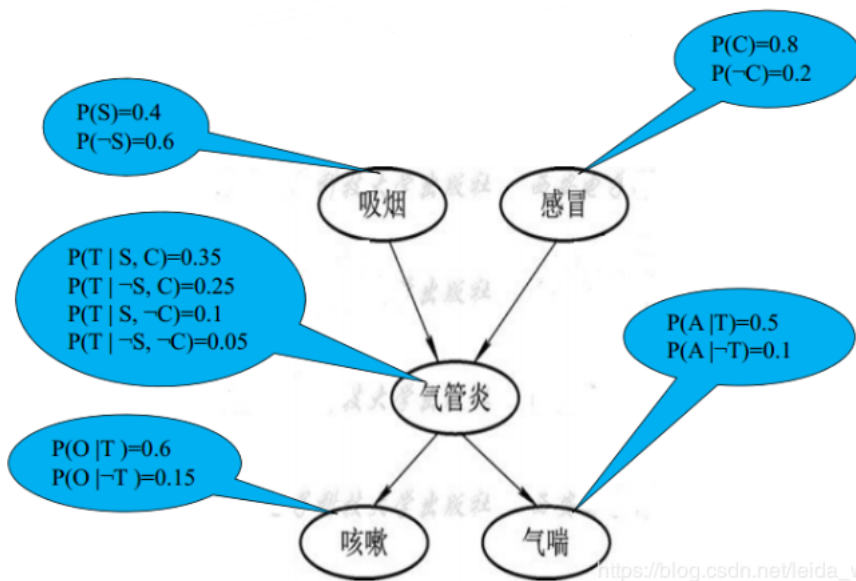
贝叶斯网络

贝叶斯网络结构如下所示, 其是有特征节点和链接构成的有向无环图。节点上是概率 $P(A), P(B)...$ 连接上是条件概率 $P(A|B) P(A|C) ...$ 即若有A指向B的边代表的就应为 $P(B|A)$, 更多信息可参考以下内容, 这里不再赘述, 贝叶斯网络结构本身不困难, 其难点主要在于推理算法等数值计算问题, 如为应用则

贝叶斯网络发展及其应用综述

《贝叶斯网络引论》@张连文

静态贝叶斯网络



贝叶斯网络的实现

相关工具一直很丰富, matlab, R上都有成熟的工具。这里使用了python下的pgmpy, 轻量好用, 不像pymc那样容易安装困难。安装:

```
conda install -c ankurankan pgmpy
```

或

```
pip install pgmpy
```

应用步骤

1.先确定以那些变量（特征）为节点，这里还包括由特征工程特征选择之类的工作。当然若有专业知识的参与会得到更合理的特征选择。

2.确定网络结构（拓扑）用以反应变量节点之间的依赖关系。也就是明确图的结构。这里既可以在有专家参与的情况下手工设计，也可以自动找到高效称为结构学习。贝叶斯网络的结构对最终网络性能很关键，若是构建所谓全连接贝叶斯网（即各个变量间两两相连），虽没有遗漏关联，但会导致严重为数据量很难支撑起全连接直接海量的条件概率。

3.明确每条边上的条件概率。和结构一样，参数也可由专家手工确定（先验），亦可通过数据自动学习（即参数学习），或两者同时进行。

下面以一个经典数据集为例展示如何利用pgmpy包进行贝叶斯网络建模

泰坦尼克数据集背景介绍

ref:<https://www.jianshu.com/p/9b6ee1fb7a60>

<https://www.kaggle.com/c/titanic>

这是kaggle经典数据集，主要是让参赛选手根据训练集中的乘客数据和存活情况进行建模，进而使用模型预测测试集中的乘客是否会存活。乘客特征由以下列出。这个数据集特征明确，数据量不大，很适合应用贝叶斯网络之类的模型来做，目前最好的结果是正确率应该有80+%（具体多少因为答案泄

