

Причинно-следственный анализ

Артём Копань

Задача



Дана таблица с данными о продажах сети кафе в промежуток времени между 5 месяцами до и 13 месяцами после введения в продажу нового напитка M. Им собираются заменить напиток S.

Требуется:

- 1) Проверить наличие причинно-следственной связи между переменными S_{sales} и M_{sales}
- 2) Восстановить как можно больше аутентичных причинно-следственных связей между переменными

Используемые алгоритмы



- Constraint-based:
 - Алгоритм Петера-Кларка (PC)
 - Fast Causal Inference (FCI)
- Score-based:
 - Алгоритм жадного поиска эквивалентности (Greedy equivalence search, GES)
- Permutation-based
 - Алгоритм Greedy relaxation of the sparsest permutation (GRaSP)

Предобработка данных



```
df = df.drop(['Unnamed: 0', 'id'], axis=1)
```

```
df = df.fillna(0)
```

```
address_values = list(df.address)
address_values = list(map(str, address_values))
address_number_values = list(map(lambda address: int(address.split()[0]), address_values))
```

```
df['address'] = address_number_values
```

Алгоритм Петера-Кларка

Пороговое значение p -уровня значимости было выбрано равным 0.2 (иначе граф получался

```
graph_pc_raw = pc(df.values, alpha=alpha, indep_test='fisherz')
```

Depth=10, working on node 16: 100%  17/17 [00:00<00:00, 1133.13it/s]

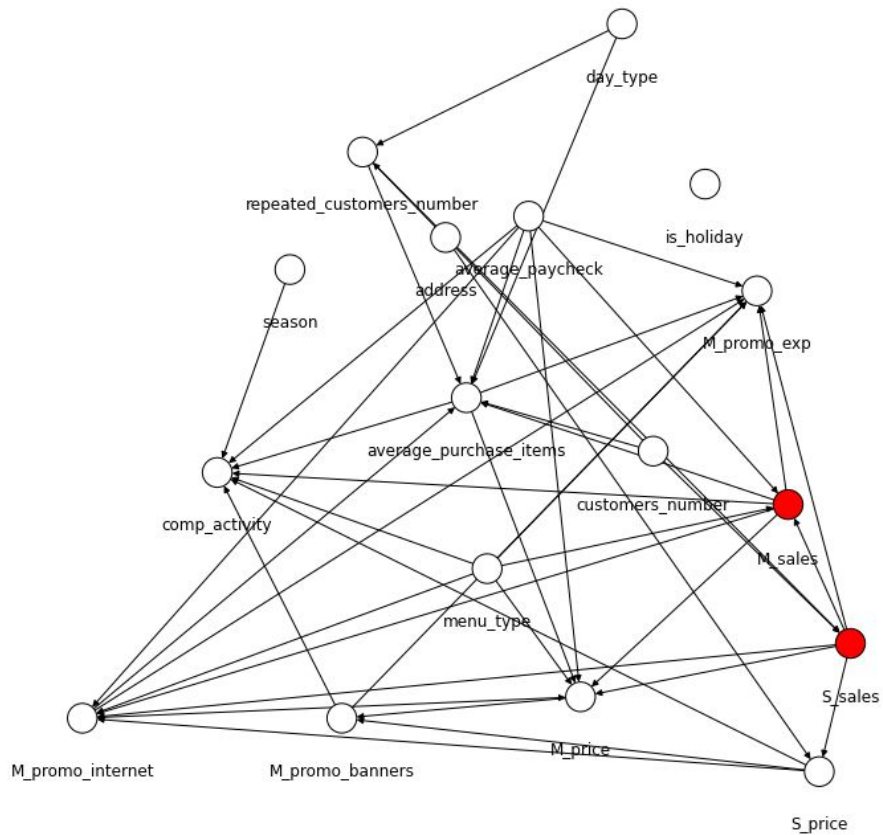
```
graph_pc_raw.to_nx_graph()
```

```
mapping = {node:i for node, i in enumerate(df.columns)}  
mapping_r = {i:node for node, i in enumerate(df.columns)}
```

```
graph_pc = nx.relabel_nodes(graph_pc_raw.nx_graph, mapping)
```

```
edges_to_remove = [(u, v) for u, v in graph_pc.edges() if v == 'is_holiday']  
edges_to_remove += [(u, v) for u, v in graph_pc.edges() if v == 'day_type']  
edges_to_remove += [(u, v) for u, v in graph_pc.edges() if v == 'season']  
edges_to_remove += [(u, v) for u, v in graph_pc.edges() if v == 'address']  
  
graph_pc.remove_edges_from(edges_to_remove)
```

