

FEDERAL STATE AUTONOMOUS EDUCATIONAL INSTITUTION
OF HIGHER EDUCATION
ITMO UNIVERSITY

REPORT

on learning practice No.3

«Sampling of multivariate random variables»

Performed by:
Putnikov Semyon
Gabrielian Mikhail

St. Petersburg

2021

code

November 26, 2021

```
[32]: # imports
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.preprocessing import KBinsDiscretizer
from sklearn.metrics import accuracy_score
from pgmpy.estimators import HillClimbSearch
from pgmpy.estimators import K2Score, BicScore
from pgmpy.models import BayesianNetwork
from pgmpy.inference import VariableElimination
import networkx as nx
from pgmpy.base import DAG
from pgmpy.models import BayesianModel
from pgmpy.sampling import BayesianModelSampling

import warnings
warnings.filterwarnings("ignore")
```

```
[2]: source_data_path = "../tcs_stock.csv"
row_df = pd.read_csv(source_data_path)
```

0.1 Step 1. Choose variables for sampling from your dataset (overall – about 10 variables, 3-4 –target variables, the rest - predictors).

```
[3]: row_df.head()
```

```
[3]:
```

	Date	Symbol	Series	Prev Close	Open	High	Low	Last	\
0	2015-01-01	TCS	EQ	2558.25	2567.0	2567.00	2541.00	2550.00	
1	2015-01-02	TCS	EQ	2545.55	2551.0	2590.95	2550.60	2588.40	
2	2015-01-05	TCS	EQ	2579.45	2581.0	2599.90	2524.65	2538.10	
3	2015-01-06	TCS	EQ	2540.25	2529.1	2529.10	2440.00	2450.05	
4	2015-01-07	TCS	EQ	2446.60	2470.0	2479.15	2407.45	2426.90	

	Close	VWAP	Volume	Turnover	Trades	Deliverable Volume	\
0	2545.55	2548.51	183415	4.674345e+13	8002	52870	
1	2579.45	2568.19	462870	1.188740e+14	27585	309350	

2	2540.25	2563.94	877121	2.248886e+14	43234	456728
3	2446.60	2466.90	1211892	2.989615e+14	84503	714306
4	2417.70	2433.96	1318166	3.208362e+14	101741	886368

	%Deliverble
0	0.2883
1	0.6683
2	0.5207
3	0.5894
4	0.6724

Target: High, Low, Trades ##### Predictors: Prev Close, Open, Close, VWAP, Volume, Deliverable Volume

```
[4]: df = row_df[['High', 'Low', 'Trades', 'Prev Close', 'Open', 'Close', 'VWAP', 'Volume', 'Deliverable Volume', 'Turnover']]
df.head()
```

```
[4]:      High      Low  Trades  Prev Close   Open   Close   VWAP   Volume \
0  2567.00  2541.00    8002    2558.25  2567.0  2545.55  2548.51   183415
1  2590.95  2550.60   27585    2545.55  2551.0  2579.45  2568.19   462870
2  2599.90  2524.65   43234    2579.45  2581.0  2540.25  2563.94   877121
3  2529.10  2440.00   84503    2540.25  2529.1  2446.60  2466.90  1211892
4  2479.15  2407.45  101741    2446.60  2470.0  2417.70  2433.96  1318166
```

	Deliverable Volume	Turnover
0	52870	4.674345e+13
1	309350	1.188740e+14
2	456728	2.248886e+14
3	714306	2.989615e+14
4	886368	3.208362e+14

Distributions:

