

# Лабораторная работа: Bayesian Networks на примере датасета Mushroom Classification

Люзина Мария

M8O-307Б-23

26.11.25

Машинное обучение

# Введение

**Байесовская сеть** (Bayesian Networks) – это вероятностная графическая модель, которая представляет: переменные как узлы, прямые зависимости как стрелки; хранит в каждом узле условные вероятности  $P(\text{узел}|\text{его родители})$ .

Байесовские сети нужны, когда: важна интерпретируемость, данных мало, есть причинно-следственные связи.

*Например*, Байесовские сети применяются в диагностике болезней в медицине:

Узлы:

- Курение  $\rightarrow$  Рак лёгких
- Курение  $\rightarrow$  Бронхит
- Рак лёгких  $\rightarrow$  Положительный рентген
- Рак лёгких  $\rightarrow$  Одышка
- Туберкулёз  $\rightarrow$  Положительный рентген
- Посещение Азии  $\rightarrow$  Туберкулёз

Пациент приходит с одышкой и положительным рентгеном. Врач вводит эти два факта  $\rightarrow$  сеть мгновенно пересчитывает:

$$P(\text{рак лёгких} \mid \text{одышка} + \text{положительный рентген}) = 0.92$$

$$P(\text{туберкулёз} \mid \dots) = 0.05$$

И выдаёт не только диагноз, но и объяснение, какие симптомы на что влияют.

# Датасет: Mushroom Classification

Цель – определить, съедобный ли гриб или ядовитый. Всего в датасете 8124 объекта, 22 категориальных признаков и 1 целевая переменная. В датасете нет пропусков, что упрощает подготовку данных. Этот датасет сбалансирован:  $\approx 51.8\%$  съедобных,  $48.2\%$  ядовитых грибов.

Признаки описывают морфологию гриба: форма и цвет шляпки, ножки, жабр, кольца, цвет спор, тип вуали, место обитания и т.д.

Несмотря на большое количество признаков, существует простое биологически осмысленное правило классификации, которое человек сразу не видит, а алгоритмы поиска структуры байесовских сетей находят автоматически.

Mushrooms – классический бенчмарк для демонстрации силы и интерпретируемости байесовских сетей.

# Загрузка и обработка датасета

Загрузка датасета:

```
path = 'mushrooms.csv'  
df = pd.read_csv(path)  
df.head(3)
```

	class	cap-shape	cap-surface	cap-color	bruises	odor	gill-attachment	gill-spacing	gill-size	gill-color	stalk-surface-below-ring	stalk-color-above-ring	stalk-color-below-ring	veil-type	veil-color	ring-number	ring-type	spore-print-color	population	habitat
0	p	x	s	n	t	p	f	c	n	k	s	w	w	p	w	o	p	k	s	u
1	e	x	s	y	t	a	f	c	b	k	s	w	w	p	w	o	p	n	n	g
2	e	b	s	w	t	l	f	c	b	n	s	w	w	p	w	o	p	n	n	m

3 rows x 23 columns

Размер датасета и информация про него:

```
print('Размер датасета: ', df.shape)  
df.info()
```

Размер датасета: (8124, 23)

```
RangeIndex: 8124 entries, 0 to 8123  
Data columns (total 23 columns):  
#   Column                Non-Null Count  Dtype    
---  ---                     -  
0   class                  8124 non-null  object   
1   cap-shape              8124 non-null  object   
2   cap-surface            8124 non-null  object   
3   cap-color              8124 non-null  object   
4   bruises                8124 non-null  object   
5   odor                  8124 non-null  object   
6   gill-attachment        8124 non-null  object   
7   gill-spacing           8124 non-null  object   
8   gill-size              8124 non-null  object   
9   gill-color             8124 non-null  object   
10  stalk-shape            8124 non-null  object   
11  stalk-root             8124 non-null  object   
12  stalk-surface-above-ring 8124 non-null  object   
13  stalk-surface-below-ring 8124 non-null  object   
14  stalk-color-above-ring  8124 non-null  object   
15  stalk-color-below-ring  8124 non-null  object   
16  veil-type              8124 non-null  object   
17  veil-color             8124 non-null  object   
18  ring-number            8124 non-null  object   
19  ring-type              8124 non-null  object   
...  
21  population             8124 non-null  object   
22  habitat                8124 non-null  object   
dtypes: object(23)  
memory usage: 1.4+ MB
```

# Загрузка и обработка датасета

Признак value-type полностью константный:

```
print("Уникальные значения в veil-type:", df['veil-type'].unique())
```

```
Уникальные значения в veil-type: ['p']
```

Поэтому удалим его, ведь он не несет никакой информации:

```
df = df.drop('veil-type', axis=1)
```

