



Department of Economics and Management

DEM Working Paper Series

**Modelling Probability of Default of Russian
Banks and Companies Using Copula Models**

Ilya Khankov

(National Research University Higher School of Economics)

Henry Penikas

(National Research University Higher School of Economics)

113 (12-15)

Via San Felice, 5

I-27100 Pavia

<http://epmq.unipv.eu/site/home.html>

December 2015

Ilya Khankov, Henry Penikas.

Modelling Probability of Default of Russian Banks and Companies Using Copula Models

Abstract.

Research is devoted to examination of the classifier, based on copula discriminant analysis (CODA). Performance of the classification of this algorithm was assessed.

On samples, modelled with some typical features of corporate default data, sensitivity of the classifier was tested, to sample size, to default rate and to different patterns of variables' interdependence.

Alternative copula families' selection method is proposed based on certain performance metric optimization. Difference in classification performance of different algorithms are investigated.

On real data of Russian corporate defaults, CODA classifier was built. It was supported by single factor analysis, based on discriminant analysis too. Final model demonstrates better classification performance than Linear Discriminant Analysis and Random Forest algorithm, and is comparable to Quadratic Discriminant Analysis.

Another experiment was set on data of Russian banks. Single factor analysis was assessed via standard procedure. CODA performance appeared to be lower than of Random Forest here, it was similar to QDA.

Contents

Introduction.....	3
1. Literature Review.....	4
1.1. Discriminant analysis.....	4
1.2. Facts from copula theory, sufficient for understanding CODA.....	6
1.3. Copula based discriminant analysis.....	8
2. Assessment of CODA classifier performance on model samples.....	11
2.1. Data considerations.....	11
2.2. General assessment.....	14
2.3. Different methods for copula family selection.....	18
2.4. Testing performance on small samples.....	20
3. Russian corporate credit defaults modelling.....	20
3.1. Data consideration.....	21
3.2. Single factor analysis.....	22
3.3. Model specification and performance assessment.....	24
4. Russian banks license withdrawal modelling.....	26
4.1. Data consideration.....	26
4.2. Single factor analysis.....	27
4.3. Model specification and performance assessment.....	27
Conclusion.....	28
References.....	30
Annex 1. Scatterplots of modelled data for bivariate case.....	31
Annex 2. Results of testing the sensitivity to sample size on bivariate data (General Assesment).....	33
Annex 3. Means and standard deviations of performance metrics in General Assessment.....	38
Annex 4. Performance characteristics of copula family selection method.....	40
Annex 5. Averaged on 4 distributions CODA performance characteristics on small samples.....	42
Annex 6. Results of single factor analysis on banks sample.....	44

Introduction.

In recent years, discriminant analysis appears to be not very popular statistical tool in credit risk modelling, especially in default predictions. However, there are some important improvements were reached in this method, which may increase the interest to discriminant analysis.

One of this improvements, application of copula theory to Bayesian classification, is a main subject of this paper. Although CODA (COpula based Discriminant Analysis) was first time introduced 9 years ago, there seems to be a lack of researches on assessment of applicability of this method to problem of credit default prediction. Given research is intended to give answers to the question: can CODA classifier make accurate enough predictions to use it in practice? There are some specific features of data used for credit risk modelling, which can potentially affect prediction quality of classifier. As an example, typical retail credit application dataset consists of limited number of factors, many of them are discretely distributed and number of observations is sufficient. In typical corporate dataset, however, number of observations can be limited to few hundreds or even less, but such dataset can include a large amount of continuously distributed factors. Moreover, some of corporate datasets can be considered as Low Default Portfolios, which means that share of defaults can be immaterial, or even equal to zero. This paper tries to give an insight into these aspects of modelling with CODA classifier.

This paper consists of 3 main chapters. In first, brief theoretical basis on discriminant analysis and copulas is provided, in volume required for understanding of copula based discriminant analysis (CODA). Then CODA classifier itself is introduced in the same chapter, and some restrictions are described. In second chapter CODA classifier behavior is assessed on modelled data, from simplest case to more complicated. Different aspects are considered,

sufficient for credit risk modelling. Third chapter includes an experimental assessment of CODA on real data. Model is built on dataset of Russian corporates, and its' defaults on bonds are modelled.

1. Literature Review.

Following chapter consists of three parts. In first part credit risk modelling problem is discussed, and main statistical models are described briefly. Second part is intended to introduce one to basics of discriminant analysis. Finally, third part of this chapter describes some facts of copula theory that is required for understanding of how copula based discriminant analysis is working.

1.1. Discriminant analysis.

Discriminant analysis is one of the classification methods, based on principle of supervised learning, i.e. it requires a training data sample with observations, assigned to classes and with observed characteristics. Based on the information from this sample, such classifier can assign any observation from the same general population to distinct class.

Discriminant analysis is a special case of Bayesian classifier, which is a classifier, based on Bayesian rule. The main idea of Bayesian classifier is to assign an observation to a class with biggest posterior probability. Posterior probability for a given class i can be calculated as following:

$$g_i(x) = P(\omega_i | x) = \frac{f(x | \omega_i) P(\omega_i)}{\sum_{j=1}^c f(x | \omega_j) P(\omega_j)}, \quad (1)$$

or it can be transformed by taking a logarithm (since monotonically increasing function not affect the classification):

$$g_i(x) = \ln f(x | \omega_i) + \ln P(\omega_i), \quad (2)$$

where $P(\cdot)$ is a cumulative distribution function, $f(\cdot)$ is a density distribution function, x — observation, which is defined as a vector of characteristics, and ω_i is a class with number i .

Two most widely used special cases of discriminant analysis are linear and quadratic discriminant analyses, which differ in assumptions. Both linear and quadratic discriminant analyses are based on the assumption of normal distribution of $f(x | \omega_i)$, which gives for quadratic discriminant analysis

$$f(x | \omega_i) = \frac{1}{|\Sigma_i|} \exp\left(-\frac{1}{2}(x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)\right), \quad (3)$$

where μ_i and Σ_i are mean and covariance matrix of class i . Linear discriminant analysis, however, is based on one additional assumption – homoscedasticity between classes.

One of the first researches, where discriminant analysis was applied to the problem of default modelling, was the article of Edward Altman [Altman, 1968]. In this article bankruptcy of large American corporates was predicted using discriminant analysis. As author writes in a conclusion: “The discriminant-ratio model proved to be extremely accurate in predicting bankruptcy...”

However, both linear and quadratic discriminant are usually underperform a linear regression, which became some kind of industrial standard in credit default modelling. According to results of recent assessment, performed by [Lessmann et. al., 2013], linear regression performs better than discriminant analysis on all tested samples (which are 7 retail credit application samples).

Although discriminant analysis is not as widely used in logistic regression, it has a potential of modifying, which is not available to latter. In theory, discriminant

analysis has no restrictions on the distribution of model factors, therefore it is more flexible and supposed to perform better than logistic regression, when distribution is chosen correctly. The problem is to model such multivariate distribution. This is the area, where copula theory can help.

1.2. Facts from copula theory, sufficient for understanding CODA.

Copula theory is one powerful instrument of multivariate dependence modelling. Copulas have applications in different areas, where it may be needed to calculate multivariate distribution of two or more random variables.

Formal definition of copula can be done as in [Penikas, 2010]. To do so, let define subcopulas first.

Definition. Let subcopula $C(X_1, X_2)$ as two-dimensional function of X_1 and X_2 , defined on $A \times B, A \in [0;1], B \in [0;1]$, with range of $[0;1]$ and satisfying following conditions:

1. If $\exists i: X_i = 0$, then $C(X_1, X_2) = 0$.
2. If $\forall j \neq i: X_j = 1$, then $C(X_1, X_2) = 1$.
3. If $x_j(1) \leq x_j(2)$, then

$$C(x_1(2), x_2(2)) + C(x_1(1), x_2(1)) - C(x_1(2), x_2(1)) - C(x_1(1), x_2(2)) \geq 0.$$

Definition. Let copula be the subcopula in case of $A = [0;1], B = [0;1]$.

Copula theory development starts from 1959 with the article of Abe Sklar [Sklar, 1959]. In this article Sklar's theorem – the foundation of all the copula theory – was formulated and proved.

Theorem [Sklar, 1959], bivariate case. Let H be the joint distribution function of two random variables x and y , having marginal distribution functions F and G respectively. Then there exists a function C :

$$\exists C : H(x, y) = C[F(x); G(y)] \quad \forall x, y \in (-\infty, +\infty). \quad (4)$$

If F and G are continuous, then copula C is unique.

Although the theorem gives an approach to model multivariate distributions, it should be mentioned that it imposes a restriction to the type of data, which can be modelled via copulas. It is important for the credit default modelling, that discretely distributed factors cannot be included into models, based on copulas. This fact can be especially important for retail models, where majority of application characteristic can be discrete. However, financial characteristics of corporates are mostly distributed continuously, and then corporate models can be the object of copula based discriminant analysis, which will be introduced below.