High - hypsecant with param 1.0910

Low - laplace with param 1.3288

Trades - gumbel_r with param 3.3611

Prev Close- laplace with parameter 7.3622

Open - laplace with parameter 8.493

Close - laplace with parameter 7.6574

VWAP - laplace with parameter 1.1143

Volume - johnsonsb with parameter 5.2133

Deliverable Volume - invgamma with parameter 9.9473

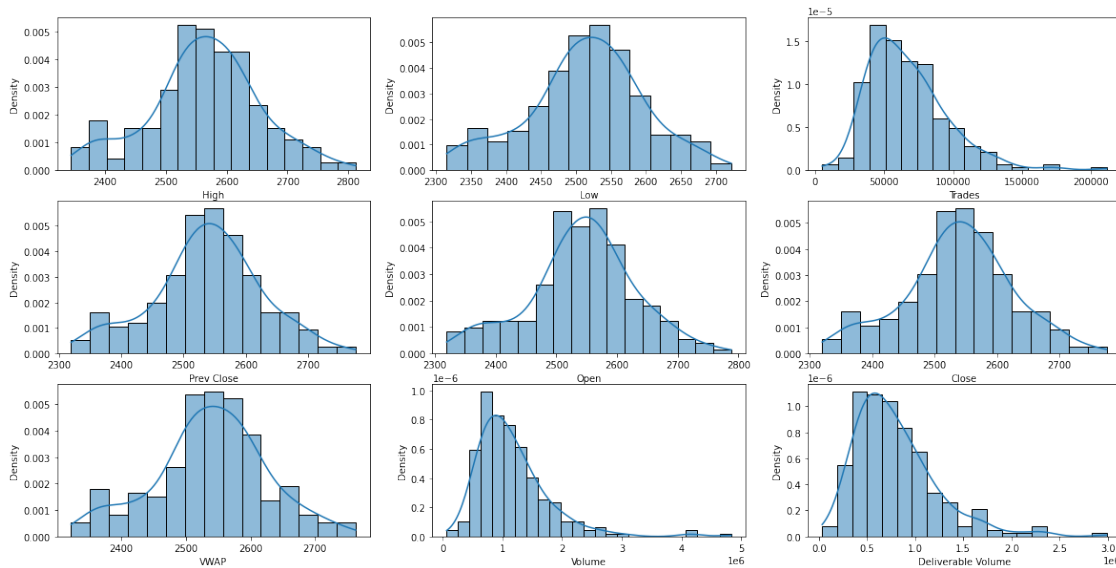
```
[5]: fig, ax = plt.subplots(3, 3, figsize=(20,10))
sns.histplot(df['High'], ax=ax[0,0], kde=True, stat="density")
sns.histplot(df['Low'], ax=ax[0,1], kde=True, stat="density")
sns.histplot(df['Trades'], ax=ax[0,2], kde=True, stat="density")
```

```

sns.histplot(df['Prev Close'], ax=ax[1,0], kde=True, stat="density")
sns.histplot(df['Open'], ax=ax[1,1], kde=True, stat="density")
sns.histplot(df['Close'], ax=ax[1,2], kde=True, stat="density")
sns.histplot(df['VWAP'], ax=ax[2,0], kde=True, stat="density")
sns.histplot(df['Volume'], ax=ax[2,1], kde=True, stat="density")
sns.histplot(df['Deliverable Volume'], ax=ax[2,2], kde=True, stat="density")

plt.show()

```



0.2 Step 2. Using univariate parametric distributions that were fitted in Lab#2 make sampling of chosen target variables. Use for this 2 different sampling methods.

```

[6]: labels = ['High', 'Low', 'Trades']

for l in labels:
    fig, axes = plt.subplots(1, 3, figsize=(20, 3))
    sns.kdeplot(ax=axes[0], data=df, x=l, shade=True)

    # Random sample
    sample_data = df.sample(frac=.5)
    sns.kdeplot(ax=axes[1], data=sample_data, x=l, shade=True)

    # Stratified random sample
    mean = df[l].mean()

    sample_data_0 = df[df[l] > mean].sample(n=60, axis=0)
    sample_data_1 = df[df[l] <= mean].sample(n=60, axis=0)

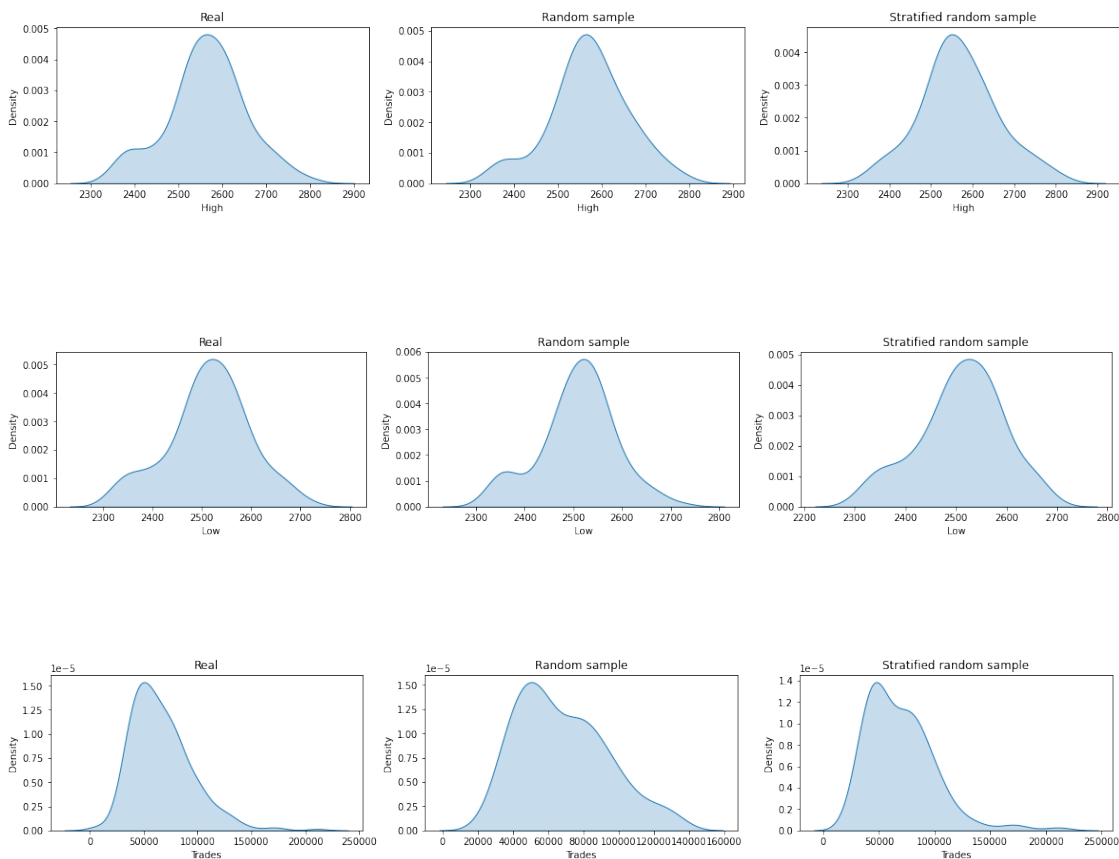
```

```

sample_data = pd.concat([sample_data_0, sample_data_1])
sns.kdeplot(ax=axes[2], data=sample_data, x=1, shade=True)

axes[0].set_title('Real')
axes[1].set_title('Random sample')
axes[2].set_title('Stratified random sample')
plt.show()