PassengerId => 乘客ID

Pclass => 客舱等级(1/2/3等舱位)

Name => 乘客姓名

Sex => 性别

Age => 年龄

SibSp => 兄弟姐妹数/配偶数

Parch => 父母数/子女数

Ticket => 船票编号

Fare => 船票价格

Cabin => 客舱号

Embarked => 登船港口

在开始建模之前，先进行下特征工程，处理原始数据集的缺项等。这里前面处理主要采用<https://www.jianshu.com/p/9b6ee1fb7a60>的方法（他应月数据的技巧很值得一学），我在他的处理后，进一步进行了一些离散化处理，以使得数据符合贝叶斯网络的要求（贝叶斯网络也有支持连续变量的版本，学习的困难，目前还用的很少），最后保留5个特征。

```

1  '''
2  PassengerId => 乘客ID
3  Pclass => 客舱等级(1/2/3等舱位)
4  Name => 乘客姓名
5  Sex => 性别 清洗成male=1 female=0
6  Age => 年龄 插补后分0,1,2 代表 幼年(0-15) 成年(15-55) 老年(55-)
7  SibSp => 兄弟姐妹数/配偶数
8  Parch => 父母数/子女数
9  Ticket => 船票编号
10 Fare => 船票价格 经聚类变0 1 2 代表少 多 很多
11 Cabin => 客舱号 清洗成有此项，并发现有的生存率高
12 Embarked => 登船港口 清洗na,填S
13  '''
14 # combine train and test set.
15 train=pd.read_csv('./train.csv')
16 test=pd.read_csv('./test.csv')
17 full=pd.concat([train,test],ignore_index=True)
18 full['Embarked'].fillna('S',inplace=True)
19 full.Fare.fillna(full[full.Pclass==3]['Fare'].median(),inplace=True)
20 full.loc[full.Cabin.notnull(),'Cabin']=1
21 full.loc[full.Cabin.isnull(),'Cabin']=0
22 full.loc[full['Sex']=='male','Sex']=1
23 full.loc[full['Sex']=='female','Sex']=0
24
25 full['Title']=full['Name'].apply(lambda x: x.split(',')[1].split('.')[0].strip())
26 nn={'Capt':'Rareman','Col':'Rareman','Don':'Rareman','Dona':'Rarewoman',
27      'Dr':'Rareman','Jonkheer':'Rareman','Lady':'Rarewoman','Major':'Rareman',
28      'Master':'Master','Miss':'Miss','Mlle':'Rarewoman','Mme':'Rarewoman',
29      'Mr':'Mr','Mrs':'Mrs','Ms':'Rarewoman','Rev':'Mr','Sir':'Rareman',
30      'the Countess':'Rarewoman'}
31 full.Title=full.Title.map(nn)
32 # assign the female 'Dr' to 'Rarewoman'