There exist several classes of copula families: elliptical, Archimedean, extremal and other. In this research it was chosen to consider only Archimedean copulas for their simplicity and ability to model different kinds of dependence (one can assess its variety in Annex 1). Copula families, used in this research are following:

Copula Family	Definition in bivariate case
Clayton copula	$(\max\{u^{-\theta} + v^{-\theta} - 1; 0\})^{-1/\theta}$
Gumbel copula	$\exp\left(-\left((-\log(u))^\theta + (-\log(v))^\theta\right)^{1/\theta}\right)$
Frank copula	$-\frac{1}{\theta} \log\left(1 + \frac{(\exp(-\theta u) - 1)(\exp(-\theta v) - 1)}{(\exp(-\theta) - 1)}\right)$

Joe copula	$1 - \left((1-u)^\theta + (1-v)^\theta - (1-u)^\theta (1-v)^\theta \right)^{1/\theta}$
------------	---

Table 1. Copula families, assessed in the research.

1.3. Copula based discriminant analysis.

The idea to use copula function in Bayesian classification first time appeared in the research of Saket Sathe in 2006 [Sathe, 2006]. In this article it is proposed to modify the formula 2 as following:

$$g_i(X) = \ln \left\{ c^i \left(F_1(x_1), \dots, F_d(x_d) \right) \right\} + \ln \left\{ \prod_{k=1}^d f_k(x_k, \theta_k^i) \right\} + \ln P(w_i). \quad (5)$$

The function $c^i(F_1(x_1), \dots, F_d(x_d))$ here is the copula density, which can be defined as:

$$c(x_1, \dots, x_d) = \frac{\partial^d C(F_1(x_1), F_2(x_2), \dots, F_d(x_d))}{\partial x_1, \dots, \partial x_d}. \quad (6)$$

In 2013 in the Journal of Machine Learning Research was published an article [Han, Zhao, Liu, 2013], proposing a different classifier with the same name – Copula Discriminant Analysis (CODA). It provided strong theoretical basis for the generalization of normal-based linear discriminant analysis to a larger class of Gaussian Copula discriminant analysis, which can be used on margins with so-called nonparanormal¹ distribution. Numerical experiments hold by authors indicated that proposed classifier is more efficient in terms of percentage of correctly classified observations, than linear discriminant analysis and sparse logistic regression.

¹ Which can be thought as distributions, able to be transformed to normal by strictly increasing transformations.

In the same year 2013, Eva Scheungrab in her master thesis [Scheungrab, 2013] provided some extensions into the model proposed by Sathe. She proposed to use kernel density estimations for estimation of margins, and vine structures to estimation of copula. However, results, obtained in real data assessments, are rather unsatisfying, by her words.

In aforementioned article, Sathe considers two methods of estimation the copula parameter. First, Exact Maximum Likelihood (EML) is the direct maximization of likelihood function by all marginal density parameters and copula parameter simultaneously. Second, so-called Canonical Maximum Likelihood (CML), is performed in three steps. At the first step, marginal distributions parameters are estimated using maximum likelihood approach. Then, using empirical marginal transformation, CDF (cumulative distribution functions) are obtained. At the third step, copula parameter is estimated on CDFs from previous step, using maximum likelihood estimation. According to Sathe, CML has advantage of robustness along with less computational costs compare to EML.

Classifier, used in this research is based on canonical maximum likelihood estimation of copula parameter. The algorithm of classification, used in base model (which is called LL-CODA here and below), can be described as following (steps 1-4 are performed on the training sample, others - on the test sample):

1. Marginal distributions of all explanatory variables are estimated using maximum likelihood estimation (parameters of the pre-defined distributions are estimated).
2. CDFs (cumulative distribution functions) of margins are calculated using parameter estimation from previous step.
3. CDFs from previous step are used to estimate parameters of Clayton, Gumbel, Frank and Joe copulas. Maximum likelihood estimation is used.

4. Estimation with maximum likelihood is selected from the list of estimations from previous step.
5. Logarithms of prior probabilities of the classes are calculated as logarithms of frequencies of the classes in the training sample.
6. On the test sample CDFs are calculated again using distribution parameters estimations from step 2.
7. Log-densities of marginal distributions are calculated as well as CDFs.
8. Copula log-density estimations are calculated on the CDFs from step 4 and copula estimations from step 3.
9. The value of the classifier for each class is calculated as sum of prior log-probabilities, marginal log-densities and copula log-densities for this class.
10. Each observation is assigned to the class with the highest value of the classifier for this observation.

Classifier, built and assessed in this research, has several restrictions:

1. In base model all marginal distributions are assumed to be Gaussian. Possibility of additional distributions estimation is expected to increase the prediction quality of CODA classifier.
2. Only 4 copula families are assessed now. Allowing more copula families to be estimated can enhance CODA performance either.
3. In case of higher dimensions, hierarchical- and vine-copulas should be considered as more flexible.
4. As it was mentioned above, discretely distributed variables were not considered, since Sklar's theorem cannot guarantee the uniqueness of copula with such variables.

All these restrictions are planned to be considered in further researches, along with more detailed testing of performance and compare with different classifiers.

2. Assessment of CODA classifier performance on model samples.

As the initial step of CODA classifier performance assessment, simple case of bivariate data will be considered, with further extension to multivariate. Although bivariate case seems unrealistic for credit risk modelling, it can be a good point to start in order to understand model behavior. In this case bivariate copulas can be used in CODA classifier, since there is no need of sophisticated copula models. We will consider classifier based on 4 Archimedean copula families: Clayton, Gumbel, Frank and Joe, and compare them to LDA and QDA classifiers.

Main object of this experiment is to assess CODA performance in different circumstances which may be the case in credit risk modeling. In order to do this, it is important to understand the classifier behavior, specifically:

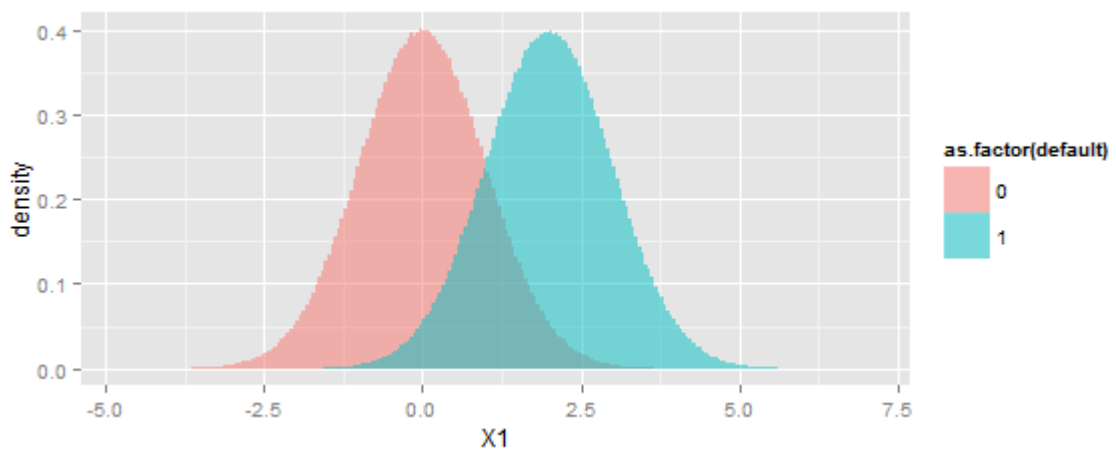
- *Sensitivity to the sample size.* It is important to understand, whether is CODA performs well on small samples (which is usually the case in evaluation of corporate borrowers).
- *Sensitivity to the default rate, i.e. to the ratio of the volumes of the classes in sample.* It should be assessed, whether can CODA classifier be appropriate (and can it perform better than LDA or QDA) on low-default portfolios.
- *Sensitivity to differences between classes.* It is important to understand, whether is CODA classifier capable to correctly classify observations with:
 - different distribution families among the classes;
 - same distribution families, but different copula parameters.

2.1. Data considerations.

On first stage of the experiment it were generated 12 samples of size 2 million observations each. These samples were generated same distributions of margins,

but with different dependence function (4 families of copulas for defaulted and non-defaulted classes, plus 4 samples for multivariate normal distribution).

Each sample consists of 2 factors and default indicator. It is expected that distribution family of margins will not affect results of comparison between CODA and LDA/QDA, because CODA algorithm, realized in this research, uses same method of margins estimation as LDA/QDA. It was considered to use Gaussian distribution of margins. For each variable means for defaulted and non-defaulted classes are different, while variances are equal. Parameters of marginal distributions were set this way in order to model the standard situation in credit risk practice, where distribution of factor for defaulted class is similar to one for non-defaulted class, but shifted for a certain amount. This amount of shift was set equal two standard deviations in the research. Histograms of margins, built on initial sample of 2 million observations, are given on plots below.