```



0.3 Step 3. Estimate relations between predictors and chosen target variables. At least, they should have significant correlation coefficients.

```

[7]: predictors = df[['Prev Close', 'Open', 'Close', 'VWAP', 'Volume', 'Deliverable_
    ↳Volume']]
    targets = df[['High', 'Low', 'Trades']]

[8]: figure, ax = plt.subplots(1, 1, figsize=(10, 10))
    sns.set_theme(style='whitegrid', palette='pastel')

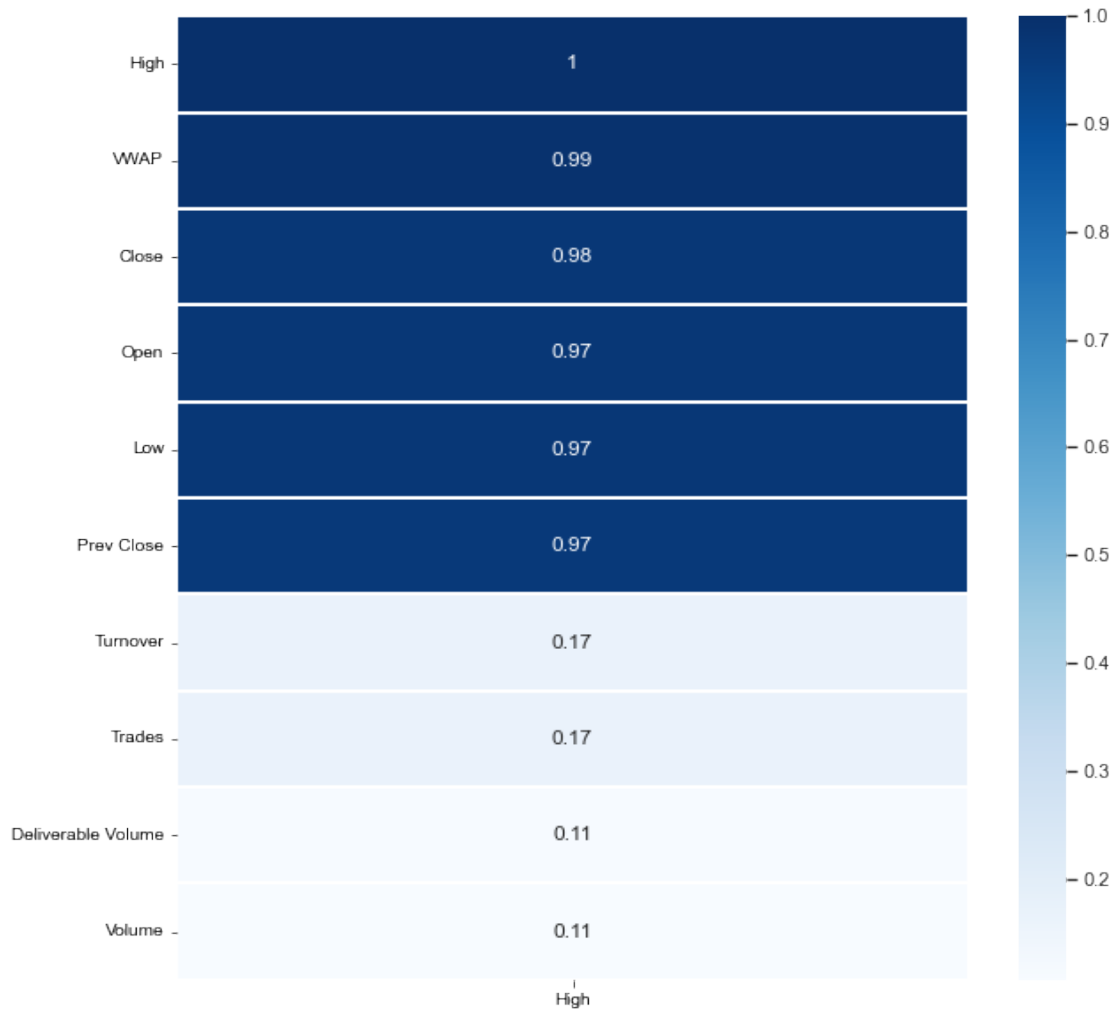
    sns.heatmap(
        df.corr()[['High']].sort_values(by='High', ascending=False),