```

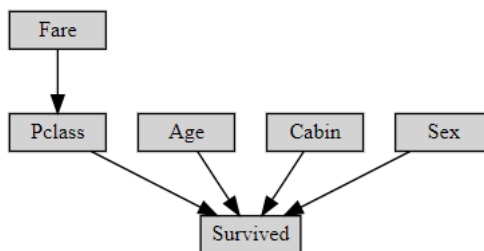
```

33 full.loc[full.PassengerId==797,'Title']='Rarewoman'
34 full.Age.fillna(999,inplace=True)
35 def girl(aa):
36     if (aa.Age!=999)&(aa.Title=='Miss')&(aa.Age<=14):
37         return 'Girl'
38     elif (aa.Age==999)&(aa.Title=='Miss')&(aa.Parch!=0):
39         return 'Girl'
40     else:
41         return aa.Title
42
43 full['Title']=full.apply(girl,axis=1)
44
45 Tit=['Mr','Miss','Mrs','Master','Girl','Rareman','Rarewoman']
46 for i in Tit:
47     full.loc[(full.Age==999)&(full.Title==i),'Age']=full.loc[full.Title==i,'Age'].median()
48
49 full.loc[full['Age']<=15,'Age']=0
50 full.loc[(full['Age']>15)&(full['Age']<55),'Age']=1
51 full.loc[full['Age']>=55,'Age']=2
52 full['Pclass']=full['Pclass']-1
53 from sklearn.cluster import KMeans
54 Fare=full['Fare'].values
55 Fare=Fare.reshape(-1,1)
56 km = KMeans(n_clusters=3).fit(Fare) #将数据集分为2类
57 Fare = km.fit_predict(Fare)
58 full['Fare']=Fare
59 full['Fare']=full['Fare'].astype(int)
60 full['Age']=full['Age'].astype(int)
61 full['Cabin']=full['Cabin'].astype(int)
62 full['Pclass']=full['Pclass'].astype(int)
63 full['Sex']=full['Sex'].astype(int)
64 #full['Survived']=full['Survived'].astype(int)
65
66
67 dataset=full.drop(columns=['Embarked','Name','Parch','PassengerId','SibSp','Ticket','Title'])
68 dataset.dropna(inplace=True)
69 dataset['Survived']=dataset['Survived'].astype(int)
70 #dataset=pd.concat([dataset, pd.DataFrame(columns=['Pri'])])
71 train=dataset[:800]
72 test=dataset[800:]
73 '''
74 最后保留如下项目,并切出800的训练集:
75 Pclass => 客舱等级(0/1/2等舱位)
76 Sex => 性别 male=1 female=0
77 Age => 年龄 插补后分0,1,2 代表 幼年(0-15) 成年(15-55) 老年(55-)
78 Fare => 船票价格 经聚类变0 1 2 代表少 多 很多
79 Cabin => 客舱号 清洗成有或无此项,并发现有的生存率高
80 '''

```

模型结构搭建

这里先手动设计网络结构。



凭借对数据的理解,先设计如下的结构

```

1 from pgmpy.models import BayesianModel
2 from pgmpy.estimators import BayesianEstimator
3
4 #model = BayesianModel([('Age', 'Pri'), ('Sex', 'Pri'), ('Pri', 'Survived'), ('Fare', 'Pclass'), ('Pclass', 'Survived'), ('Cabin', 'Survived')])
5 model = BayesianModel([('Age', 'Survived'), ('Sex', 'Survived'), ('Fare', 'Pclass'), ('Pclass', 'Survived'), ('Cabin', 'Survived')])

```

其中('Age' , 'Survived')表示Age指向Survived

pgmpy没有提供可视化, 这里简单用graphviz实现了一下。

```

1 def showBN(model, save=False):
2     '''传入BayesianModel对象, 调用graphviz绘制结构图, jupyter中可直接显示'''
3     from graphviz import Digraph
4     node_attr = dict(
5         style='filled',
6         shape='box',
7         align='left',
8         fontsize='12',
9         ranksep='0.1',
10        height='0.2'
11    )
12    dot = Digraph(node_attr=node_attr, graph_attr=dict(size="12,12"))
13    seen = set()
14    edges=model.edges()
15    for a,b in edges:
16        dot.edge(a,b)
17    if save:
18        dot.view(cleanup=True)
19    return dot
20 showBN(model)

```

模型参数构建

接下来就是确定网络的参数, 也就是各个边上的条件概率。

若手工填入, 可这样写

Examples

For a distribution of $P(\text{grade}|\text{diff}, \text{intel})$

diff	easy			hard		
	dumb	avg	smart	dumb	avg	smart
intel	0.1	0.1	0.1	0.1	0.1	0.1
gradeA	0.1	0.1	0.1	0.1	0.1	0.1
gradeB	0.1	0.1	0.1	0.1	0.1	0.1
gradeC	0.8	0.8	0.8	0.8	0.8	0.8