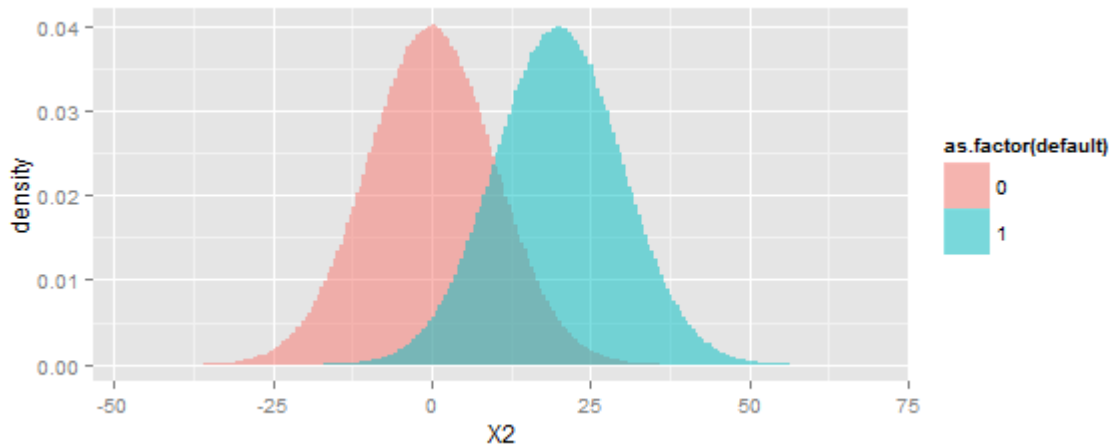


Diagram 1. Histograms of marginal distributions of factors.

As it was mentioned above, marginal distributions were set the same in all testing cases on this step. Diversity of samples was achieved by using different copula functions. Four Archimedean copula families were used: Clayton, Gumbel, Frank and Joe. These families were chosen in order to assess performance on the data with different dependence patterns (to get the picture of these patterns, please see Annex 1). Additionally to these copulas, multivariate normal distribution was assessed as the best suited dataset for LDA/QDA classifiers. Two versions of multivariate normal data were used: with correlation between variables on the level 0.8 and without correlation. Scatterplots of all samples used for testing are provided in the Annex 1.

All samples for learning and testing classifiers were generated from initial 12 samples by random subsampling and combination of required samples for defaulted / non-defaulted class. Samples, generated this way, were randomly split to training and test samples of equal size.

2.2. General assessment.

First part of the experiment performance of classifiers was assessed on 240 samples of different distributions, different sizes and different default rates, considering both extremal and typical cases. Parameters were varied as following:

- *Distribution*: all combinations of Clayton, Gumbel, Frank and Joe families for defaulted and non-defaulted classes;
- *Default Rate*: 1%, 10% and 30%;
- *Sample size*: 100, 500, 1000, 5000, and 10000 observations.

This part was designed this way to obtain rough, but broad picture of CODA performance in different circumstances. General purpose of this part of the experiment is to roughly mark boundaries of use of this classifier (i.e. sample size, default rate and differences between classes), which will be estimated more precisely in following parts of the experiment.

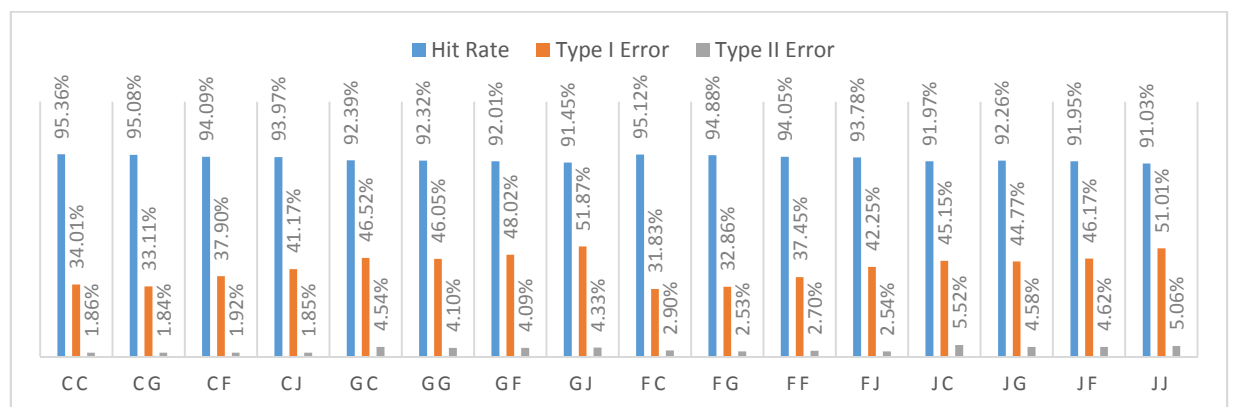
In order to assess classifiers performance, three metrics were used: Hit Rate (share of correctly classified observations in the sample), Type I Error (share of misclassified defaulted observations) and Type II Error (share of misclassified non-defaulted observations). Full dynamics of these metrics (relative to sample size) on samples, generated by different copula families, and with different default rate are provided in the Annex 2.

Sensitivity to distribution family can be assessed, using averaged by copula combination performance of the classifier. On diagram 2 following averaged metrics are provided: Hit Rate, Type I and Type II Errors. As one can see from the diagram 2.a, performance of CODA depends on copula families, used for sample generation, despite the fact, that all families, used in testing, were available for classifier to select. The reason for unstable results is a difference, appearing in combination of two classes. As one can see in Appendix 1, combinations with high performance (like Clayton-Clayton or Frank-Gumbel) are much easier to recognize

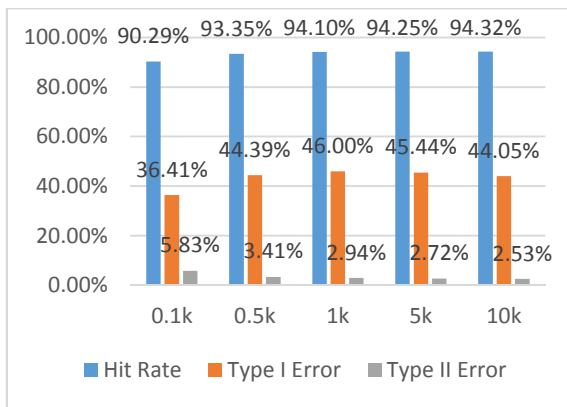
even visually, then combinations with low performance (like Gumbel-Frank or Joe-Joe).

Difference in performance is highlighted best in the level of Type I Error on different distributions, which varies from 31.83% on Frank-Clayton sample to 51.87% on Gumbel-Joe sample. Type II Error in general behave the same way, although it reaches minimum on Clayton-Clayton sample. Hit Rate varies from 91.03% on Joe-Joe sample to 95.3% on Clayton-Clayton sample, so there is a similarity to behavior of other metrics. Pearson correlation coefficient between Hit Rate and Type I Error is -96.9%, between Hit Rate and Type II Error is -89.5%, and between two errors is 79%.

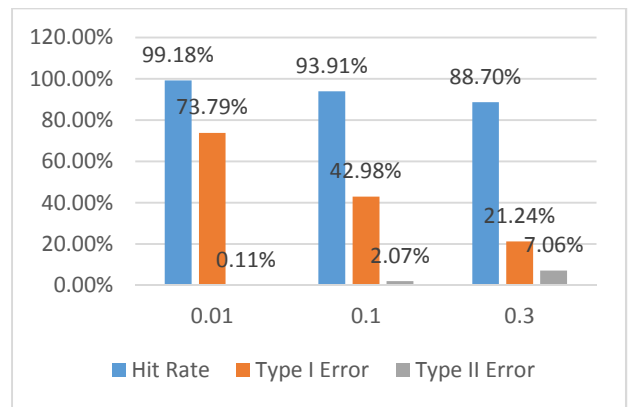
Main result of this analysis is that samples, better in terms of Hit Rate of classifier built on it, will generally be better in terms of error rates too. Also, based on these results, we can reduce number of assessed copula families' combinations. In next assessments, only 4 combinations will be used: Clayton-Clayton as the best case, Joe-Joe as the worst, and Gumbel-Gumbel and Frank-Frank as middling cases.



a) by copulas used for sample generation.



b) by sample size.



c) by default rate

Diagram 2. Averaged performance statistics of LL-CODA.

As one can see from diagram 2.b, *sensitivity on sample size* differs on small and large samples. Once sample reaches 500 observations, there is no major increase in performance metrics with further increasing in sample size. Although, performance dynamics on the interval between 100 and 500 observations should be assessed more precisely.