```

```

cmap='Blues',
annot=True,
linewidths=0.25
)
plt.show()

```



We have got troubles with correlation of Volume and Deliverable Volume. So we will remove this predictors from dataset. Other predictors with ban correlation will be used, because we have small number of columns $\backslash () /$

Another point - use Trades feature as target. But it has a problem - only Turnover has a great correlation with it. We can check out make decision after all.

```

[9]: df = row_df[['High', 'Low', 'Trades', 'Prev Close', 'Open', 'Close', 'VWAP',
↪ 'Turnover']]

```

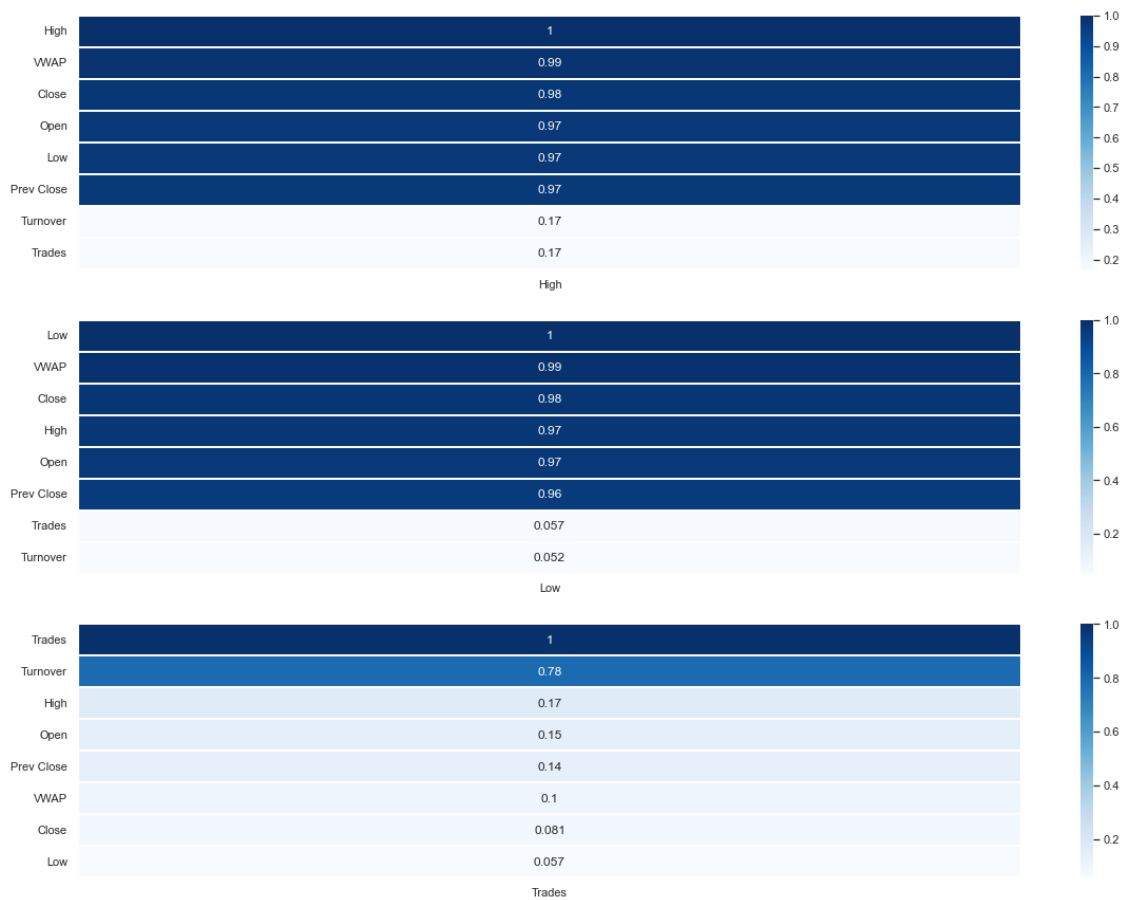
```
[10]: labels = ['High', 'Low', 'Trades']

figure, ax = plt.subplots(3, 1, figsize=(20, 15))
sns.set_theme(style='whitegrid', palette='pastel')

count = 0

for l in labels:
    sns.heatmap(
        df.corr()[[l]].sort_values(by=l, ascending=False),
        cmap='Blues',
        annot=True,
        linewidths=0.25,
        ax=ax[count]
    )
    count += 1

plt.show()
```



0.4 Step 4. Build a Bayesian network for chosen set of variables. Choose its structure on the basis of multivariate analysis and train distributions in nodes using chosen algorithm.

```
[11]: df_transformed = df.copy()
discretizer = KBinsDiscretizer(n_bins=10, encode='ordinal', strategy='kmeans')
df_discretized = discretizer.fit_transform(df.values[:])
```

```
[12]: df_transformed[:] = df_discretized
df_transformed
```

```
[12]:
```