Для рgmру не будем применять LabelEncoder, потому что он работает со строковыми или категориальными метками. Все 22 оставшихся признака уже дискретные категориальные, поэтому дискретизация не требуется

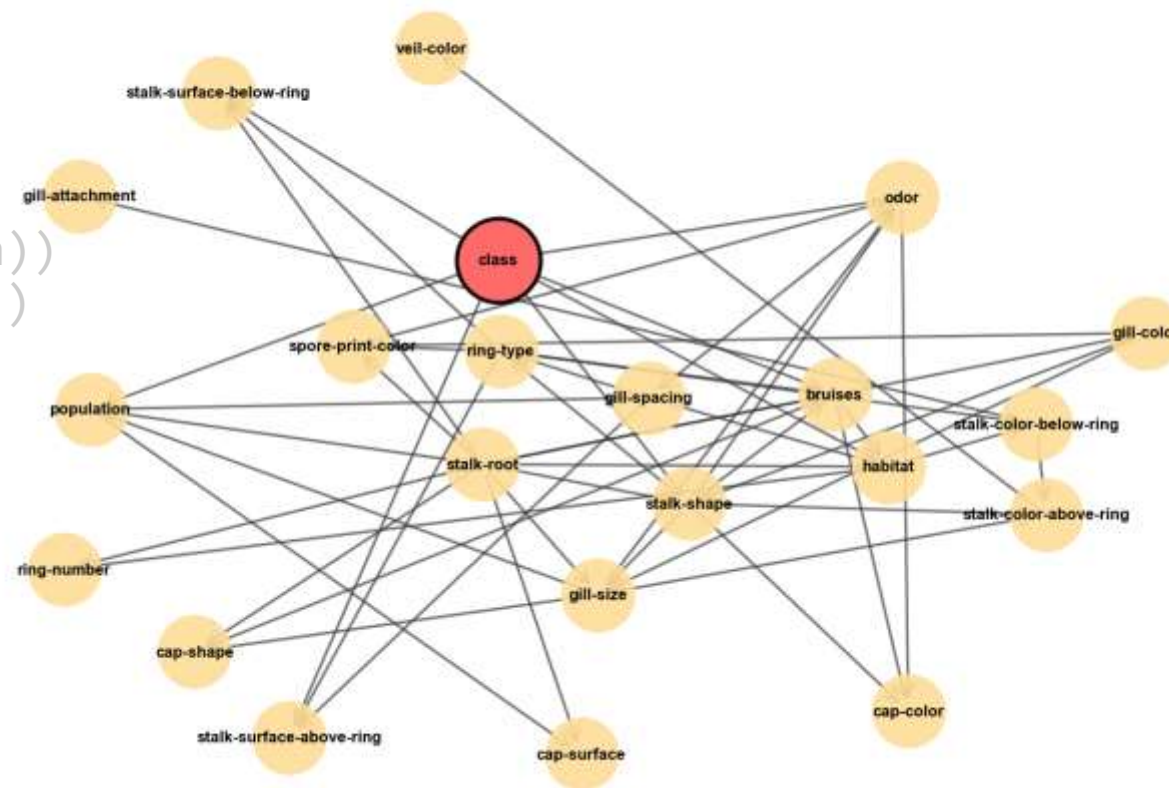
# Построение структуры сети

Автоматическое построение:

```
hc = HillClimbSearch(data)
best_dag
hc.estimate(scoring_method=BicScore(data))
model = BayesianNetwork(best_dag.edges())
```

Найдено узлов: 22  
Найдено рёбер: 52  
Структура построена.

Автоматически найденная структура байесовской сети  
(HillClimbSearch + BIC, 52 рёбер)



# Оценка параметров и СРТ: код model fit

```
model.fit(  
    df,      - наш датафрейм  
    estimator=BayesianEstimator, - используем  
байесовский оценщик, а не частотный  
    prior_type='BDeu', - тип априорного распределения  
    equivalent_sample_size=10 - «виртуальный» размер  
выборки для сглаживания  
)
```

# Оценка параметров и СРТ: просмотр вероятностных таблиц

```
for node in model.nodes():  
    cpt = model.get_cpds(node)  
    print(f"\nCPT для узла: {node}")  
    print(cpt)
```

СРТ для узла: class

+-----+-----+-----+		
odor	...	odor(y)
+-----+-----+-----+		
stalk-shape	...	stalk-shape(t)
+-----+-----+-----+		
class(e)	...	0.0004817883985353632
+-----+-----+-----+		
class(p)	...	0.9995182116014646
+-----+-----+-----+		