Analysis of different default rates, conversely, indicates that there is a material dependence of all three assessed metrics on the default rate. Classifier shows best Hit Rate on samples with default rate of 1%, however, it did not recognize almost $\frac{3}{4}$ of all defaults. This performance cannot be considered as satisfactory. Increase in default rate cause decrease in Hit Rate and increase in Type II Error, but considerable decrease in Type I Error. In other words, classifier becomes to less precisely estimates non-defaulters, and more precisely – defaulters.

It may seem from diagrams above, that CODA classifier is not conservative enough to use it “as is”. Indeed, Type I Error on a level of 40% is not satisfactory in case of credit default prediction. However, as one can see from Annex 2.B, these levels are not specific classifier characteristic, but it is characteristic of data. Both Linear and Quadratic Discriminant Analyses shows comparable results in terms of all calculated metrics.

In order to understand classifier performance in more details, it is useful to assess separately its parts. Two special steps of CODA classifier fitting are selection of copula family for each class and fitting copula parameter of selected families.

In base model (which is called LL-CODA in this paper), copula family was selected using likelihood maximization criterion. It was expected, that this method will correctly select copula families, if sample size is big enough. It was also expected, that copula families, selected this way, in general will provide best predictions in terms of Type I and Type II Errors. However, maximum likelihood does not guarantee this. Analysis of estimations, obtained on this step, indicates that these expectations are partly correct.

As one can see from Annex 4.A, copula families for non-defaulted class were predicted correctly in 100% of cases. Share of correctly predicted families for default class depends on copulas combination, sample size and default rate. However, average result is 88%, and on sample of size more than 1'000 observations it reaches 100%. Therefore, this part of classifier performance can be assessed as satisfactory.

The other part is precision of copula parameter estimation. In Annex 4.B averaged t-scores are provided. It can be seemed, that despite the expectations, estimates of copula parameters are become less precisely with sample size and default rate increase. However, this is not the case. As it is demonstrated on graphs d and e of this annex, average deviations of estimates from true values of copula parameters are decreasing with sample size or default rate increase. However, deviations are decreasing slower than standard errors of parameter estimations, which cause increase of t-scores. Therefore, there is no statistical evidence that sample size increase cause improvement of copula parameter estimates precision.

Summarizing aforementioned, main results of general assessment are following:

- Performance of CODA classifier is relatively stable on all tested distributions, however, Type I Error is volatile: it varies among distributions by 20%, while Hit Rate – only by 4%, and Type II Error – by 3%. All metrics are highly correlated: distributions, assessed as good by one metric, will be assessed as good by other two.
- There is no major increase in performance with increase in sample size, since it reaches 500 observations. Most probable reason for this result is the behavior of copula parameter estimates: precision of copula parameters estimates increases very slow.
- There is a material impact of default rate in CODA performance. There is a trade-off between Type I- and Type II Errors. Once default rate increases, Type I Error begins to decrease significantly, but Type II Error grows at the same time. Hit Rate behave the same way.

2.3. Different methods for copula family selection.

As it was mentioned above, in base CODA model in this research, likelihood maximization criterion was used. Although this method showed satisfactory share of correctly selected copula families, it turns to be underperforming compare to other copula families. It was indicated, that in some cases incorrect copula family generate a better prediction, than correct one.

As one can see from the Annex 4.C, share of cases, where LL-CODA showed best results among all available combinations, depends materially from particular distribution. Dependence on sample size and default rate is not that important, it mostly ceases as number of defaults reached required volume (which is roughly equals 100 observations, however will be assessed more precisely). Higher values on small samples are caused by overall poor performance of all copula combinations on these samples with little difference among it.

We propose other method of copula family selections, which is based on performance metric optimization. Algorithm of classifier is similar to one for LL-CODA, except Step 4. Instead of selection for each class copula family with maximum likelihood, in this version of CODA classifier particular performance metric is calculated on learning sample for each available copula families' combination, and combination with best value of metric is selected.

In tests we used three metrics, mentioned above, as a target characteristic in process of family selection. Results of assessment are provided in the Diagram 3.

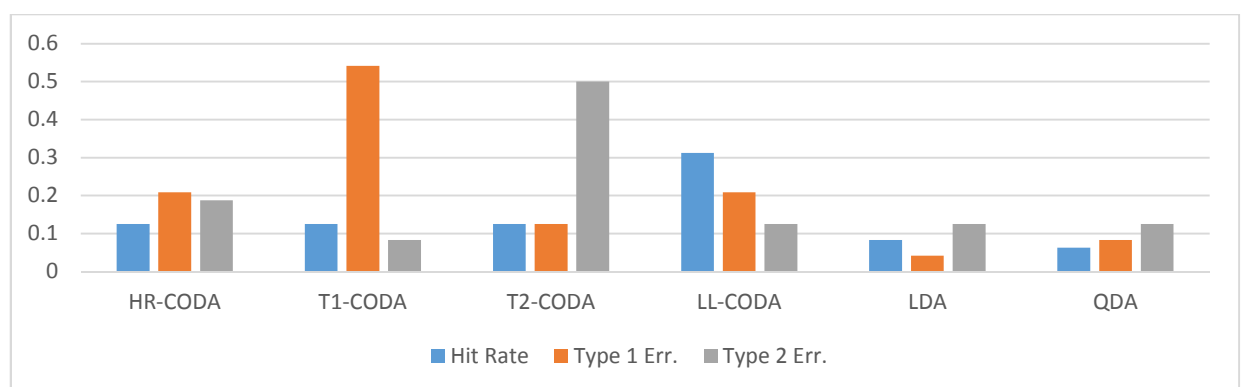


Diagram 3. Share of cases, where given classifiers showed best performance among all available copula families' combinations.

Interesting result of this assessment is the performance of T1-CODA, which is CODA, selecting copula families with minimal Type I Error. As one can see from the diagram, performance of all CODA classifiers is different, what means that different methods of family selection materially affects the performance. What is more important, is that T1-CODA allows to increase the performance of the classifier in terms of Type I Error, which is an object of interest in credit default modelling. Average Type I Error of T1-CODA is 34.6%, while for LL-CODA it is 39.6%, and for LDA it is 42.9%. In the same time, it provides comparable Hit Rate: 92.4% for T1-CODA, 92.7% for LL-CODA and 92.1% for LDA.

Summarizing these results, T1-CODA can be considered as good alternative to LL-CODA, if latter performs unsatisfactory in terms of Type I Error, or if one need conservative estimations.

2.4. Testing performance on small samples.

To assess more precisely performance of CODA classifier on small samples, it was decided to vary number of defaults, with fixed default rate equal 10%. This decision is caused by the fact that in CODA classifier defaulted and non-defaulted observations are assessed separately. Four copula family combinations, mentioned above, were assessed. Sample size was varied from 50 observations to 500 observations with 50 observations step. That gave assessment on samples with number of default from 5 to 50. Results are provided in Annex 5.

LL-CODA demonstrates stable high level of Hit Rate: median level is 92.65%, which is higher than other tested classifiers, and difference between minimum and maximum is shortest. Results in terms of Type II Error are comparable. Unfortunately, LL-CODA turns to be non-conservative: Type I Error is highly volatile and its' median level is higher than of any other classifier (53% for LL-CODA, 37% for T1-CODA, 47% for LDA and 44% for QDA).

T1-CODA demonstrates best result in terms of Type I Error, but worst in terms of Type II Error. Its' median Hit Rate is lowest (90.5%, against 92.5% for LDA, 91.2% for QDA and 92.6% for LL-CODA).

3. Russian corporate credit defaults modelling.

In this section we describe the results of the experimental credit default model building for Russian corporates. Data for this experiment were kindly provided by Maria Ermolova and Henry Penikas. These data were used in their recently published paper [Ermolova, Penikas, 2015], in which interdependence between

two credit risk components – PD and LGD² – is assessed. One of the steps in their assessment is PD model construction. Authors use for this purpose the statistical instrument of logistic regression. The experimental modelling exercise, described below, is intended to use CODA classifier in the process of model development and assess its performance. Although multidimensional cases and different marginal distributions were not carefully assessed in previous chapter on modelled data because of the volume of calculations required for this, preliminary results of its implementations were rather satisfying, and therefore it was decided to consider these aspects in experimental model construction in this section.

3.1. Data consideration.

Sample consists of 4050 observations on 4045 Russian corporates, issued bonds traded on the market. There are 214 defaulted observations in the sample which gives a default rate of 5.3%.

Default definition used here is in line with Basel II requirements, i.e. it considers all sufficient overdue by 90 days or more.

There are 51 variables in the samples, which were considered in this experiments, all representing financial characteristics of company's performance. Despite sample includes some discrete variables, the majority of factors are continuous, which is important, since CODA classifier, used in this research can't handle discretely distributed variables.