	High	Low	Trades	Prev	Close	Open	Close	VWAP	Turnover
0	4.0	5.0	0.0		5.0	5.0	5.0	5.0	0.0
1	5.0	5.0	1.0		5.0	5.0	5.0	5.0	1.0
2	5.0	5.0	1.0		5.0	6.0	4.0	5.0	2.0
3	3.0	3.0	3.0		4.0	4.0	2.0	3.0	3.0
4	2.0	2.0	4.0		2.0	3.0	2.0	2.0	3.0
..
243	2.0	2.0	0.0		2.0	2.0	2.0	2.0	1.0
244	2.0	2.0	1.0		2.0	2.0	2.0	3.0	5.0
245	2.0	3.0	1.0		2.0	2.0	2.0	2.0	2.0
246	2.0	2.0	2.0		2.0	2.0	2.0	2.0	2.0
247	2.0	2.0	1.0		2.0	2.0	2.0	2.0	1.0

[248 rows x 8 columns]

```
[13]: blacklist = [(x, y) for x in df_transformed.columns.to_list() for y in
    ↳ ['Turnover', 'Trades'] if x != y]
```

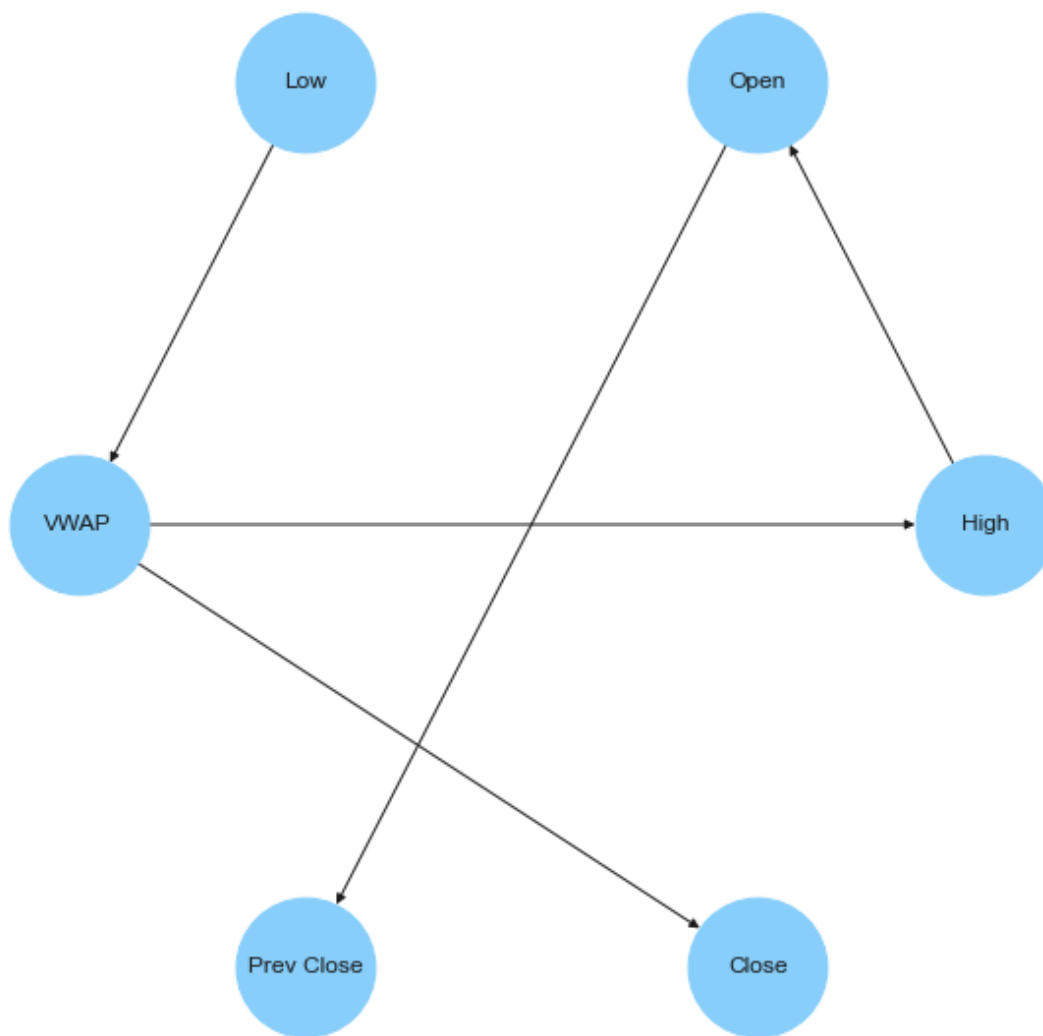
0.5 K2Score

```
[14]: hc_search = HillClimbSearch(data=df_transformed)
model_k2 = hc_search.estimate(scoring_method=K2Score(df_transformed),
    ↳ black_list=blacklist, show_progress=False)
```

```
[15]: figure, ax = plt.subplots(1, 1, figsize=(10, 10))

graph = nx.DiGraph()
graph.add_edges_from(model_k2.edges())
positions = nx.layout.circular_layout(graph)
nx.draw(graph, positions, with_labels=True, node_color='lightskyblue',
    ↳ node_size=5000)

plt.show()
```

```
[16]: def accuracy_params_restoration(bn: BayesianNetwork, data: pd.DataFrame):
    bn.fit(data)
    result = pd.DataFrame(columns=['Parameter', 'accuracy'])
    bn_infer = VariableElimination(bn)
    for j, param in enumerate(data.columns):
        accuracy = 0
        test_param = data[param].copy()
        test_data = data.drop(columns=param)
        evidence = test_data.to_dict('records')
        predicted_param = []
        for element in evidence:
            prediction = bn_infer.map_query(variables=[param],
            ↪evidence=element, show_progress=False)
```

```

        predicted_param.append(prediction[param])
    accuracy = accuracy_score(test_param.values, predicted_param)
    result.loc[j, 'Parameter'] = param
    result.loc[j, 'accuracy'] = accuracy
    return result

