

FEDERAL STATE AUTONOMOUS EDUCATIONAL INSTITUTION OF HIGHER
EDUCATION
ITMO UNIVERSITY

Report on learning practice #3
“Sampling of multivariate random variables“

Performed by
Tyulkov Nikita

St. Petersburg
2021

1. Substantiation of chosen sampling.

On the first step we have to choose target and predictor variables. From avocado dataset the next variables was chosen:

large_bags	target
xlarge_bags	target
4770	target
4046	predictor
4225	predictor
total_volume	predictor
small_bags	predictor
average_price	predictor
total_bags	predictor

2. Sampling of chosen target variables using univariate parametric distributions (from practice 2) with 2 different sampling methods.

For the sampling two algorithms were chosen: inverse transform and accept-reject samplings. On the Figure 1 we can see the results of inverse transform sampling algorithm.

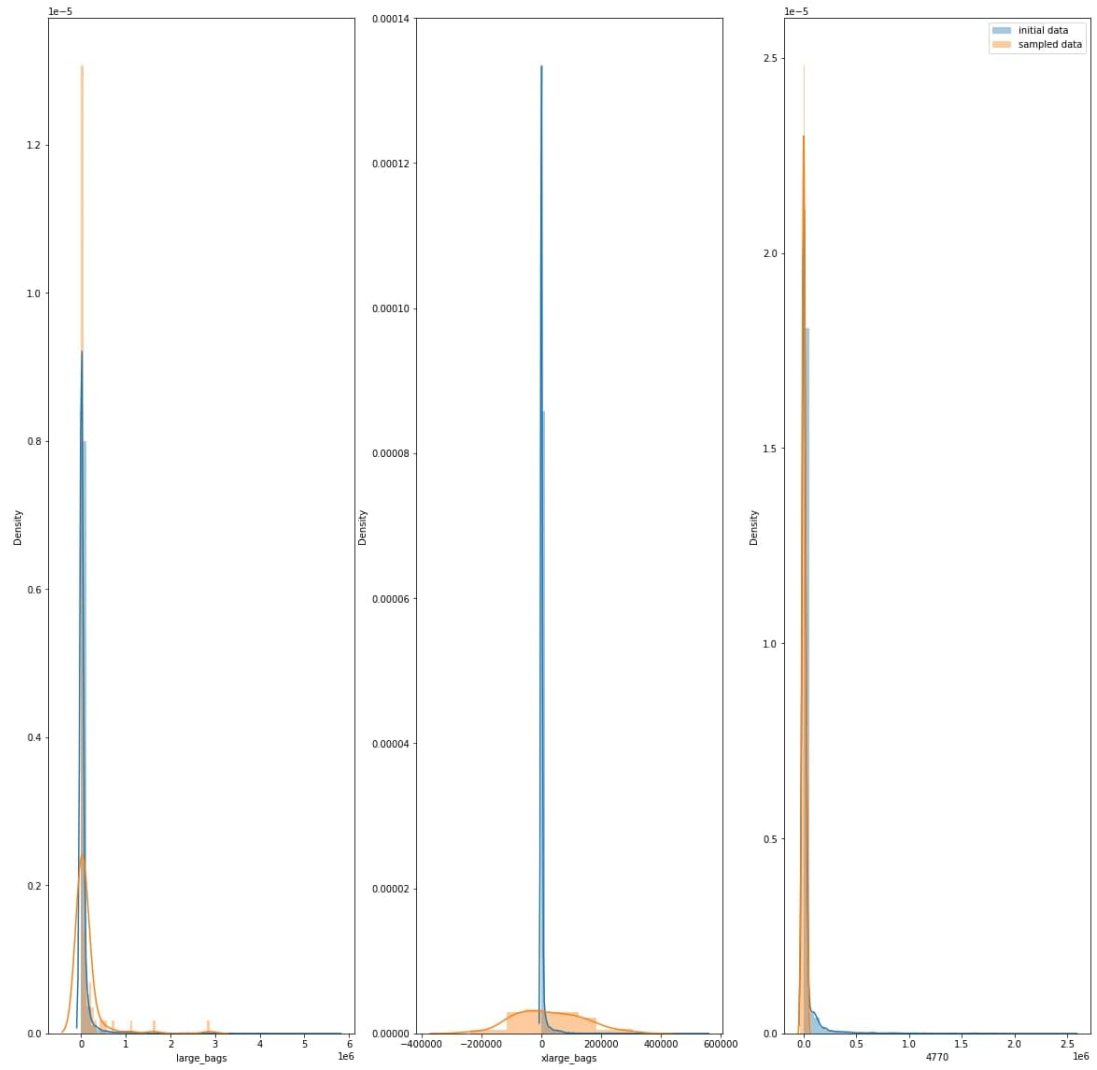


Figure 1: Inverse transform sampling algorithm

On the Figures 2,3,4 we can see the results of accept-reject sampling algorithm. By red line we mean KDE and by green - our sampling.

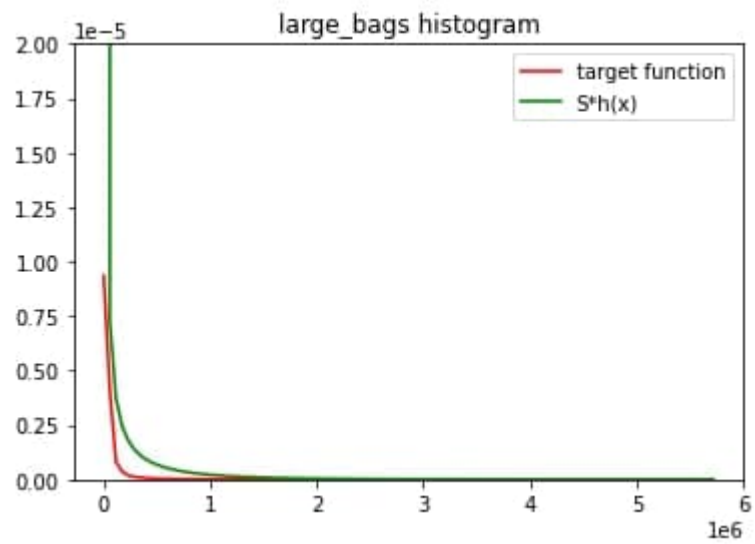


Figure 2: Accept-reject sampling algorithm (large_bags)

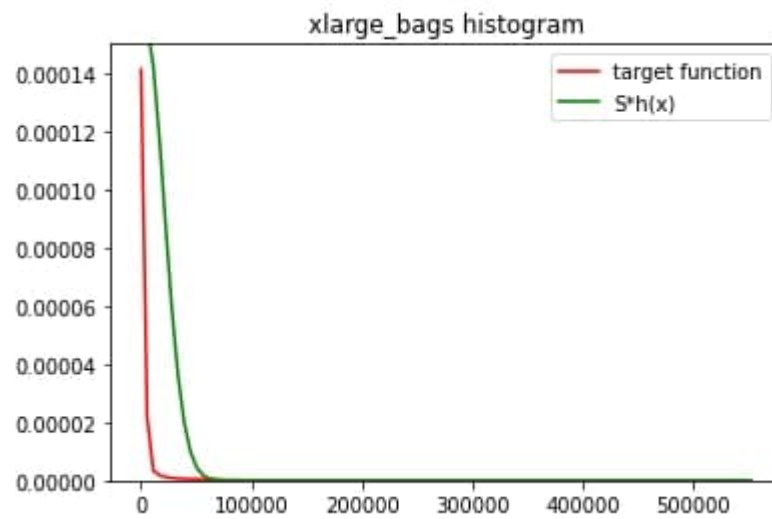


Figure 3: Accept-reject sampling algorithm (xlarge_bags)