Data were divided on 2 samples of equal size as it was done in research by Ermolova and Penikas (out-of-sample split). One sample was used for model building (training sample) and other – for model's performance assessment (test sample).

² PD – probability of default. LGD – loss given default.

3.2. Single factor analysis.

Factors from the long list of 46 continuous variables were selected using single factor assessment of discriminatory power. For each variable discriminant analysis with only one variable was estimated. The algorithm is as following (steps 1–2 on training sample, other – on test sample):

1. All available distribution families are fitted on the data on the data of given variable separately for defaulted and non-defaulted observations. In this experiment normal and exponential distributions were assessed. Distinct family is selected using maximum likelihood criterion. Distribution parameters are saved.
2. Logarithms of prior probabilities of the classes are calculated as logarithms of frequencies of the classes in the training sample.
3. On test sample log-densities of marginal distributions are calculated using distribution family and parameters, estimated on step 1.
4. The value of the classifier for each class is calculated as sum of prior log-probabilities and marginal log-densities for this class.
5. Each observation is assigned to the class with the highest value of the classifier for this observation.

Given this, algorithm of variables selection is in line with the methodology of classification, used in this research. It is expected, that variables with highest differences between marginal distribution for defaulted and non-defaulted classes will be selected.

Applying this algorithm to 46 variables from the long list, following results were obtained. For 33 variables were not able to perform a classification, assigning all observations to one class (non-defaulted). For other variables results are following.

	HR	T1	T2	family for defaulted class	family for non-defaulted class
ar_turnover	94.3%	95.8%	0.3%	<i>Normal</i>	<i>Normal</i>

clca	60.3%	54.3%	38.8%	<i>Exponential</i>	<i>Exponential</i>
Cfta	59.0%	34.5%	41.4%	<i>Exponential</i>	<i>Exponential</i>
cap_ap	35.1%	5.0%	68.7%	<i>Normal</i>	<i>Normal</i>
salta	32.1%	17.6%	71.0%	<i>Normal</i>	<i>Exponential</i>
conc_equity	31.4%	9.5%	72.2%	<i>Normal</i>	<i>Normal</i>
cap_ta	31.1%	9.2%	72.6%	<i>Normal</i>	<i>Normal</i>
ebit_tot_debt	20.6%	1.7%	84.8%	<i>Normal</i>	<i>Normal</i>
ar_debt	16.7%	2.6%	88.2%	<i>Normal</i>	<i>Normal</i>
cash_lr_liab	16.6%	1.8%	89.1%	<i>Exponential</i>	<i>Exponential</i>
ebit_net_debt	11.8%	2.6%	93.4%	<i>Normal</i>	<i>Normal</i>
cash_gain	10.9%	5.6%	94.5%	<i>Normal</i>	<i>Normal</i>
cur_ass_liab	10.4%	0.9%	95.1%	<i>Normal</i>	<i>Normal</i>

Table 2. Results of single factor analysis.

For all variables with marginal distribution assessed as normal, outlier analysis was performed, considering as an outlier all observations outstanding from the mean for more than 5 standard deviations. For margins, considered as exponential, visual analysis was performed, however, no outliers was found in these margins. As a result of an analysis, 14 observations were considered as outliers and deleted from the sample. Results were materially changed for 3 variables from the list. Variables *cap_ap* and *cap_ta* loosed the discriminatory power, starting to assign all observations to non-defaulted class. Variable *clca* increased Type II Error, but decreased Type I Error.

Variables having Type II Error more than 60% Type I Error or 80% Type II Error were taken out of the consideration as not efficient. Other 6 variables have formed short list, given below.

	HR	T1	T2	family for defaulted class	family for non-defaulted class
clca	59.1%	53.8%	40.1%	<i>Exponential</i>	<i>Exponential</i>
cfta	59.0%	34.5%	41.4%	<i>Exponential</i>	<i>Exponential</i>
salta	32.1%	17.6%	71.0%	<i>Normal</i>	<i>Exponential</i>
conc_equity	31.4%	9.5%	72.2%	<i>Normal</i>	<i>Normal</i>

Table 3. Variables in the short list.

3.3. Model specification and performance assessment.

For variables of the short list all possible combinations were considered and models were assessed. Final specification was selected by expertly set criterion of maximization of linear combination of three performance metrics, controlled here:

$$\text{Hit Rate} - 1.25\text{Type I Error} - 0.5\text{Type II Error} \rightarrow \max \quad (7)$$

Weights of metrics in this combination represents its' subjective importance for default forecasting.

In final specification three variables were chosen: ratio of accounts payable to current assets (*clca*), ratio of cash flow to total assets (*cfta*) and ratio of sales to total assets (*salta*). This list sounds reasonable and is partly in line with the model, recently developed by CMASF (Center for Macroeconomic Analysis and Short-term Forecasting) to predict bankruptcy of Russian real sector companies [Mogilat, 2015]. This fact can be an evidence of soundness of the algorithm used for variable selection on the previous step.

In Chapter 2 several different algorithms for CODA were considered. In Table 4 one can find results for each of them, applied to the model specification, selected on previous step. It is interesting, that unlike on the modelled data in section 2.3 of this paper, "LL-best" algorithm (base algorithm, described in section 1.3) performs better in terms of Type I Error, than "T1-best" algorithm, selecting copula families' combination, which minimizes Type I Error on development sample. "HR-best" and "T2-best" algorithms, optimizing respective metrics on development sample, performs equally, and better in terms of both these metrics than "LL-best" and "T2-best" algorithms.

However, percentage of correctly classified may seem unsatisfactory, as it reaches only 36% provided by best algorithms. But the case itself appears to be extremely difficult not only for CODA classifier. Two classic classification algorithms – LDA and QDA – showed inaccurate results. LDA assigned all observations to non-

defaulted class, reaching Hit rate of 95.3% and Type I Error of 100%, which is completely unacceptable. QDA performed better, as it correctly recognized all defaulted observations and about 12% of non-defaulted, which gives 15.6% of the sample.

Following the recommendations of [Lessmann et. al., 2013], we compared these results to one, obtained with Random Forest classification algorithm. Not getting into details, this classifier builds a voting committee of randomly generated decision trees. Full description of the method can be found in [Breiman, 2001], where it was first time introduced by its' author, Leo Breiman. For this assessment, Random Forest as it implemented in R in package 'randomForest' was used. As one can see from the Table 5, this algorithm shows comprehensive results. With Type I Error fixed approximately on same level as CODA provides, it can reach higher level of Hit Rate. Therefore, there is no evidence that CODA classifier performs better than other modern classifiers, but its' performance can be assessed as comprehensive.

model	good_cop	bad_cop	good_par	good_sd	bad_par	bad_sd	hitrate	type1	type2
HR_best	joe	joe	1.105	0.000	1.433	0.009	<u>0.360</u>	0.295	<u>0.657</u>
T1_best	gumbel	clayton	1.166	0.000	0.443	0.010	0.352	0.284	0.666
T2_best	joe	joe	1.105	0.000	1.433	0.009	<u>0.360</u>	0.295	<u>0.657</u>
LL_best	clayton	frank	0.491	0.001	2.229	0.155	0.324	<u>0.263</u>	0.696
QDA	NA	NA	NA	NA	NA	NA	0.156	0.000	0.885
LDA	NA	NA	NA	NA	NA	NA	0.953	1.000	0.000
RF	NA	NA	NA	NA	NA	NA	0.438	0.274	0.576

Table 4. Performance metrics of final model specifications.

Although all considered classifications were limited to only three variables from the sample, obtained results can prove that there exist real credit default modelling cases, where CODA classifier can demonstrate comprehensive level of performance. These results provides a justification for further assessment and improvement of this classification algorithm, including extensions, described in Section 1.3.

Since CODA was considered in this research only as binary classifier, it was not possible to directly compare results obtained here with ones, obtained by Ermolova and Penikas in their aforementioned article.

4. Russian banks license withdrawal modelling.

Last section is dedicated to another real data experiment. Here we build a model of license withdrawal of Russian commercial banks. Defaults of commercial banks are expected to be an area of interest during my thesis preparation in several next years, and this section will begin the research.

Defaults of commercial banks have their features, many of which are not considered in this experiment. A simple example of this features is a default definition. Here, in this experiment, we build a model of license withdrawal forecasting, which is not equal to default, since license may be withdrawn for several reasons except default. Nevertheless, even the experiment built this way may be insightful in understanding of CODA performance on bank data.

4.1. Data consideration.

Sample of Russian banks was formed using data of Mobile agency, database “Banks and Finance”. One section was taken, as of January, 1 of 2013. Information of license withdrawals was collected manually from the official CBR announces. Observation period was set to one year, i.e. license withdrawals were tracked from January, 1 of 2013 to January, 1 of 2014.