```

```

[17]: accuracy_params_restoration(BayesianNetwork(model_k2.edges()),
                                df_transformed[['High', 'Low', 'Prev_
↪Close', 'Open', 'Close', 'VWAP']]).sample(frac=0.3))

```

```

[17]:
   Parameter  accuracy
0        High    0.72973
1         Low    0.77027
2  Prev Close    0.797297
3         Open    0.905405
4        Close    0.716216
5         VWAP    0.851351

```

0.6 Step 5. Build a Bayesian network for the same set of variables but using 2 chosen algorithms for structural learning.

0.6.1 K2 algorithm implemented above!

0.7 BicScore

```

[18]: model_bic = hc_search.estimate(scoring_method=BicScore(df_transformed),
↪black_list=blacklist, show_progress=False)

```

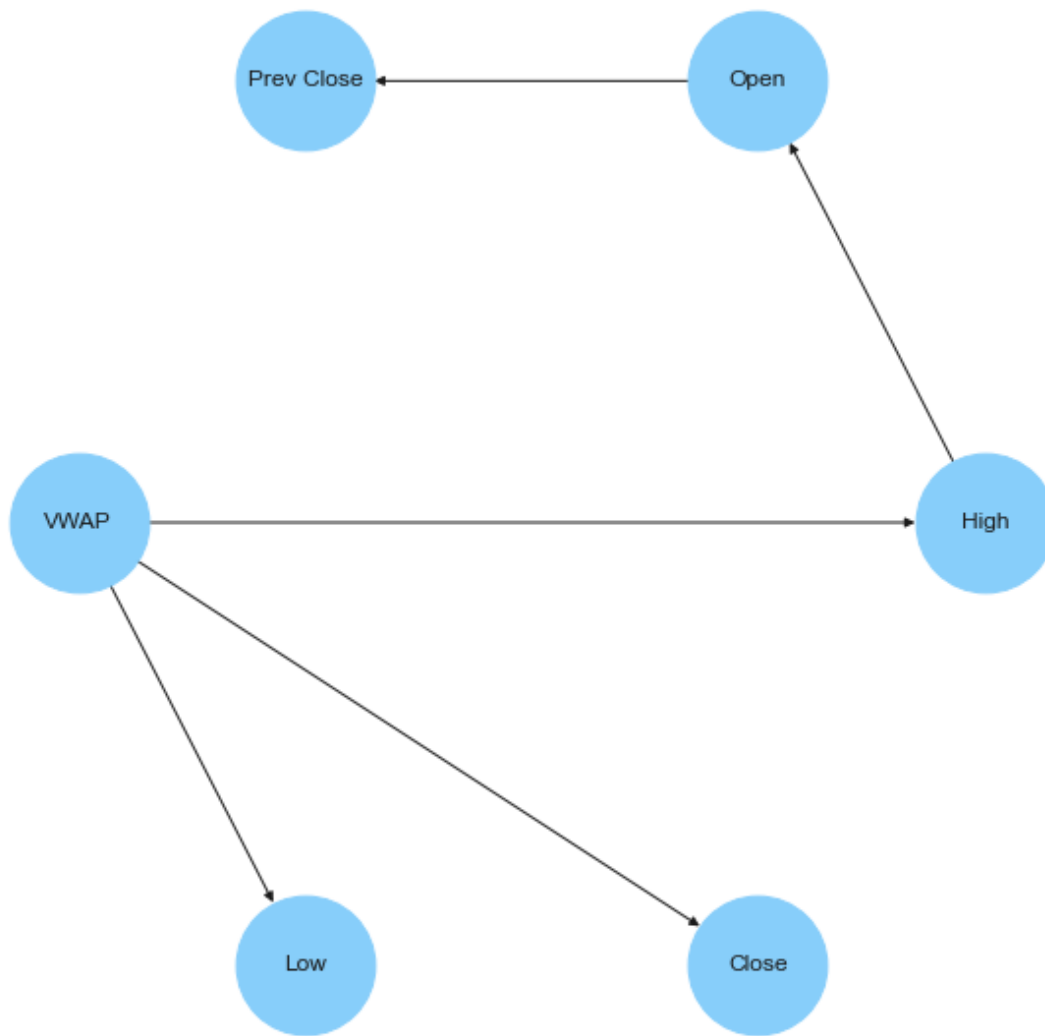
```

[19]: figure, ax = plt.subplots(1, 1, figsize=(10, 10))

graph = nx.DiGraph()
graph.add_edges_from(model_bic.edges())
positions = nx.layout.circular_layout(graph)
nx.draw(graph, positions, with_labels=True, node_color='lightskyblue',
↪node_size=5000)

plt.show()