Алгоритм Петера-Кларка



Алгоритм Петера-Кларка



Связь между S_{sales} и M_{sales} есть.

Оценка эффекта: -0.6276425022932131

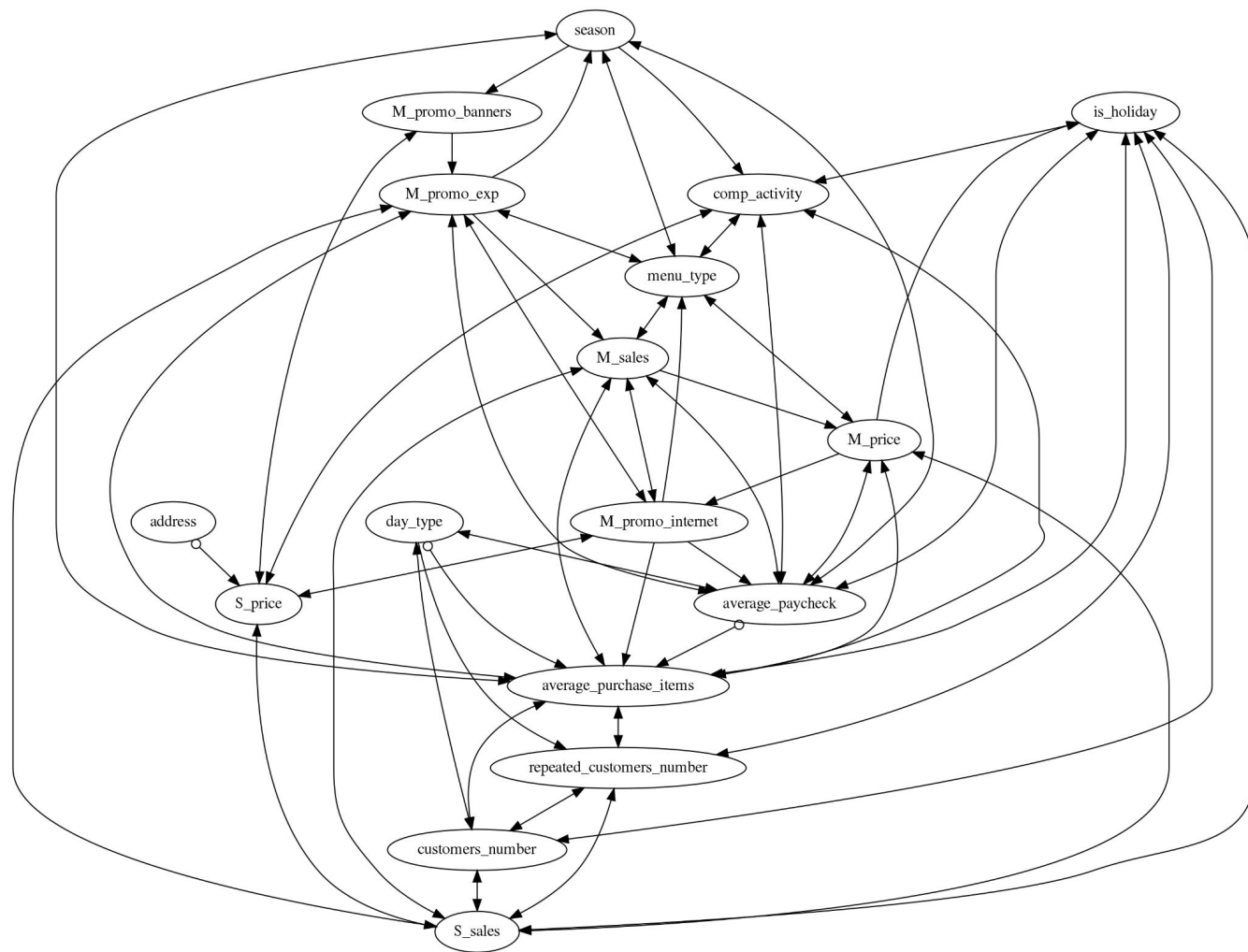
Алгоритм FCI

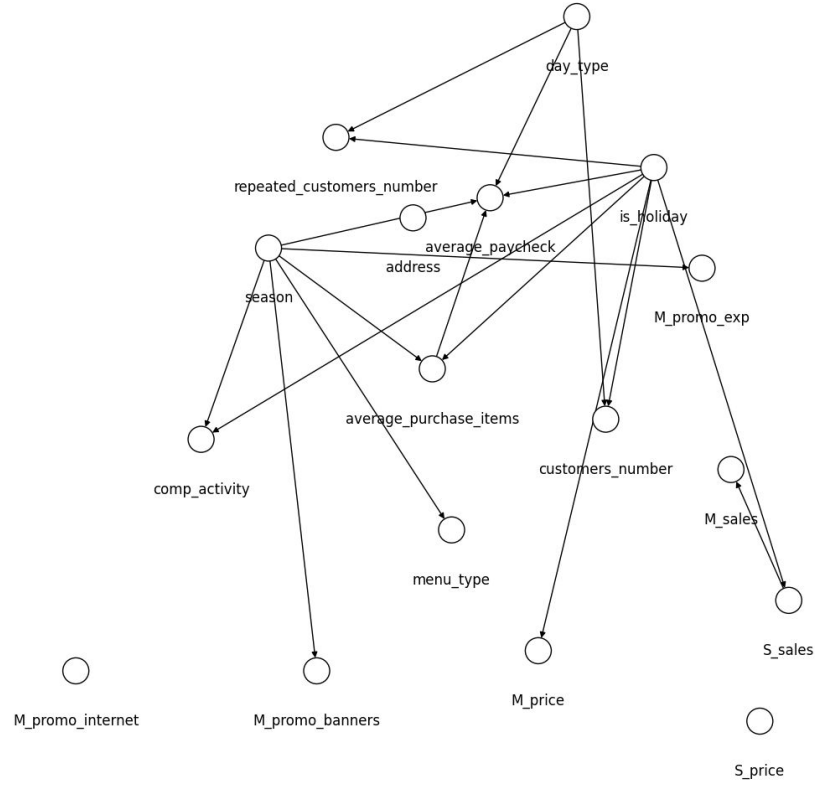
```
from causallearn.search.ConstraintBased.FCI import fci  
fci_graph, fci_edges = fci(df.values, alpha=alpha, indep_test='fisherz')
```

```
0%|          | 0/17 [00:00<?, ?it/s]
```

```
X10 --> X2  
X3 --> X14  
X4 --> X5  
X6 --> X4  
X4 --> X8  
X8 --> X6  
X6 --> X16  
X7 --> X12  
X16 --> X10
```

```
pdy = GraphUtils.to_pydot(fci_graph, labels=df.columns)  
pdy.write_png('fci_graph.png')
```