Initial sample was built of 944 Russian banks with 30 license withdrawals which gives a default rate of 3.17%. Database provided 177 variables, but many of them contained high level of missing values. Variables with more than 199 missing

values (21%) were excluded from the following consideration. After this filtering, only 80 variables are left in the sample.

4.2. Single factor analysis.

Single factor analysis, assessed in section 3 had indirectly prove it's soundness by comparison with another model, built by experts of CMASF. However, this procedure may affect the result of classification algorithms comparison, since factors are selected using discriminant analysis. Therefore, in this experiment it was decided to use standard metric for single factor analysis, AUC (Area under the ROC Curve). This metric is unspecific to the form of interconnection between default flag and factor.

All financial factors were scaled by division to the value of current assets of banks. This transformation is commonly used in banks defaults forecasting (e.g. [Peresetsky, 2013], [Karminsky, Kostrov, 2012]) as it helps to avoid distortions in model coefficients caused by modelling defaults for large banks along with very small ones.

After the scaling, AUC metric was assessed on all 80 variables. Full list of variables with corresponding AUC values may be found in Annex 6. It should be noted, that unlike the previous experiments, on this dataset all factors are comprehensive, difference in AUC between the best and the worst factor is less than 22%.

4.3. Model specification and performance assessment.

Final specification was selected by considering factors with highest values of AUC (higher than 60%) and representing different aspects of bank's condition. Selected factors sound reasonable and considered as acceptable for this model. In order to reduce computational costs, which are relatively high for CODA algorithm, described and realized in this research, three-factor model was assessed again. This

number can't be considered as acceptable for the model, which forecasts will be used for decision making, this may be sufficient at this phase of research.

First factor is the ratio of balance sheet profit to current assets. This factor seems reasonable, as it is considered as rough measure of financial success of a bank. Second factor is a net position of a bank on the interbank market. Third factor is a normative ratio called CAR (capital adequacy ratio).

Results obtained during this experiment are worse than in previous one. While CODA demonstrates modest results, Random Forest predicts very accurately, providing higher Hit Rate and lower errors of both types simultaneously. Performance of CODA, however, is similar to QDA in terms of Hit Rate. All algorithms of CODA, considered in this research, demonstrated equal results in terms of quality metrics on this sample.

These results are proving that CODA as it was realized here, is very sensitive to data. More extensions should be realized before compare its results to comprehensive classifiers.

model	good_cop	bad_cop	good_par	good_sd	bad_par	bad_sd	hitrate	type1	type2
HR_best	clayton	joe	0.294	0.002	1.906	0.146	0.284	0.333	0.726
T1_best	clayton	joe	0.294	0.002	1.906	0.146	0.284	0.333	0.726
T2_best	clayton	joe	0.294	0.002	1.906	0.146	0.284	0.333	0.726
LL_best	joe	joe	1.701	0.004	1.906	0.146	0.284	0.333	0.726
QDA	NA	NA	NA	NA	NA	NA	0.272	0.250	0.741
LDA	NA	NA	NA	NA	NA	NA	0.975	1.000	0.000
RF	NA	NA	NA	NA	NA	NA	<u>0.745</u>	<u>0.167</u>	<u>0.257</u>

Table 5. Performance metrics of final model specifications for banks.

Conclusion.

Main result of this research is that CODA classifier, as it was generally proposed by Sathe in 2006, can be appropriate for credit risk modelling, and can demonstrate comprehensive performance both on modelled and real data.

Modelled data assessment indicated that CODA classifier is capable to provide stable predictions on samples of 500 observations or more. Performance on small

sample can be unstable, however it was indicated, that these results are better, in general, than ones obtained by LDA and QDA. Impact of default rate is material, and performance on Low Default Portfolio is rather unstable.

New method of copula family selection within CODA classifier was introduced, which is based on performance metrics optimization on development sample, rather than likelihood maximization. Although maximization of Hit Rate give approximately equal to LL-CODA results, minimization of Type I- or Type II Error helps in building of more conservative (in terms of defaulters or non-defaulters, respectively) model with minor decrease in Hit Rate. CODA with minimization of Type I Error is especially useful in default prediction problem, since Type I Error costs more in this case. However, these results cannot be supported by the evidence obtained on real data. Further assessment both on real and modelled data is required to get the insight into reasons of such contradictory results.

As general results of this assessment are rather inspiring, further steps in CODA improvement will include extensions, described in Section 1.3, which includes implementation of more copula families and marginal distributions into the classifier, along with the implementation of vine-copulas or hierarchical copulas for flexibility on multidimensional datasets. In the same time, all these extensions should be carefully assessed on both real and modelled datasets. Moreover, CODA classifier was assessed within this research only as a binary classifier. Assessment of it as a probabilistic classifier will be the subject of following exercises.

References.

Breiman, 2001. Random Forests. // Machine Learning, 45, 5–32, 2001.

Ermolova, Penikas, 2015. М.Д. Ермолова, Г.И. Пеникас. Исследование взаимосвязи параметров моделей внутренних рейтингов оценки кредитного риска – вероятности дефолта и доли убытка при дефолте. // Управление финансовыми рисками, 1, 18–42, 2015.

Altman, 1968. E. Altman. Financial ratios, discriminant analysis and the prediction of corporate bankruptcy. // The Journal of Finance, No.4, 1968.

Han, Zhao, Liu, 2013. CODA: High Dimensional Copula Discriminant Analysis. // Journal of Machine Learning Research, 14, 629-671, 2013.

Lessmann et. al., 2013. S. Lessmann, H.-V. Seow, B. Baesens, L.C. Thomas. Benchmarking state-of-the-art classification algorithms for credit scoring: A ten-year update. // Pre-print. Available at: http://www.business-school.ed.ac.uk/waf/crc_archive/2013/42.pdf

Mogilat, 2015. А. Могилат. Банкротство компаний реального сектора в России: тенденции, структурные характеристики и основные факторы. Доступно по адресу: http://www.forecast.ru/_ARCHIVE/Presentations/HSE/2015/bank042015.pdf

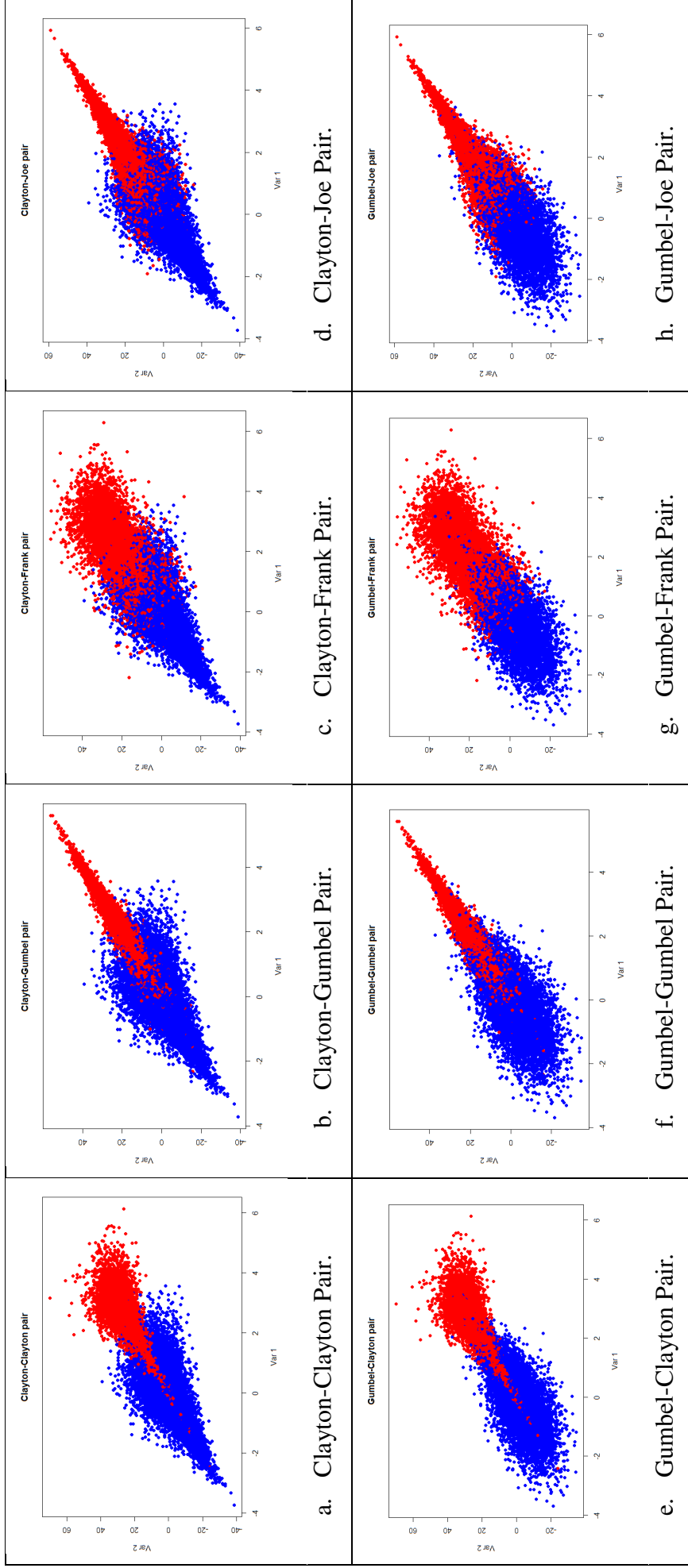
Penikas, 2010. Г.И. Пеникас. Модели «копула» в приложении к задачам финансов. // Журнал Новой Экономической Ассоциации, №7, стр. 24–44, 2010.