```



Step 6. Analyze a quality of sampled target variables from the point of view of problem statement (e.g. prediction, gap filling, synthetic generation).

```
[29]: def sampling (bn: DAG, data: pd.DataFrame, n: int = 100):
    bn_new = BayesianModel(bn.edges())
    bn_new.fit(data)
    sampler = BayesianModelSampling(bn_new)
    sample = sampler.forward_sample(size=n)
    return sample
```

```
[30]: def draw_comparative_hist (parametr: str, original_data: pd.DataFrame,
    ↪data_sampled: pd.DataFrame, axes=None):
    final_df = pd.DataFrame()
```

```

df1 = pd.DataFrame()
df1[parametr] = original_data[parametr]
df1['Data'] = 'Original data'
df1['Probability'] = df1[parametr].apply(lambda x: (df1.
→groupby(parametr)[parametr].count()[x])/original_data.shape[0])

df2 = pd.DataFrame()
df2[parametr] = data_sampled[parametr]
df2['Data'] = 'Synthetic data'
df2['Probability'] = df2[parametr].apply(lambda x: (df2.
→groupby(parametr)[parametr].count()[x])/data_sampled.shape[0])
final_df = pd.concat([df1, df2])

sns.barplot(ax=axes, x=parametr, y="Probability", hue="Data", data=final_df)

```

```

[33]: sample_K2 = sampling(model_k2, df_transformed, df_transformed.shape[0])
sample_Bic = sampling(model_bic, df_transformed, df_transformed.shape[0])

```

```
HBox(children=(IntProgress(value=0, max=6), HTML(value='')))
```

```
HBox(children=(IntProgress(value=0, max=6), HTML(value='')))
```