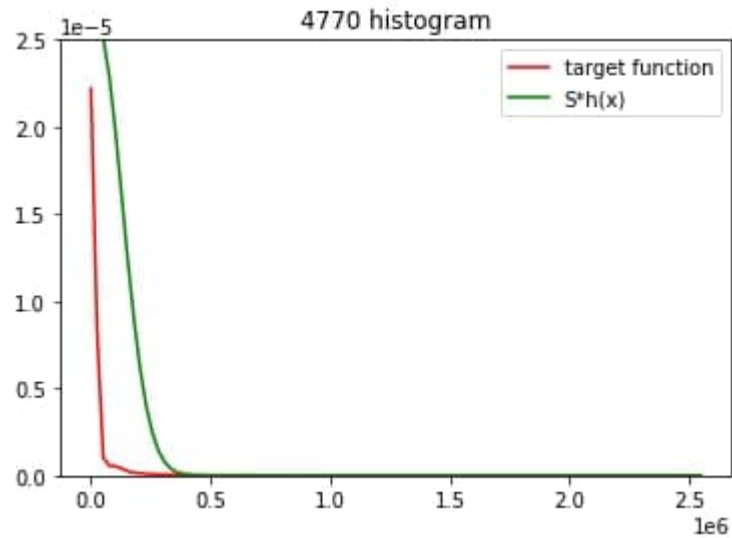


Figure 4: Accept-reject sampling algorithm (4770)

3. Estimation of relations between predictors and chosen target variables.

On the third step we calculate the correlation between our variables for next choosing of them to build the Bayesian network. The results shown on Figures 5, 6, 7.

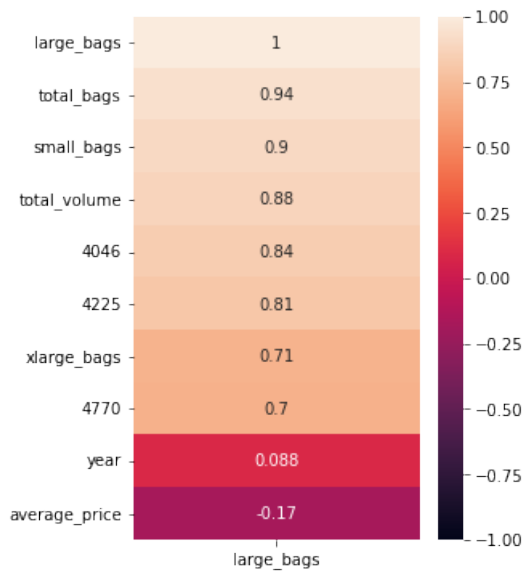


Figure 5: “large_bags” correlations

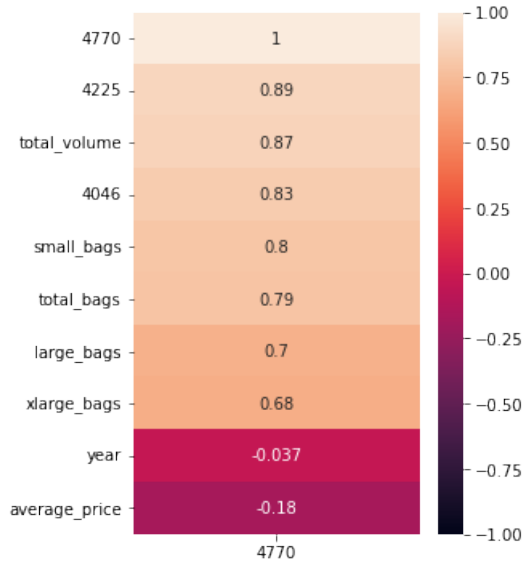


Figure 6: “4770” correlations

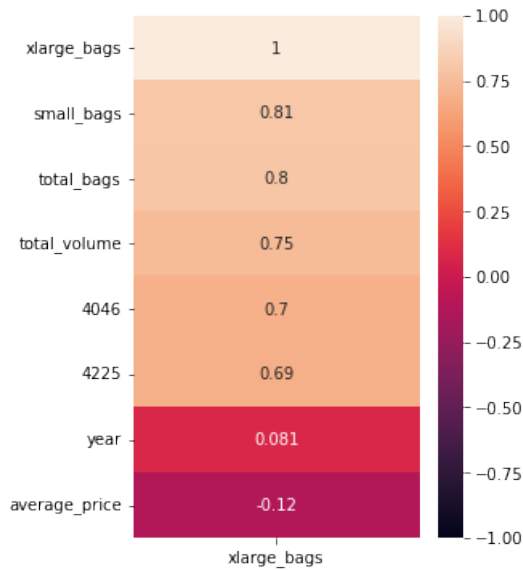


Figure 7: “xlarge_bags” correlations

4. Bayesian network.

In that section we going to build, firstly, manual Bayesian network and than 2 algorithms for structural learning.