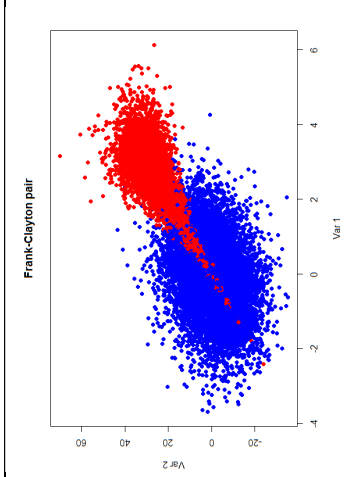
Sathe, 2006. S. Sathe. A Novel Bayesian Classifier using Copula Functions. // Pre-print. Available at: <http://arxiv.org/pdf/cs/0611150v1.pdf>

Scheungrab, 2013. E. Scheungrab. Copula based discriminant analysis with application. // Master Thesis, Technische Universität München. Available at: <https://mediatum.ub.tum.de/doc/1166376/1166376.pdf>

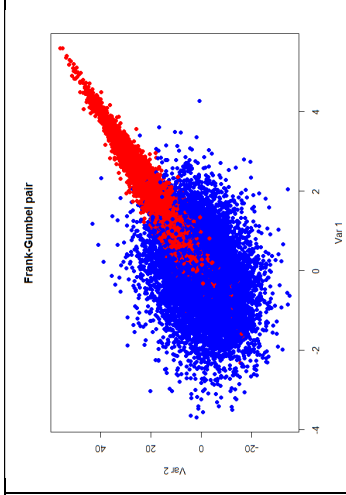
Sklar, 1959. A. Sklar. "Fonctions de répartition à n dimensions et leurs marges", Publ. Inst. Statist. Univ. Paris 8, 229–231, 1959.

Annex 1. Scatterplots of modelled data for bivariate case.

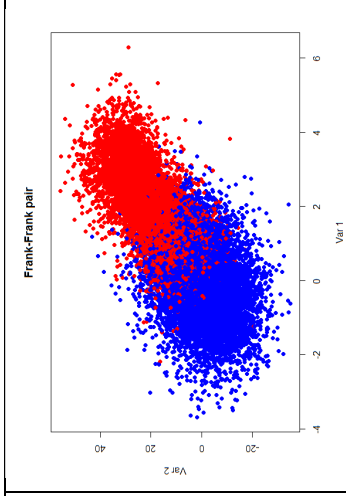




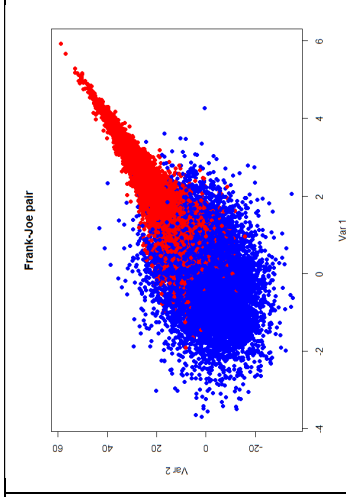
i. Frank-Clayton Pair.



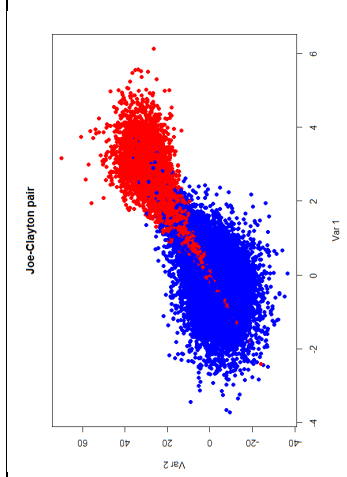
j. Frank-Gumbel Pair.



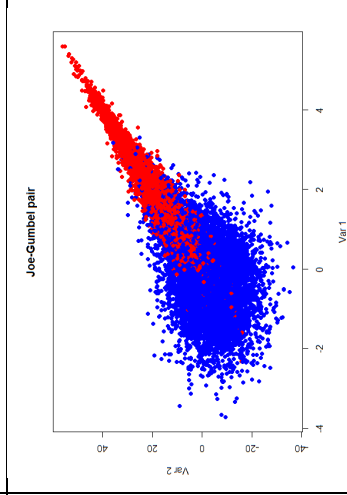
k. Frank-Frank Pair.



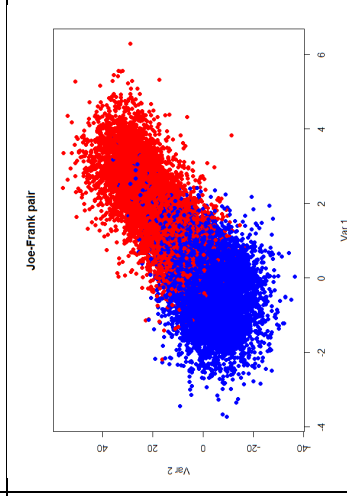
l. Frank-Joe Pair.



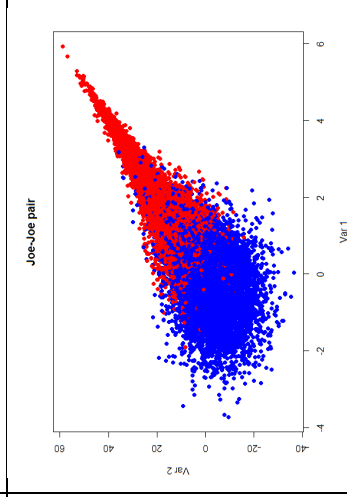
m. Joe-Clayton Pair.



n. Joe-Gumbel Pair.



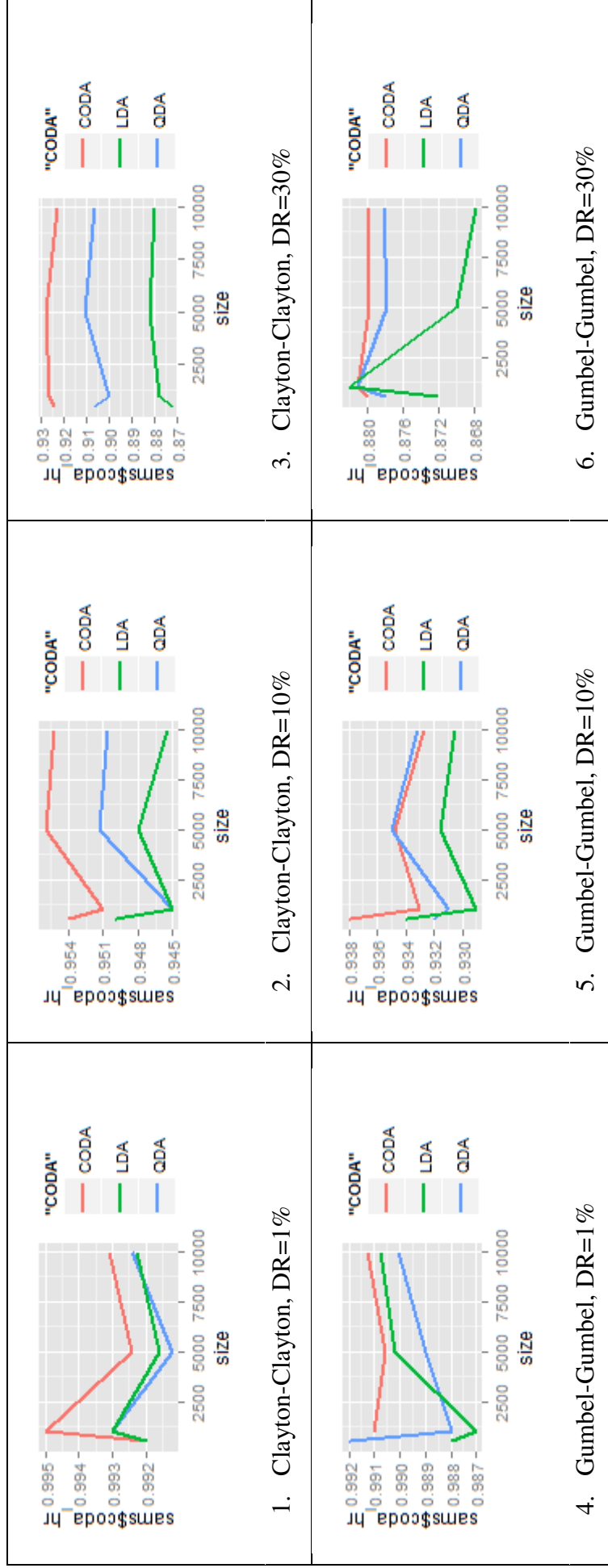
o. Joe-Frank Pair.



p. Joe-Joe Pair.