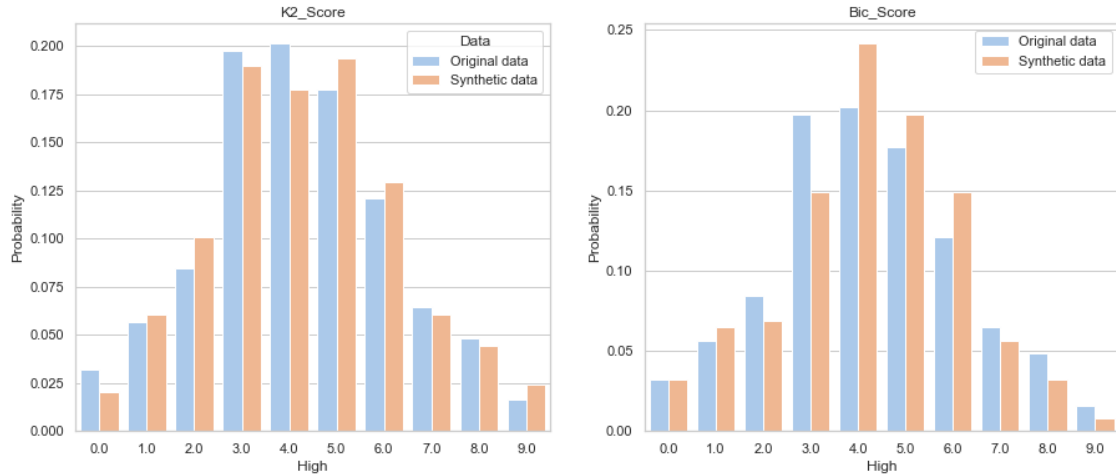
```

[38]: fig, axes = plt.subplots(1, 2, figsize=(15, 6))

draw_comparative_hist('High', df_transformed, sample_K2, axes=axes[0])
draw_comparative_hist('High', df_transformed, sample_Bic, axes=axes[1])

axes[0].set_title('K2_Score')
axes[1].set_title('Bic_Score')
plt.legend()
plt.show()

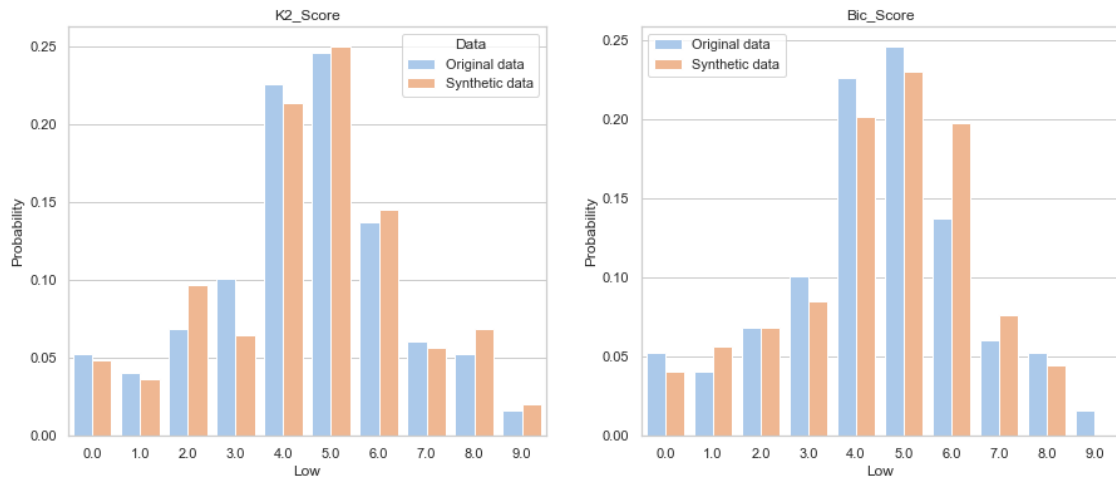
```



```
[40]: fig, axes = plt.subplots(1, 2, figsize=(15, 6))

draw_comparative_hist('Low', df_transformed, sample_K2, axes=axes[0])
draw_comparative_hist('Low', df_transformed, sample_Bic, axes=axes[1])

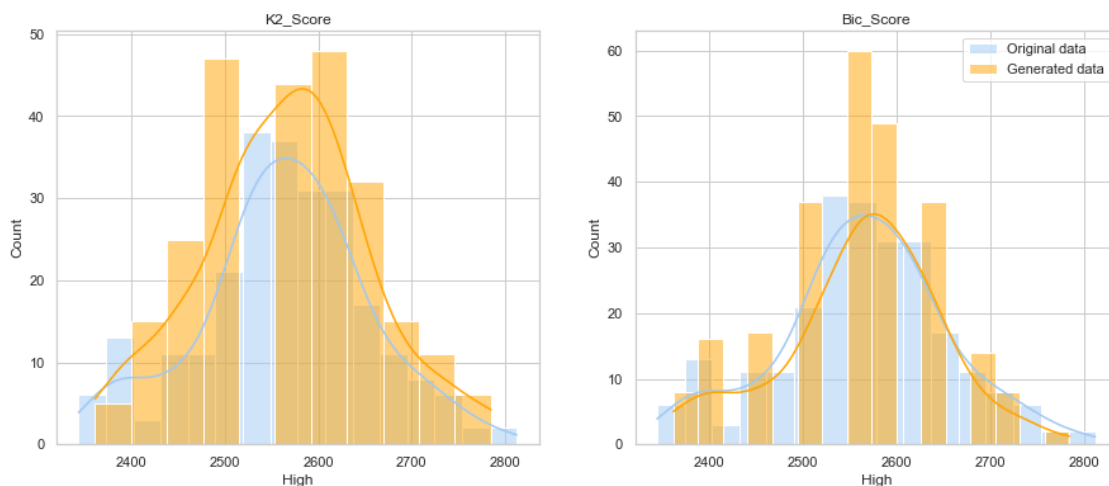
axes[0].set_title('K2_Score')
axes[1].set_title('Bic_Score')
plt.legend()
plt.show()
```



```
[42]: for in_blacklist in list(set(df_transformed.columns) - set(sample_K2.columns)):
    sample_K2[in_blacklist] = df_transformed[in_blacklist].copy()
    sample_Bic[in_blacklist] = df_transformed[in_blacklist].copy()
```

```
[44]: sample_K2[sample_K2.columns] = discretizer.  
      ↪inverse_transform(sample_K2[sample_K2.columns].values)  
      sample_Bic[sample_K2.columns] = discretizer.  
      ↪inverse_transform(sample_Bic[sample_K2.columns].values)
```

```
[47]: fig, axes = plt.subplots(1, 2, figsize=(15, 6))  
  
sns.histplot(df['High'], label='Original data', kde=True, ax=axes[0])  
sns.histplot(sample_K2['High'], label='Generated data', kde=True,   
      ↪color='orange', ax=axes[0])  
  
sns.histplot(df['High'], label='Original data', kde=True, ax=axes[1])  
sns.histplot(sample_Bic['High'], label='Generated data', kde=True,   
      ↪color='orange', ax=axes[1])  
  
axes[0].set_title('K2_Score')  
axes[1].set_title('Bic_Score')  
plt.legend()  
plt.show()
```



```
[55]: fig, axes = plt.subplots(1, 2, figsize=(15, 6))  
  
sns.histplot(df['Prev Close'], label='Original data', kde=True, ax=axes[0])  
sns.histplot(sample_K2['Prev Close'], label='Generated data', kde=True,   
      ↪color='orange', ax=axes[0])  
  
sns.histplot(df['Prev Close'], label='Original data', kde=True, ax=axes[1])  
sns.histplot(sample_Bic['Prev Close'], label='Generated data', kde=True,   
      ↪color='orange', ax=axes[1])
```

```

axes[0].set_title('K2_Score')
axes[1].set_title('Bic_Score')
plt.legend()
plt.show()

```