Алгоритм FCI

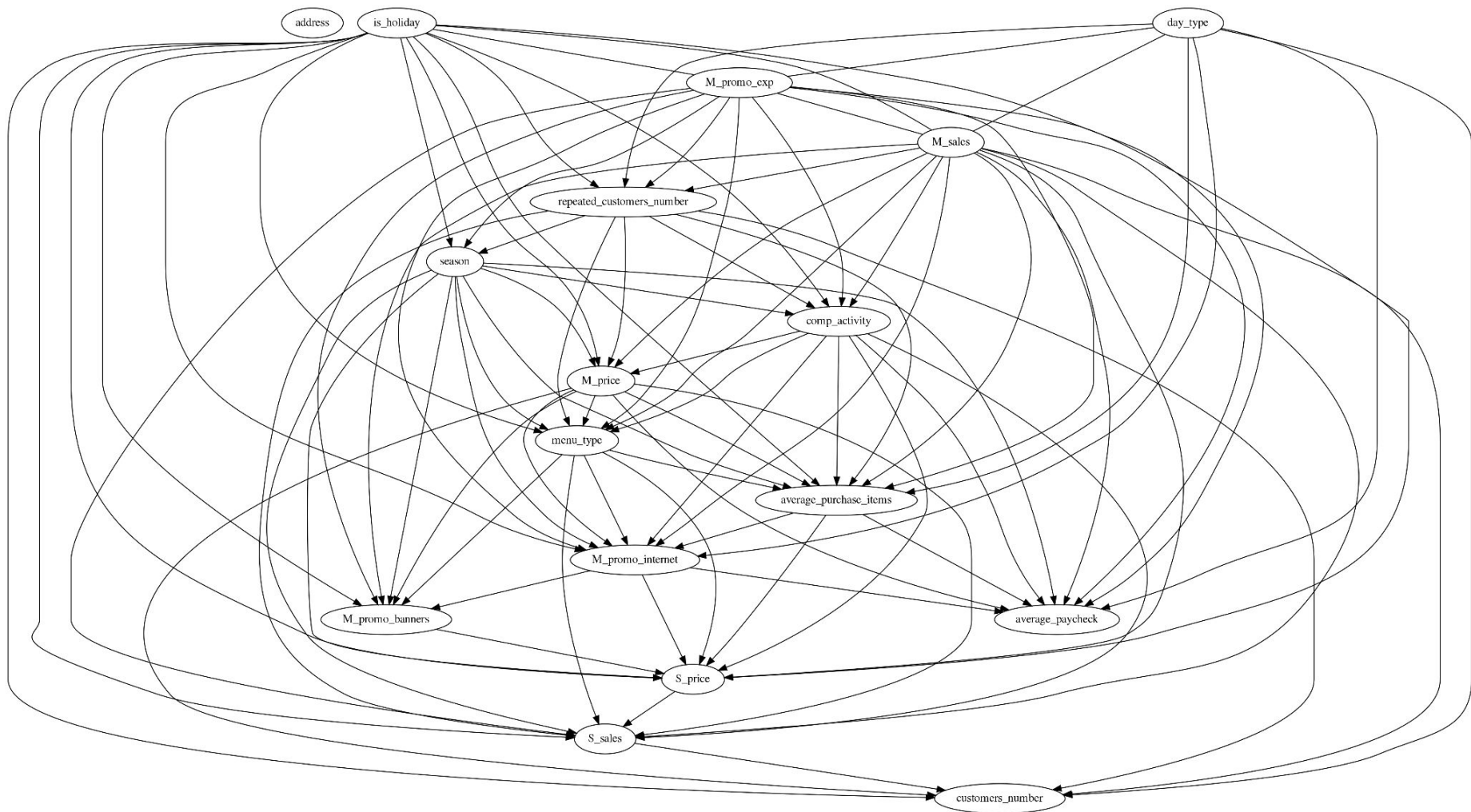


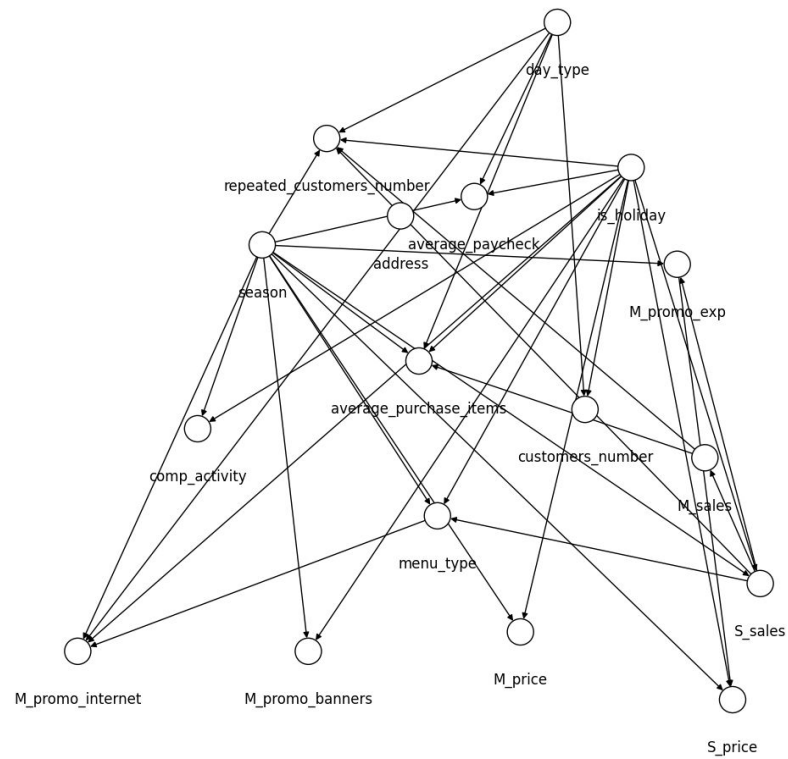
- Есть связь между S_{sales} и M_{sales} , но её характер из графа неясен (ребро двунаправленное)
- Оценка эффекта: -0.9311992463210004

Алгоритм GES

```
from causallearn.search.ScoreBased.GES import ges  
graph_ges = ges(df, score_func='local_score_BIC')
```

```
import io  
from causallearn.utils.GraphUtils import GraphUtils  
import matplotlib.image as mimg  
  
pyd = GraphUtils.to_pydot(graph_ges['G'], labels=df.columns)  
pyd.write_png('graph_ges.png')
```





Алгоритм GES



- Причинно-следственная связь между S_{sales} и M_{sales} есть
- Оценка эффекта: -0.9311992463210004

Алгоритм GRaSP



```
from causallearn.search.PermutationBased.GRaSP import grasp
```

```
G = grasp(df.values)
```

```
GRaSP edge count: 85
```