СРТ для узла: habitat

+-----+-----+-----+		
bruises	...	bruises(t)
+-----+-----+-----+		
class	...	class(p)
+-----+-----+-----+		
gill-spacing	...	gill-spacing(w)
+-----+-----+-----+		
stalk-root	...	stalk-root(r)
+-----+-----+-----+		
habitat(d)	...	0.14285714285714288
+-----+-----+-----+		
habitat(g)	...	0.14285714285714288
+-----+-----+-----+		
habitat(l)	...	0.14285714285714288
+-----+-----+-----+		
habitat(m)	...	0.14285714285714288
+-----+-----+-----+		
habitat(p)	...	0.14285714285714288
+-----+-----+-----+		
habitat(u)	...	0.14285714285714288
+-----+-----+-----+		
habitat(w)	...	0.14285714285714288
+-----+-----+-----+		



# Визуализация сети

```
import matplotlib.pyplot as plt
import networkx as nx

G = nx.DiGraph()
G.add_edges_from(model.edges())

plt.figure(figsize=(14, 10))

pos = nx.spring_layout(G, k=1.5, iterations=50, seed=42)

nx.draw_networkx_nodes(G, pos, node_size=3000, node_color="#ffdd99", alpha=0.9)
nx.draw_networkx_nodes(G, pos, nodelist=['class'], node_size=4000,
                        node_color="#ff6b6b", edgecolors='black', linewidths=3)
nx.draw_networkx_edges(G, pos, width=2, alpha=0.7, edge_color="#333333",
                        arrows=True, arrowsize=25, arrowstyle='->')

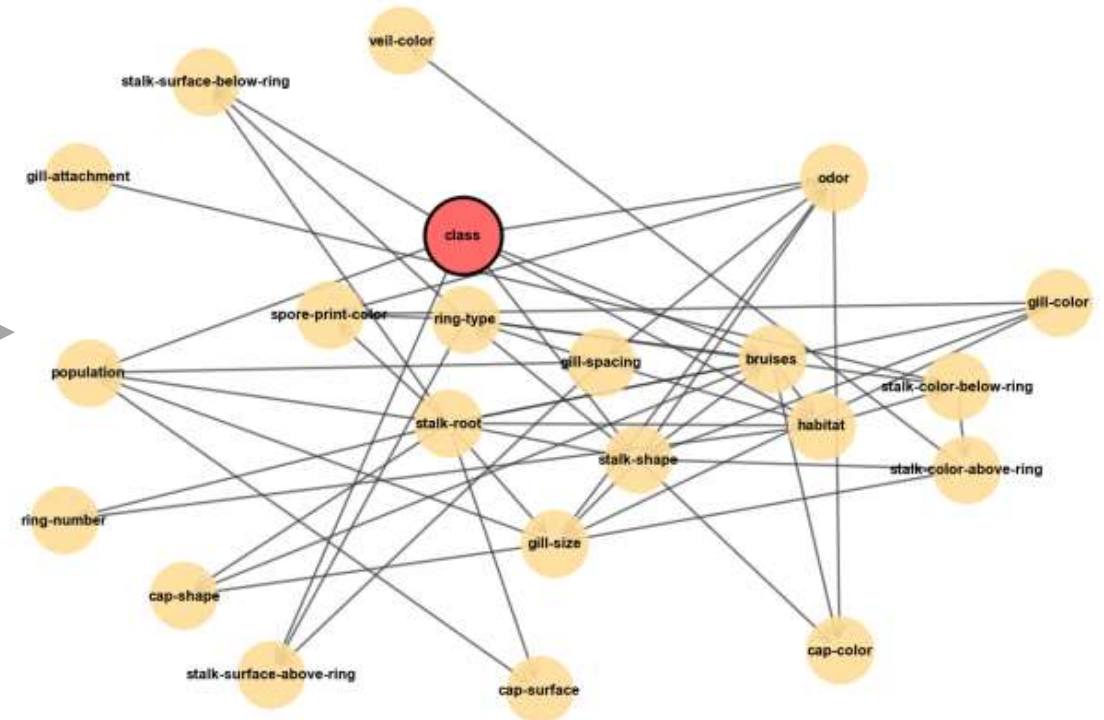
nx.draw_networkx_labels(G, pos, font_size=12, font_weight="bold", font_family="Arial")

plt.title("Автоматически найденная структура байесовской сети\n"
          f"(HillClimbSearch + BIC, {len(G.edges())} рёбер)",
          fontsize=18, pad=20, fontweight='bold')

plt.axis('off')
plt.tight_layout()
plt.show()

print("Топ-10 найденных зависимостей:")
for i, (u, v) in enumerate(list(G.edges())[:10], 1):
    print(f"{i:2}. {u} → {v}")
```