values should be $[[0.1,0.1,0.1,0.1,0.1,0.1], [0.1,0.1,0.1,0.1,0.1,0.1], [0.8,0.8,0.8,0.8,0.8,0.8]]$

```

>>> cpd = TabularCPD('grade',3,[[0.1,0.1,0.1,0.1,0.1,0.1],
                                [0.1,0.1,0.1,0.1,0.1,0.1],
                                [0.8,0.8,0.8,0.8,0.8,0.8]],
                                evidence=['diff', 'intel'], evidence_card=[2,3])

>>> print(cpd)
+-----+-----+-----+-----+-----+-----+
| diff  | diff_0 | diff_0 | diff_0 | diff_1 | diff_1 | diff_1 |
+-----+-----+-----+-----+-----+-----+
| intel | intel_0| intel_1| intel_2| intel_0| intel_1| intel_2|
+-----+-----+-----+-----+-----+-----+
| grade_0| 0.1   | 0.1   | 0.1   | 0.1   | 0.1   | 0.1   |
+-----+-----+-----+-----+-----+-----+
| grade_1| 0.1   | 0.1   | 0.1   | 0.1   | 0.1   | 0.1   |
+-----+-----+-----+-----+-----+-----+
| grade_2| 0.8   | 0.8   | 0.8   | 0.8   | 0.8   | 0.8   |
+-----+-----+-----+-----+-----+-----+
>>> cpd.values
array([[ [ 0.1, 0.1, 0.1],
        [ 0.1, 0.1, 0.1]],

```

https://blog.csdn.net/leida_wt

```

1 from pgmpy.factors.discrete import TabularCPD
2 #构建cpd表
3 my_cpd= TabularCPD(variable='Pclass', variable_card=3,
4                    values=[[0.65, 0.3], [0.30, 0.6],[0.05,0.1]],
5                    evidence=['Fare'], evidence_card=[2])
6 # 填cpd表
7 model.add_cpds(my_cpd)
8
9 # 执行检查 (可选, 用于检查cpd是否填错)
10 cancer_model.check_model()

```

但在此例子里, 使用参数学习的办法从数据里自动学习。

参数学习有两种典型方法, 极大似然估计和贝叶斯估计。因为前者的过拟合严重, 一般都使用后者进行参数学习。pgmpy提供的贝叶斯估计器提供三

支持，'dirichlet'，'BDeu'，'K2'，实际上都是dirichlet分布，这里解释下贝叶斯估计器的工作原理。

贝叶斯估计器

在贝叶斯分析的框架下，待求参数 θ 被看做是随机变量，对他的估计就是在其先验上，用数据求后验，因此先要有对 θ 的先验假设。而我们通常取的先验

$$\theta_i = P(X = x_i), i = 1, 2, \dots, r$$

dirichlet（狄利克雷）分布。对于一个含有 i 个离散状态的节点，我们设其参数为 α_i ，并令其先验为狄利克雷分布 $D[\alpha_1, \alpha_2, \dots, \alpha_i]$ （ $i=2$ 时也称beta分布）

$$p(\theta) = \frac{\Gamma(\alpha)}{\prod_{i=1}^r \Gamma(\alpha_i)} \prod_{i=1}^r \theta_i^{\alpha_i - 1}$$

这个先验有 i 个参数，数学上证明了，这些参数就相当于将先验表达成了 α 个虚拟样本，其中满足 $X=x_i$ 的样本数为 α_i ，这个 α 就成为等价样本量（equivalent_sample_size）。这个巧合其实正式先验函数取这个函数的缘由，另外，其计算后的后验分布也是狄利克雷分布（称这种行为叫共轭先验），这个分布可参考 <https://zhuanlan.zhihu.com/p/37976562>