For manual Bayesian network the variables were chosen based on correlations that were calculated in Section 3. For target variable “large.bags” predictors like “total_volume”, “small_bags”, and “total_bags” were chosen. For target variable “4770” predictors like “4225”, “total_volume”, “4046” were chosen. For target variable “xlarge.bags” predictors like “total_volume”, “small_bags”, and “total_bags” were chosen. The approach for choosing was simple: first variables with biggest correlation were chosen.

After that the manual Bayesian network was built. On Figure 8 we can see the result.

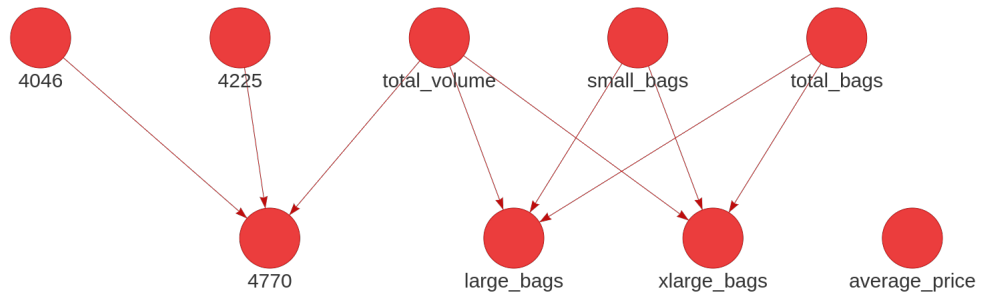


Figure 8: Manual Bayesian network

Then Bayesian network for the same set of variables based on algorithms for structural learning was built. On Figure 9 Bayesian network based on Hill-Climbing with Mutual Information function is shown:

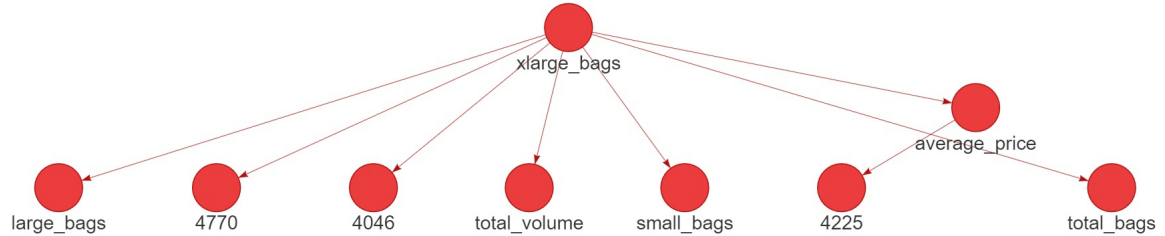


Figure 9: Bayesian network based on Hill-Climbing with Mutual Information function

The last one is Evolutionary algorithm with Mutual Information function. On Figure 10 the result is presented:

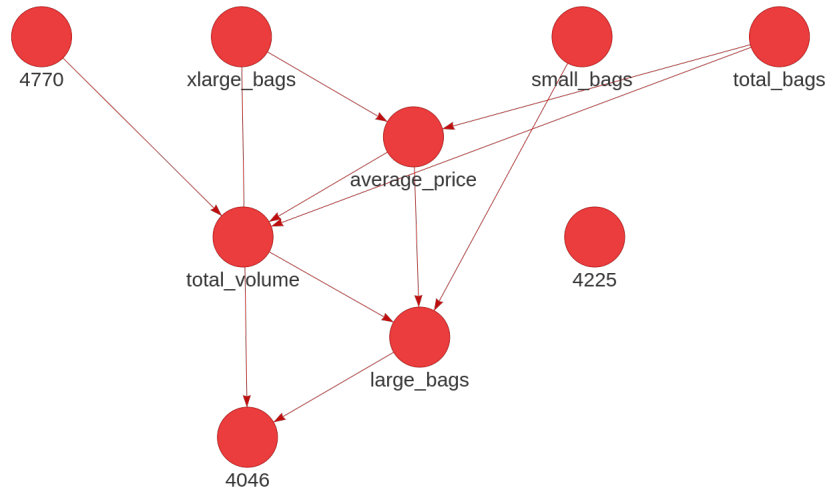


Figure 10: Evolutionary algorithm with Mutual Information function

5. Quality analysis.

In last section we want analyze the quality of our networks. Via methods from BAMT synthetic data was generated and drawn on real data. On Figures 11, 12, 13 results are presented.