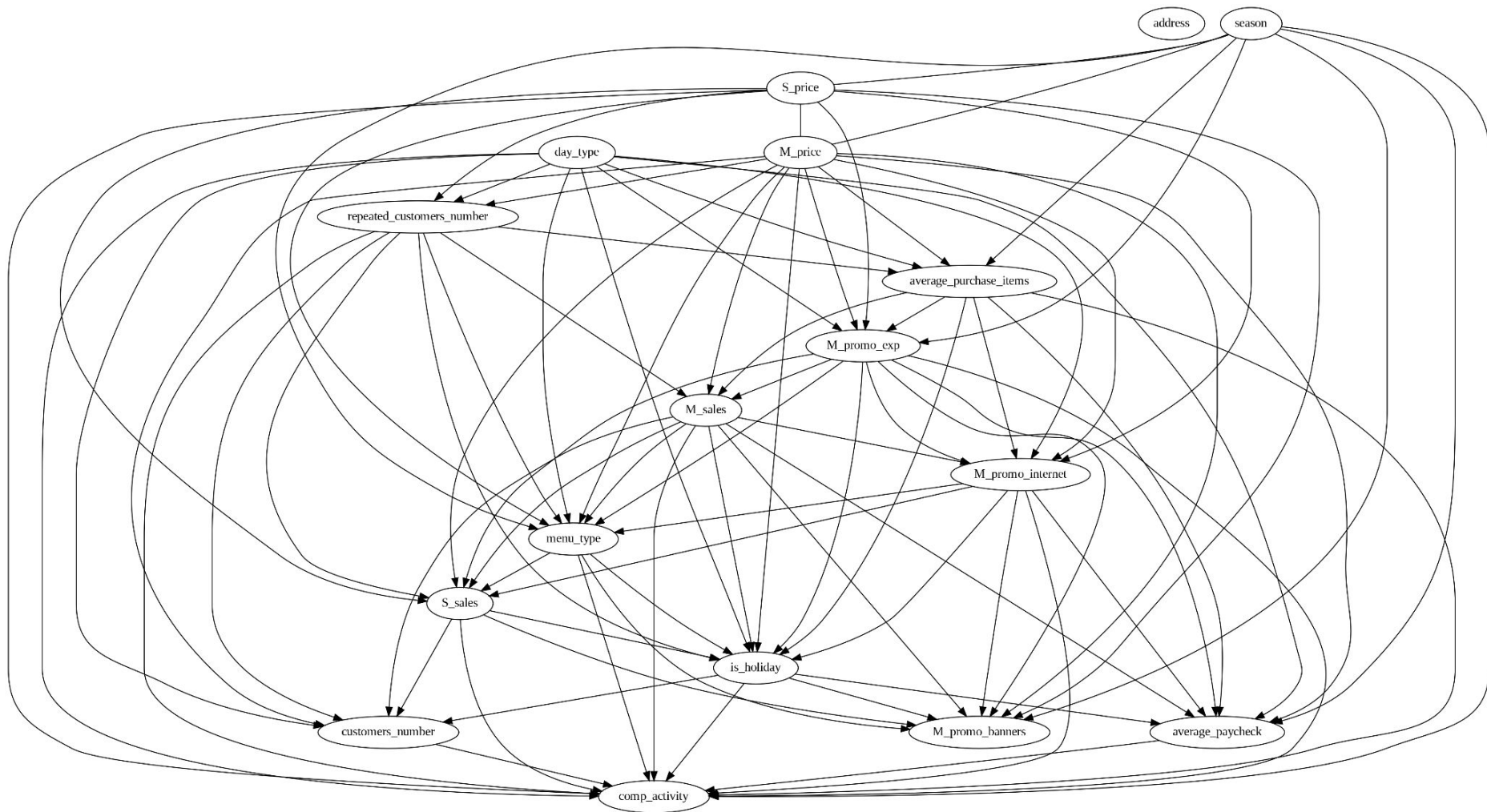
```
GRaSP completed in: 722.11s
```

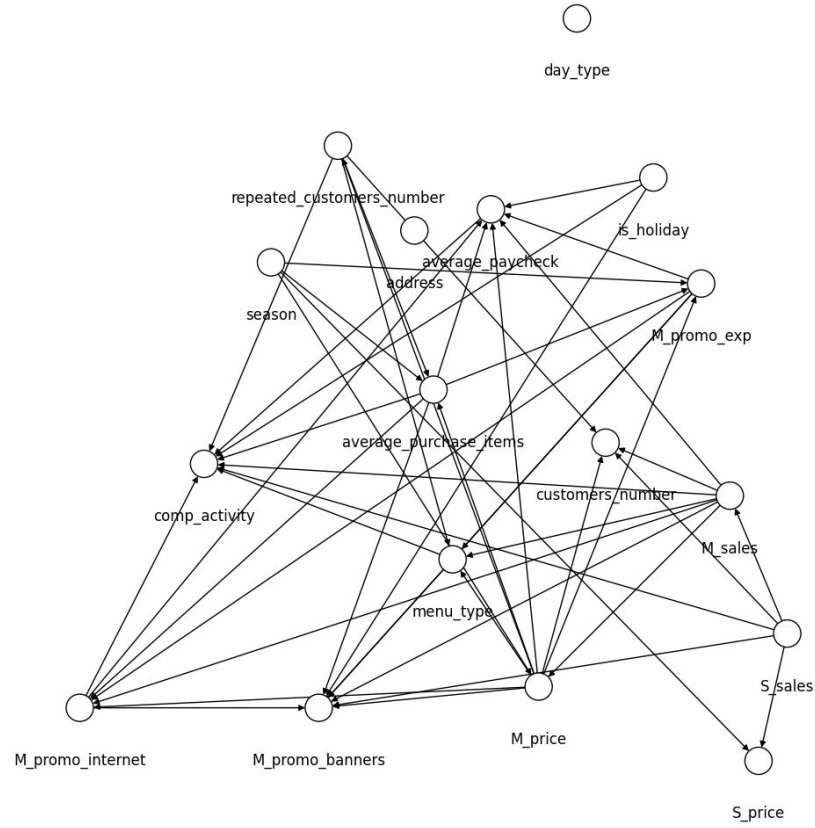
```
import io
```

```
from causallearn.utils.GraphUtils import GraphUtils
```

```
pyd = GraphUtils.to_pydot(G, labels=df.columns)
```

```
pyd.write_png('graph_grasp.png')
```



Алгоритм GRaSP



- Имеется связь между S_{sales} и M_{sales}
- Оценка эффекта: -0.9311992463210004

Desbordante



Desbordante — это библиотека для поиска функциональных зависимостей в данных. Я решил применить её в данной задаче.

Находим функциональные зависимости при помощи Desbordante, составляем на их основе граф и применяем background knowledge.

Desbordante

```
import desbordante

TABLE = 'ci_data_internship_2.csv'

algo = desbordante.fd.algorithms.Default()
algo.load_data(table=TABLE, ',', True)
algo.execute()
result = algo.get_fds()
print('FDs:')
for fd in result:
    print(fd)
```

```
adjacency_matrix = [[0 for j in range(len(df.columns))] for i in range(len(df.columns))]
nodes = list(df.columns)
d = {nodes[i]: i for i in range(len(df.columns))}

for fd in result:
    lhs, rhs = str(fd).split(' -> ')
    lhs = lhs[1:-1].split()
    if rhs != '' and rhs != '' and len(lhs) != 0:
        for node in lhs:
            if node == '':
                continue
            adjacency_matrix[d[node]][d[rhs]] = 1

graph = nx.from_numpy_array(np.array(adjacency_matrix), create_using=nx.DiGraph)
```