至此可以理解pgmpy提供的贝叶斯估计器的参数的含义了。

其定义为estimate_cpd(node, prior_type='BDeu', pseudo_counts=[], equivalent_sample_size=5)

node是节点名

当prior_type='BDeu' 就表示选择了一个equivalent_sample_size=5的无差别客观先验，认定各个概率相等，不提供信息，但并不是没用，这个神经网络里头控制过拟合的正则项的作用。

当prior_type='dirichlet' 表示选择一般的狄利克雷分布，这时候要主动填入 $[\alpha_1, \alpha_2, \dots, \alpha_i]$

当prior_type='K2' 意为 'dirichlet' + setting every pseudo_count to 1

具体使用如下：

```
1 | data = pd.DataFrame(data={'A': [0, 0, 1], 'B': [0, 1, 0], 'C': [1, 1, 0]})
2 | model = BayesianModel([('A', 'C'), ('B', 'C')])
3 | estimator = BayesianEstimator(model, data)
4 | cpd_C = estimator.estimate_cpd('C', prior_type="dirichlet", pseudo_counts=[1, 2])
5 | model.add_cpds(cpd_C)
6 |
```

上面是一个一个填进去的，在本例中有更简单的方法，就是利用提供的fit函数，一并估计各个cpd（条件概率），即

```
1 | model.fit(train, estimator=BayesianEstimator, prior_type="BDeu") # default equivalent_sample_size=5
```

直接把前面处理得到dataframe 传入即可。这里记录一个bug：pgmpy目前将离散变量命名限制为从0开始，所以本例子里的Pclass 项从（1/2/3等级）成了（0/1/2等级）以解决此问题。

dirichlet也可在fit函数里使用，只要传入pseudo_counts字典即可，如下面这样

```
1 | pseudo_counts = {'D': [300, 700], 'I': [500, 500], 'G': [800, 200], 'L': [500, 500], 'S': [400, 600]}
2 | model.fit(data, estimator=BayesianEstimator, prior_type='dirichlet', pseudo_counts=pseudo_counts)
```

到此为止，模型已经完全构建完毕，下面可以开始使用其进行推理了。

推理

首先可以通过一些方法查看模型

```
1 | #输出节点信息
2 | print(model.nodes())
3 | #输出依赖关系
4 | print(model.edges())
5 | #查看某节点概率分布
6 | print(model.get_cpds('Pclass').values)
```

当然我们更关心的是给定某些节点后，感兴趣节点的概率等，这就是推理。

贝叶斯网络推理分成：

1.后验概率问题：

表达为求 $P(Q|E=e)$ 其中Q为查询变量 E为证据变量

即如本例子里已知一个人，女，<15岁，高票价，问生还几率是多少？

2.最大后验假设问题 (MAP) :

$$h^* = \arg \max_{\mathbf{h}} P(\mathbf{H} = \mathbf{h} | \mathbf{E} = \mathbf{e}).$$

已知证据E时，对某些变量的转态组合感兴趣（称假设变量H），找最可能组合就是最大后验假设问题。如本例子里，问一个活下来的女乘客最可能有几岁年龄？

3.最大可能解释问题 (MPE) : 是2的特例，即假设包含网络里所有非证据变量（同时也可包含证据变量）

贝叶斯网络推理主要有两类方法，精确推理（变量化简Variable Elimination和置信传播）和近似推理(如mcmc采样)，一般精确推理足以解决pgmpy解决1可以用query函数 解决2, 3可以用map_query函数