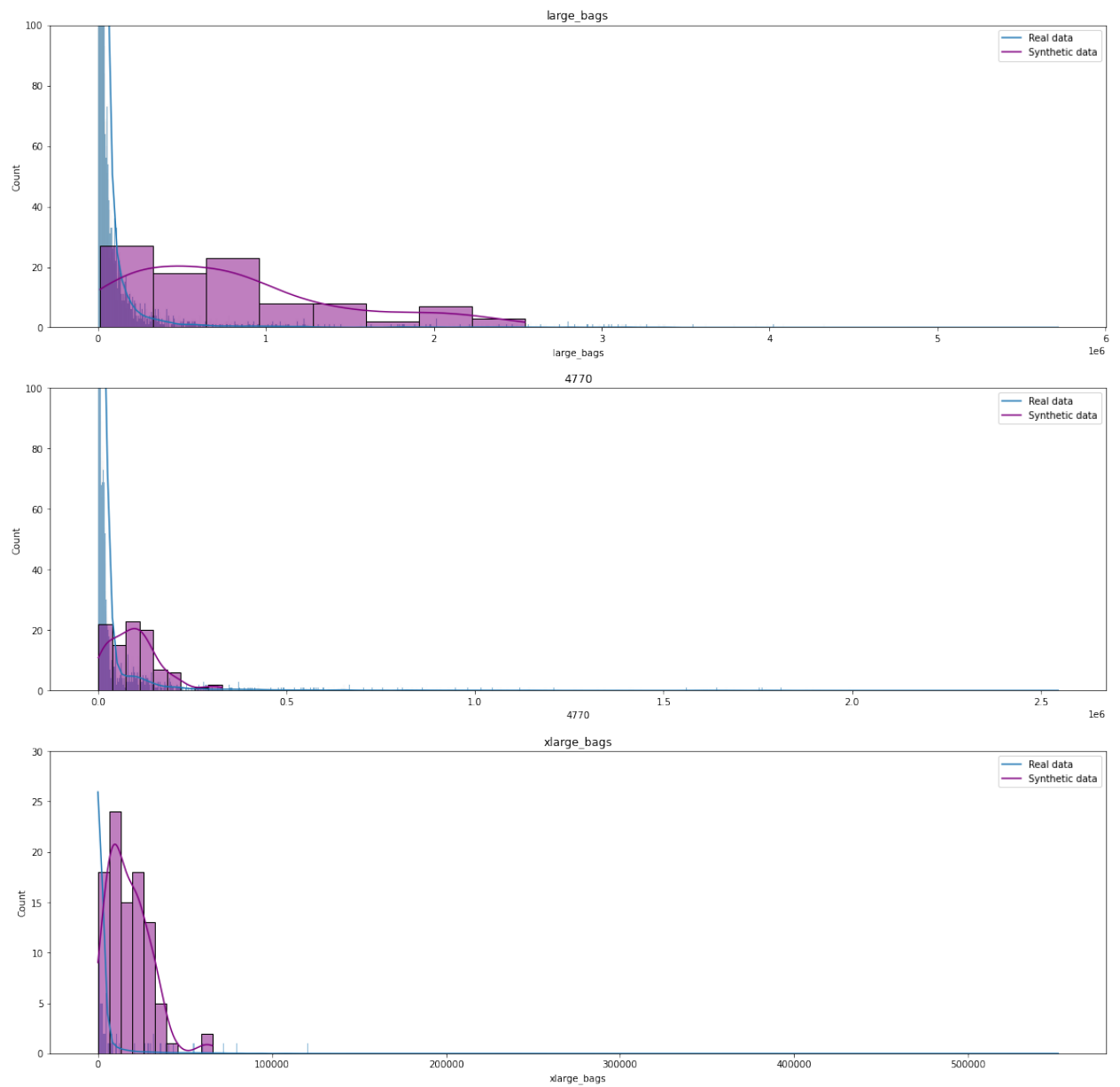


Figure 11: Mutual

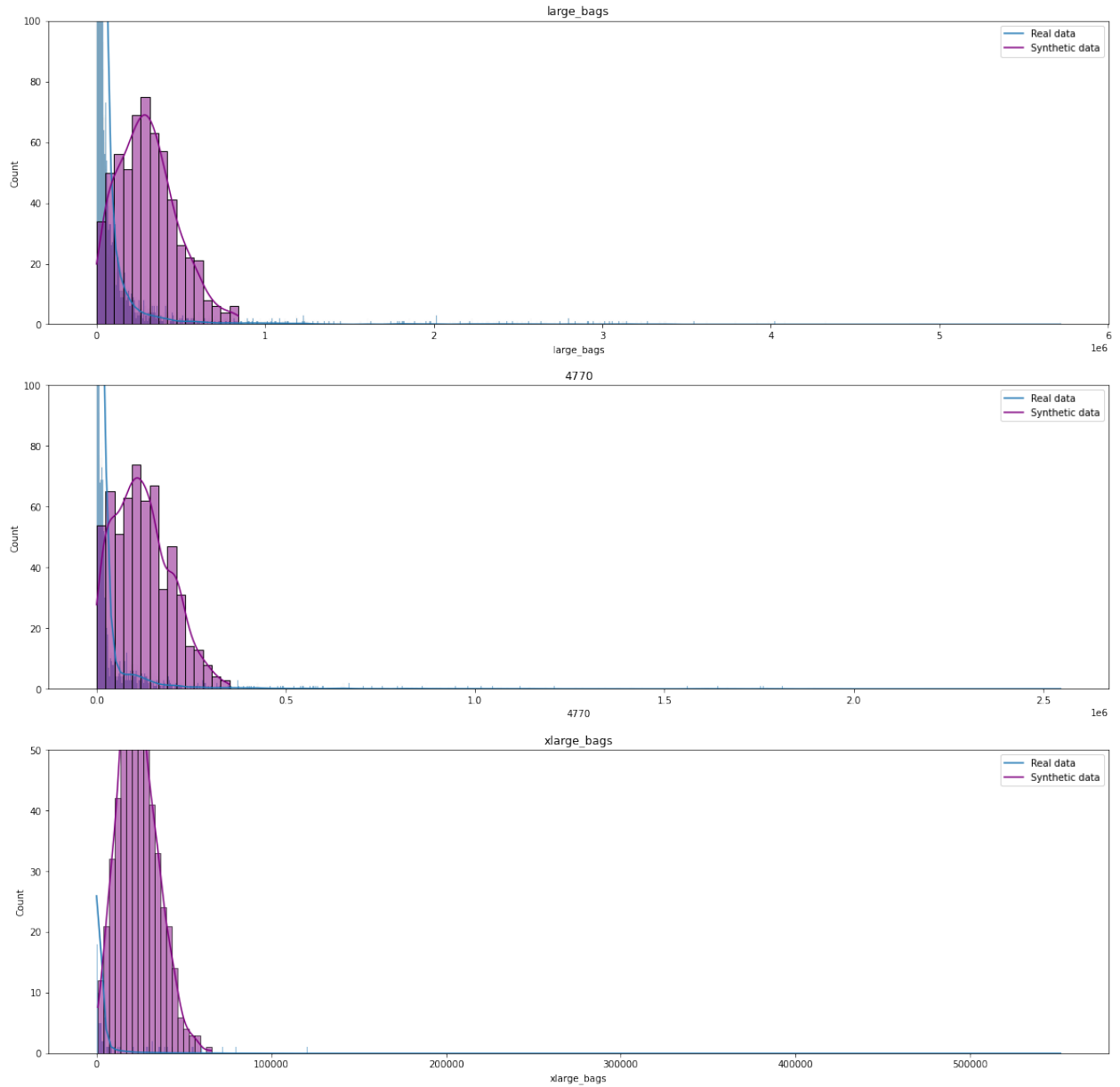


Figure 12: Hill-Climbing with Mutual Information function

The results are not perfect (perhaps via the chosen dataset). Generated data does not match with real. Also we can compare RMSE values for our networks. As we can see, the best result was provided by manual Bayesian network, since RMSE values for variables are the smallest. So the best model in our case is manual Bayesian network.