Автоматически найденная структура байесовской сети  
(HillClimbSearch + BIC, 52 рёбер)



# Inference и результаты

```
from pgmpy.inference import VariableElimination
infer = VariableElimination(model)

print("=== ПРОВЕРКА ИНФЕРЕНСА ===")
test_odors = ['n', 'a', 'l', 'f', 'c', 'y', 'm', 'p', 's']

for odor in test_odors:
    try:
        prob = infer.query(['class'], evidence={'odor': odor})
        p_poisonous = prob.values[1]
        status = "ЯДОВИТЫЙ" if p_poisonous > 0.8 else "СЪЕДОБНЫЙ"
        print(f"odor = {odor:>2} → P(ядовитый) = {p_poisonous:.10f} → {status}")
    except Exception as e:
        print(f"odor = {odor:>2} → ОШИБКА: {e}")
```

```
=== ПРОВЕРКА ИНФЕРЕНСА ===
odor =  n → P(ядовитый) = 0.0359568491 → СЪЕДОБНЫЙ
odor =  a → P(ядовитый) = 0.0011584760 → СЪЕДОБНЫЙ
odor =  l → P(ядовитый) = 0.0011584760 → СЪЕДОБНЫЙ
odor =  f → P(ядовитый) = 0.9997332099 → ЯДОВИТЫЙ
odor =  c → P(ядовитый) = 0.9468853966 → ЯДОВИТЫЙ
odor =  y → P(ядовитый) = 0.9975347155 → ЯДОВИТЫЙ
odor =  m → P(ядовитый) = 0.8060939732 → ЯДОВИТЫЙ
odor =  p → P(ядовитый) = 0.9589843578 → ЯДОВИТЫЙ
odor =  s → P(ядовитый) = 0.9975347155 → ЯДОВИТЫЙ
```

# Сравнение результатов с baseline-моделью

```
scikit-learn CategoricalNB:
```

```
Accuracy: 0.949549
```

```
Log-loss: 0.138412
```

```
Байесовская сеть (pgmpy):
```

```
Accuracy: 1.000000
```

```
Log-loss: 0.000003
```

```
=== СРАВНЕНИЕ ===
```

Модель	Accuracy	Log-loss
pgmpy (HillClimb + BDeu)	1.000000	0.000003
CategoricalNB (sklearn)	0.949549	0.138412

Байесовская сеть, построенная с помощью жадного поиска структуры (HillClimbSearch) и BDeu-приора, достигла 100% accuracy и практически нулевого log-loss (0.000003) на тестовом наборе.

Это значительно превосходит классический наивный байесовский классификатор из scikit-learn (accuracy 94.95%, log-loss 0.138).

Превосходство объясняется тем, что алгоритм поиска структуры автоматически выявил ключевой предиктор — признак odor (запах гриба), который является достаточным для детерминированного разделения классов.

Полученная модель не только идеально классифицирует, но и даёт биологически осмысленную и полностью интерпретируемую зависимость: гриб ядовитый  $\Leftrightarrow$  у него есть любой неприятный запах.

# Как сеть моделирует реальные зависимости?

Моя Байесовская сеть моделирует реальные биолого-химические зависимости в датасете.

## Реальная структура, которую нашла сеть

Узел `class` имеет родителя `odor`  
Остальные 21 признак либо вообще не связаны с class напрямую, либо связаны только через odor.

То есть сеть говорит:

Гриб ядовитый  $\Leftrightarrow$  у него есть неприятный запах  $\Leftrightarrow \text{odor} \in \{f, c, y, m, p, s\}$

Все остальные признаки (форма шляпки, цвет спор, жабры и т.д.) — либо следствия запаха, либо вообще не важны.

И это абсолютно соответствует реальной биологии и химии грибов.

## Почему это так?

1) Запах — это летучие вещества (вторичные метаболиты).

- Ядовитые грибы рода *Amanita*, *Chlorophyllum*, *Entoloma* и др. содержат токсины (аматоксины, мускарин, коприн и др.).

- При этом они выделяют характерные летучие соединения → неприятный запах (тухлая рыба, креозот, резкий, плесень и т.д.).

2) Съедобные грибы либо не пахнут, либо пахнут приятно.

- *Agaricus*, *Boletus*, *Cantharellus* и др. → запах отсутствует или анисовый/миндальный (благодаря бензальдегиду и др. безопасным соединениям).

3) Все остальные признаки — вторичны.

- Цвет спор (*spore-print-color*) — часто коррелирует с запахом, потому что зависит от того же вида.

- Цвет жабр (*gill-color*) — тоже видовая особенность.

- Форма ножки, кольцо, вуаль — тоже. Они меняются вместе с запахом, а не являются независимыми причинами ядовитости.