通过这些查询可以获得我们感兴趣的关于因果关系信息，这是贝叶斯网络模型的一大优势。此处的因果关系并不可解释为一般意义上的逻辑因果，而率上的相关，比如我们不能将 $P(\text{天亮了}|\text{公鸡打鸣})$ 很高解释为是因为公鸡打鸣天才亮的。

```

1 | from pgmpy.inference import VariableElimination
2 | model_infer = VariableElimination(model)
3 | q = model_infer.query(variables=['Survived'], evidence={'Fare': 0})
4 | print(q['Survived'])
5 | '''
6 | +-----+-----+
7 | | Survived | phi(Survived) |
8 | +-----+-----+
9 | | Survived_0 |          0.6341 |
10 | +-----+-----+
11 | | Survived_1 |          0.3659 |
12 | +-----+-----+
13 | '''
14 | q = model_infer.map_query(variables=['Fare','Age','Sex','Pclass','Cabin'], evidence={'Survived': 1})
15 | print(q#{'Sex': 0, 'Fare': 0, 'Age': 1, 'Pclass': 2, 'Cabin': 0})
16 |

```

上面的代码使用了VariableElimination方法，亦可用BeliefPropagation，其有相同的接口。

与fit函数类似，也提供了输入dataframe的简便推理方法predict，如下

只要剔除想预测的列输入predict函数，就可以得到预测结果的dataframe

```

1 | predict_data=test.drop(columns=['Survived'],axis=1)
2 | y_pred = model.predict(predict_data)
3 | print((y_pred['Survived']==test['Survived']).sum()/len(test))
4 | #测试集精度0.8131868131868132

```

只是用了泰坦尼克数据集的一部分特征随手设计网络就可以达到不错的效果了，精度81.3%，上传kaggle的正确率0.77990。

自动设计网络结构->使用结构学习方法

ref:

https://github.com/pgmpy/pgmpy_notebook

《贝叶斯网络引论》

自动设计网络结构的核心问题有两个，一个是评价网络好坏的指标，另一个是查找的方法。穷举是不可取的，因为组合数太大，只能是利用各种启发式搜索条件以减少搜索空间，因此产生两大类方法，Score-based Structure Learning与constraint-based structure learning 以及他们的结合hybrid learning。

1.Score-based Structure Learning

这个方法依赖于评分函数，常用的有bdeu k2 bic，更合理的网络评分更高，如下面的例子

此例子随机产生x y 并令 $z=x+y$ ，显然 $X \rightarrow Z \leftarrow Y$ 的结构合理

```

1 | import pandas as pd
2 | import numpy as np
3 | from pgmpy.estimators import BdeuScore, K2Score, BicScore
4 | from pgmpy.models import BayesianModel
5 |
6 | # create random data sample with 3 variables, where Z is dependent on X, Y:
7 | data = pd.DataFrame(np.random.randint(0, 4, size=(5000, 2)), columns=list('XY'))
8 | data['Z'] = data['X'] + data['Y']
9 |

```

```

10 | bdeu = BdeuScore(data, equivalent_sample_size=5)
11 | k2 = K2Score(data)
12 | bic = BicScore(data)
13 |
14 | model1 = BayesianModel([('X', 'Z'), ('Y', 'Z')]) # X -> Z <- Y
15 | model2 = BayesianModel([('X', 'Z'), ('X', 'Y')]) # Y <- X -> Z
16 |
17 |
18 | print(bdeu.score(model1))
19 | print(k2.score(model1))
20 | print(bic.score(model1))
21 |
22 | print(bdeu.score(model2))
23 | print(k2.score(model2))
24 | print(bic.score(model2))
25 | '''
26 | -13936.101051153362
27 | -14326.88012027081
28 | -14292.1400887
29 | -20902.744280734016
30 | -20929.567083476162
31 | -20946.7926535
32 | '''

```

X -> Z <- Y 的评分更高

而依据评分函数进行搜索的搜索方法常用的有穷举 (5个节点以下可用) 和爬山算法 (一个贪婪算法) pympy的实现如下:

```

1 | from pgmpy.estimators import HillClimbSearch
2 |
3 | # create some data with dependencies
4 | data = pd.DataFrame(np.random.randint(0, 3, size=(2500, 8)), columns=list('ABCDEFGH'))
5 | data['A'] += data['B'] + data['C']
6 | data['H'] = data['G'] - data['A']
7 |
8 | hc = HillClimbSearch(data, scoring_method=BicScore(data))
9 | best_model = hc.estimate()
10 | print(best_model.edges())
11 | #[('A', 'C'), ('A', 'B'), ('C', 'B'), ('G', 'A'), ('G', 'H'), ('H', 'A')]
12 |