	Manual	Hill-Climbing	Evolutionary
large_bags	11453	234066	122737
xlarge_bags	12253	22461	25233
4770	66361	98871	134250

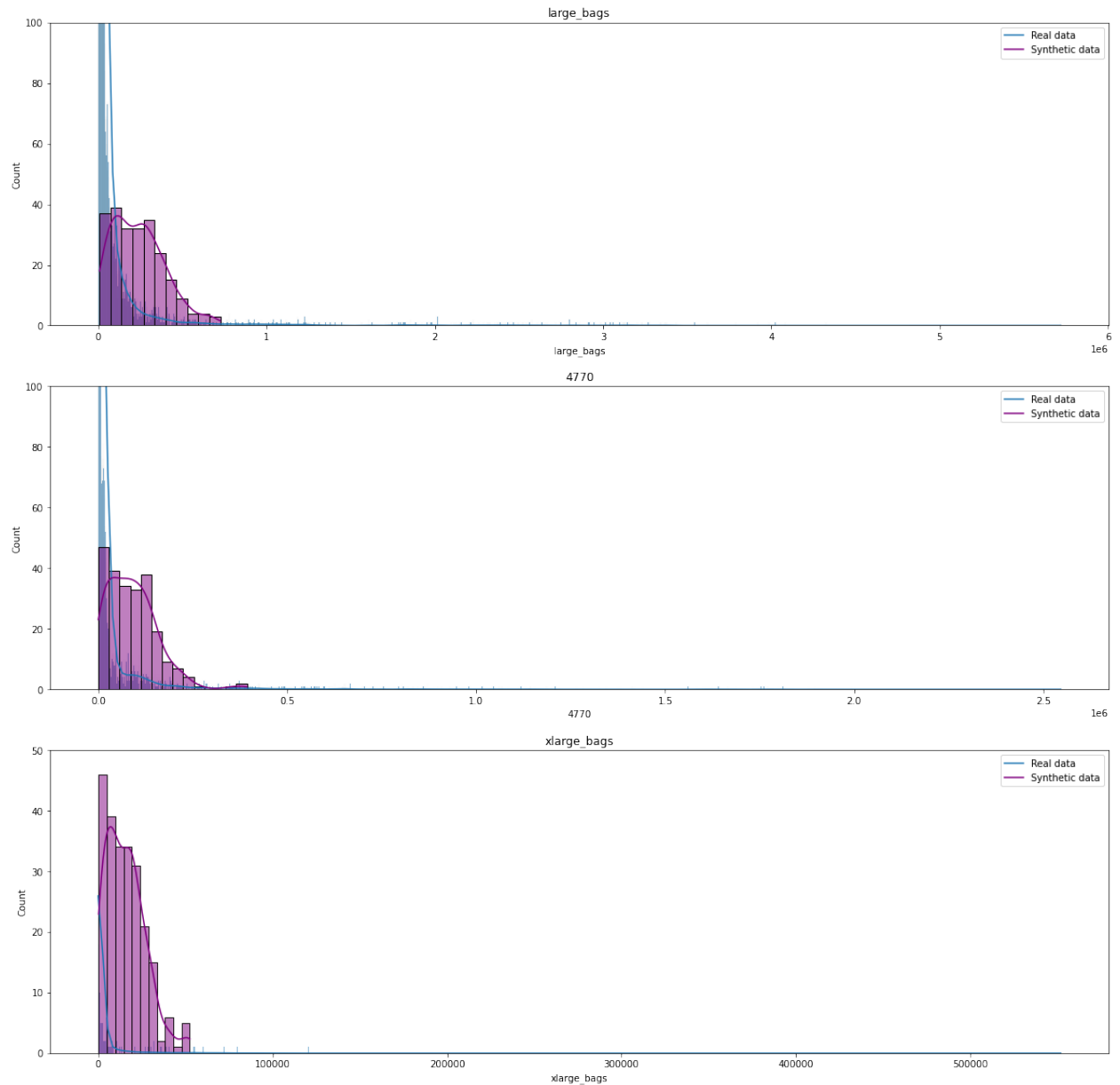


Figure 13: Evolutionary algorithm

Conclusions

In this laboratory work we used several algorithms for sampling data. Then built Bayesian networks using different algorithms and provide quality analysis for these algorithms. The best model was chosen based on minimum of RMSE.

Source code

Source code is located on GitHub: <https://github.com/DmitryPogrebnoy/multivariate-data-analysis>