Desbordante

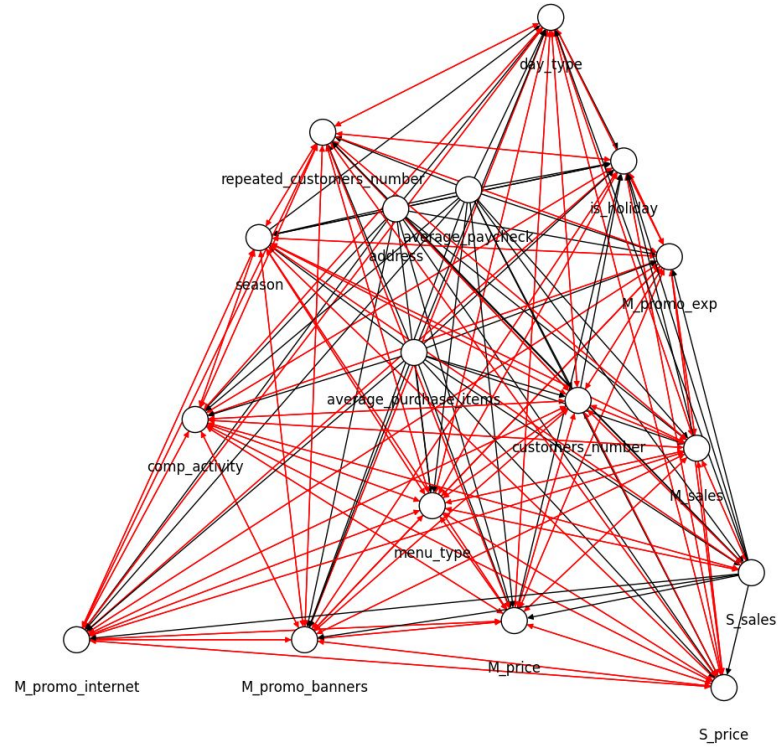


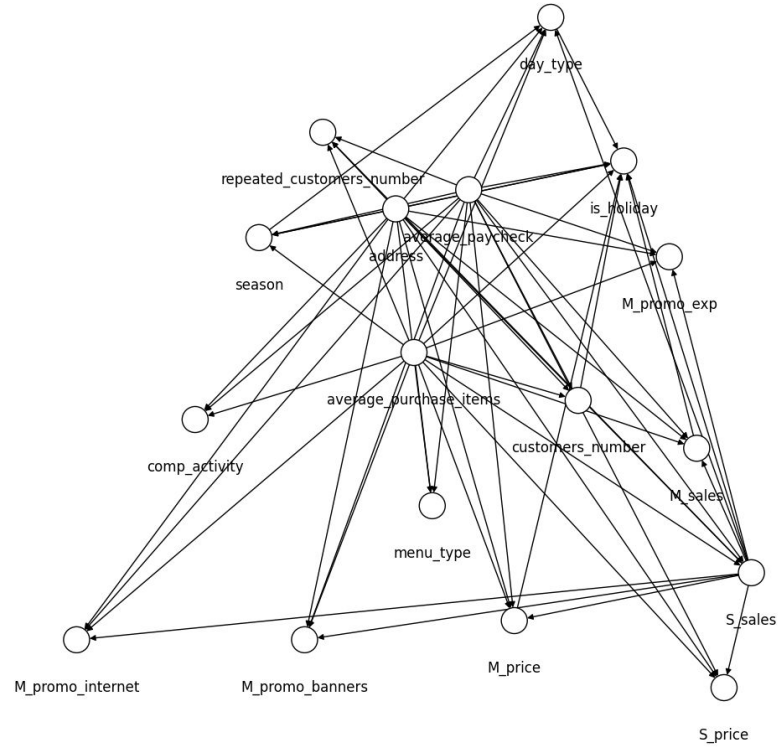
Background knowledge:

```
edges_to_remove = [(u, v) for u, v in graph.edges() if v == 'is_holiday']
edges_to_remove += [(u, v) for u, v in graph.edges() if v == 'day_type']
edges_to_remove += [(u, v) for u, v in graph.edges() if v == 'season']
edges_to_remove += [(u, v) for u, v in graph.edges() if v == 'address']
graph.remove_edges_from(['M_sales', 'S_sales'])
edges_to_remove = [(u, v) for u, v in graph.edges if (u, v) in get_non_directed_edges(graph)]

graph.remove_edges_from(edges_to_remove)
```

- Имеется связь между S_sales и M_sales
- Оценка эффекта: -0.6085022649245388





Выводы



- Все алгоритмы обнаруживают причинно-следственную связь между S_{sales} и M_{sales}
- Эта связь отрицательна, а оценка её эффекта при использовании разных алгоритмов колеблется в диапазоне $[-0.932, -0.609]$.
- Алгоритмы в большинстве случаев хорошо восстанавливают причинно-следственные связи между различными переменными, но требуется удаление нелогичных связей с помощью background knowledge.