```

2.Constraint-based Structure Learning

比如根据独立性得到最优结构的方法, 相对来讲前一种更有效

```

1 | from pgmpy.independencies import Independencies
2 |
3 | ind = Independencies(['B', 'C'],
4 |                      ['A', ['B', 'C'], 'D'])
5 | ind = ind.closure() # required (!) for faithfulness
6 |
7 | model = ConstraintBasedEstimator.estimate_from_independencies("ABCD", ind)
8 |
9 | print(model.edges())
10 | #[('A', 'D'), ('B', 'D'), ('C', 'D')]

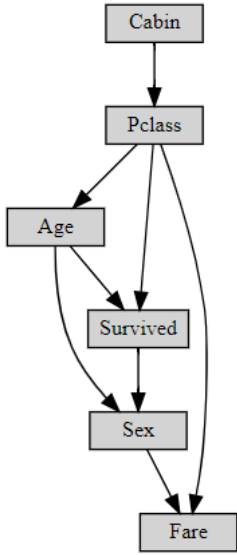
```

回到泰坦尼克的例子, 使用HillClimbSearch试了以下, 在自己的测试集得到了和之前手工设计网络相同的精度 (kaggle测试成绩略低一些), 但是模型, 不过看看给出的模型可以发现一些有趣的东西, 比如Cabin Pclass Fare 相关, Age还影响了Sex等等

```

1 | from pgmpy.estimators import HillClimbSearch
2 | from pgmpy.estimators import BdeuScore, K2Score, BicScore
3 | hc = HillClimbSearch(train, scoring_method=BicScore(train))
4 | best_model = hc.estimate()
5 | #print(best_model.edges())
6 | showBN(best_model)

```

```

1 | best_model.fit(train, estimator=BayesianEstimator, prior_type="BDeu") # default equivalent_sample_size=5
2 | predict_data=test.drop(columns=['Survived'],axis=1)
3 | y_pred = best_model.predict(predict_data)
4 | (y_pred['Survived']==test['Survived']).sum()/len(test)#测试集精度
5 | #0.8131868131868132
  
```

模型保存

略

先验

这里摘录一些下文对先验的介绍

ref:Bayesian Methods for Hackers chp6

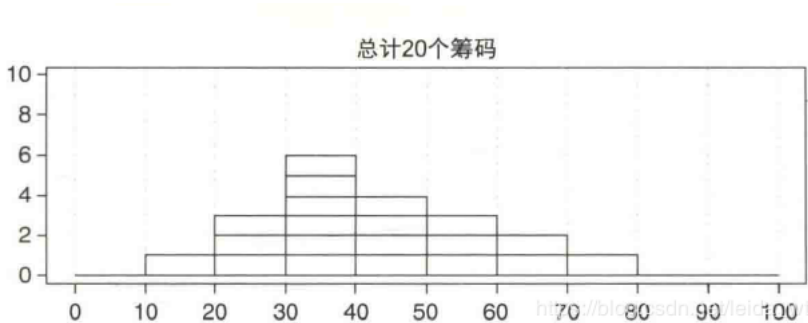
贝叶斯先验可以分为两类:客观先验和主观先验。客观先验的目的是让数据对后验影响最大, 主观先验的目的是让从业者对前验表达自己的观点。事实上, 从选择先验分布开始就已经开始搭建模型了, 属于建模的一部分。如果后验不符合要求, 自然可修改更换先验, 这是无关紧要的。没有正确的用的模型。

经验贝叶斯方法

一种叫经验贝叶斯方法融合了频率派的做法, 采用 数据>先验>数据>后验 的方法, 从观测数据的统计特征构建先验分布。这种方法实则违背了贝叶斯据>后验 的思路。

从专家获得先验分布

可以考虑实验转盘赌法捕获专家认为是的先验分布形状。做法是让专家吧固定总数的筹码放在各个区间上以此来表达各个区段的概率, 如下所示。之后可对此进行建模拟合, 得到专家先验。



判断先验是否合适

只要先验在某处概率不为零, 后验就有机会在此处表达任意的概率。当后验分布的概率堆积在先验分布的上下界时, 那么很肯先验是不大对的。(比如 Uniform (0,0.5) 作为真实值p=0.7的先验, 那么后验推断的结果会堆积在0.5一侧, 表明真实值极可能大于0.5)

练手数据集

Binary Classification

- Indian Liver Patient Records
- Synthetic Financial Data for Fraud Detection
- Business and Industry Reports
- Can You Predict Product Backorders?
- Exoplanet Hunting in Deep Space
- Adult Census Income

Multiclass Classification

- Iris Species
- Fall Detection Data from China
- Biomechanical Features of Orthopedic Patients

下载

链接: <https://pan.baidu.com/s/1lR3-tiSt-NJgogg4al8u-w> 提取码: 2yc1

👍 点赞 8 ☆ 收藏 📄 分享 ...



leida_wt

发布了40 篇原创文章 · 获赞 51 · 访问量 7万+

私信

机器学习sklearn之贝叶斯网络实战 (一)

阅读数 1176

贝叶斯网络贝叶斯网络、信念网络、贝叶斯模型或概率定向无环图模型是一种概率图模型(一种统计模型),通过有... 博文 来自: [weixin_41599977...](#)



想对作者说点什么

查看 23 条热评

数学推导+纯Python实现机器学习算法12: 贝叶斯网络

阅读数 48

Python机器学习算法实现Author: louwill 在上一讲中, 我们讲到了经典的朴素贝叶斯算法。朴素贝叶斯的一大特点... 博文 来自: [weixin_37737254...](#)

贝叶斯网络应用实例一: 胸部疾病诊所

阅读数 2294

转自: <http://blog.sciencenet.cn/blog-82650-255141.html> 以下内容摘录自www.norsys.com, 根据实例内容意... 博文 来自: [sueji的专栏](#)

使用 Python 第三方库 daft 绘制 PGM 中的贝叶斯网络

阅读数 2333

daft 的官方文档请见 DAFT: BEAUTIFULLY RENDERED PROBABILISTIC GRAPHICAL MODELS.from matplotlib i... 博文 来自: <https://space.bilibi...>

pgmpy包构建贝叶斯网络架构 (Bayesian Network Structure)

阅读数 2536

import pandas as pdimport numpy as npfrom pgmpy.estimators import ExhaustiveSearch, HillClimbSearch,... 博文 来自: [csdn_47的博客](#)

用Python构建概率图模型

06-18

利用Python的pgmpy包构建概率图模型

下载

用于贝叶斯网络的python库--pgmpy安装配置

阅读数 8734

1.软件准备pgmpy:<https://github.com/pgmpy/pgmpy/archive/dev.zip>minconda(最新版):<https://mirrors.tuna...> 博文 来自: [yangbei1993的博客](#)

全栈必备 贝叶斯方法

阅读数 253

数据的重要性毋庸置疑, 但是如何让数据产生价值呢? 对一个全栈老码农而言, 经常在开发或者研发管理的时候遇到... 博文 来自: [喔家ArchiSelf](#)

python机器学习——贝叶斯方法

阅读数 649

贝叶斯公式百度百科贝叶斯定理由英国数学家贝叶斯(ThomasBayes1702-1761)发展, 用来描述两个条件概率之间的... 博文 来自: [乐亦亦乐的博客](#)

用于贝叶斯网络的python库--pgmpy安装配置 - yangbei19..._CSDN博客