Annex 2. Results of testing the sensitivity to sample size on bivariate data (General Assessment).

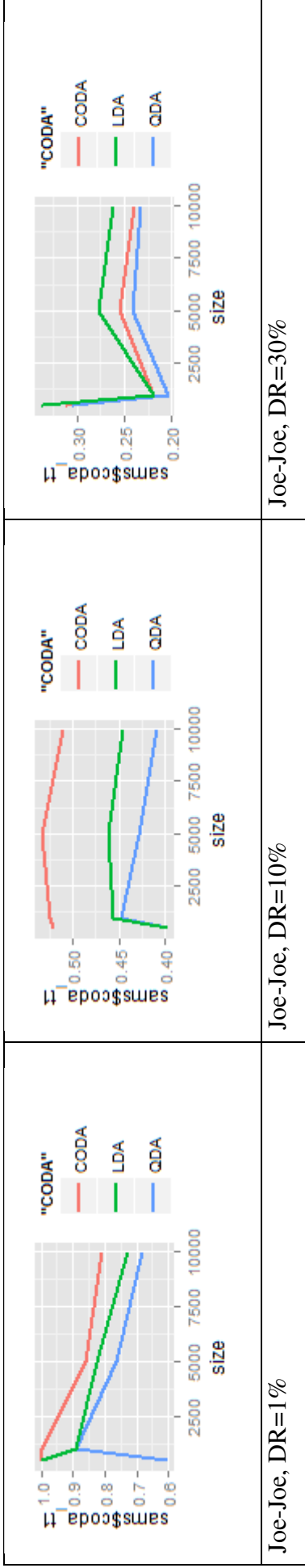
A. Hit Rate Dynamics.



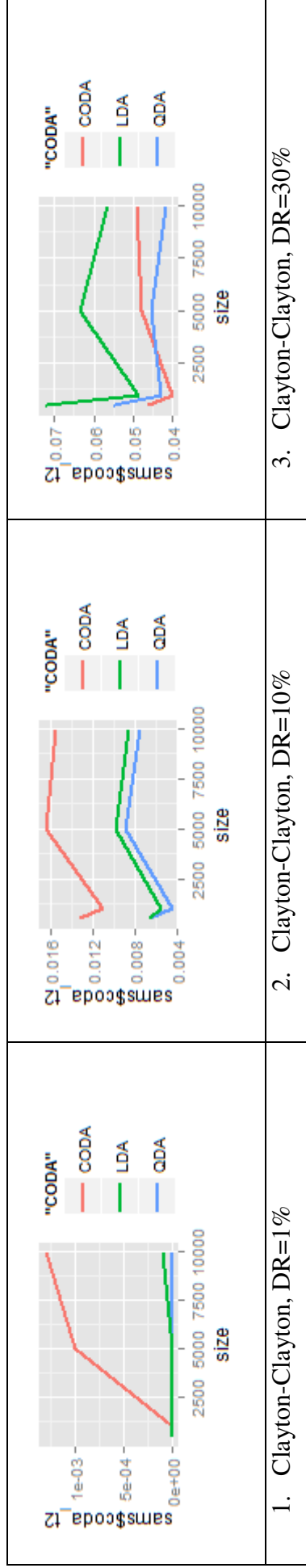
<p>7. Frank-Frank, DR=1%</p>	<p>8. Frank-Frank, DR=10%</p>	<p>9. Frank-Frank, DR=30%</p>
<p>10. Joe-Joe, DR=1%</p>	<p>11. Joe-Joe, DR=10%</p>	<p>12. Joe-Joe, DR=30%</p>

B. Type I Error Dynamics.

<p>Line graph showing sam\$ccda vs size for Clayton-Clayton, DR=1%. The y-axis ranges from 0.6 to 0.8. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.8, while QDA starts lower at ~0.65 and increases to ~0.75 at size 10000.</p>	<p>Line graph showing sam\$ccda vs size for Clayton-Clayton, DR=10%. The y-axis ranges from 0.30 to 0.50. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.45, while QDA starts lower at ~0.35 and increases to ~0.45 at size 10000.</p>	<p>Line graph showing sam\$ccda vs size for Clayton-Clayton, DR=30%. The y-axis ranges from 0.15 to 0.30. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.25, while QDA starts lower at ~0.15 and increases to ~0.25 at size 10000.</p>
<p>1. Clayton-Clayton, DR=1%</p>	<p>2. Clayton-Clayton, DR=10%</p>	<p>3. Clayton-Clayton, DR=30%</p>
<p>Line graph showing sam\$ccda vs size for Gumbel-Gumbel, DR=1%. The y-axis ranges from 0.75 to 0.85. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.85, while QDA starts lower at ~0.75 and increases to ~0.85 at size 10000.</p>	<p>Line graph showing sam\$ccda vs size for Gumbel-Gumbel, DR=10%. The y-axis ranges from 0.35 to 0.50. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.45, while QDA starts lower at ~0.35 and increases to ~0.45 at size 10000.</p>	<p>Line graph showing sam\$ccda vs size for Gumbel-Gumbel, DR=30%. The y-axis ranges from 0.20 to 0.25. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.25, while QDA starts lower at ~0.20 and increases to ~0.25 at size 10000.</p>
<p>Gumbel-Gumbel, DR=1%</p>	<p>Gumbel-Gumbel, DR=10%</p>	<p>Gumbel-Gumbel, DR=30%</p>
<p>Line graph showing sam\$ccda vs size for Frank-Frank, DR=1%. The y-axis ranges from 0.5 to 0.8. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.75, while QDA starts lower at ~0.65 and increases to ~0.75 at size 10000.</p>	<p>Line graph showing sam\$ccda vs size for Frank-Frank, DR=10%. The y-axis ranges from 0.36 to 0.41. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.40, while QDA starts lower at ~0.35 and increases to ~0.40 at size 10000.</p>	<p>Line graph showing sam\$ccda vs size for Frank-Frank, DR=30%. The y-axis ranges from 0.20 to 0.25. The x-axis ranges from 0 to 10000. Three methods are compared: CODA (red), LDA (green), and QDA (blue). CODA and LDA maintain high values around 0.25, while QDA starts lower at ~0.20 and increases to ~0.25 at size 10000.</p>
<p>Frank-Frank, DR=1%</p>	<p>Frank-Frank, DR=10%</p>	<p>Frank-Frank, DR=30%</p>



C. Type II Error Dynamics.



<p>Gumbel-Gumbel, DR=1%</p>	<p>Gumbel-Gumbel, DR=10%</p>	<p>Gumbel-Gumbel, DR=30%</p>
<p>Frank-Frank, DR=1%</p>	<p>Frank-Frank, DR=10%</p>	<p>Frank-Frank, DR=30%</p>
<p>Joe-Joe, DR=1%</p>	<p>Joe-Joe, DR=10%</p>	<p>Joe-Joe, DR=30%</p>

Annex 3. Means and standard deviations of performance metrics in General Assessment.

By copula families combination

Set	Hit Rate	HR std. dev.	Type I Error	T1 std. dev.	Type II Error	T2 std. dev.
cc	95.36%	2.60%	34.01%	19.10%	1.86%	1.99%
cg	95.08%	3.23%	33.11%	15.46%	1.84%	1.75%
cf	94.09%	3.87%	37.90%	13.39%	1.92%	2.09%
cj	93.97%	4.22%	41.17%	16.73%	1.85%	1.96%
gc	92.39%	4.38%	46.52%	27.26%	4.54%	3.94%
gg	92.32%	4.57%	46.05%	25.11%	4.10%	3.66%
gf	92.01%	5.00%	48.02%	23.44%	4.09%	3.61%
gj	91.45%	5.15%	51.87%	25.91%	4.33%	3.72%
fc	95.12%	2.84%	31.83%	18.51%	2.90%	2.55%
fg	94.88%	3.24%	32.86%	15.65%	2.53%	2.18%
ff	94.05%	3.92%	37.45%	13.79%	2.70%	2.60%
fj	93.78%	4.25%	42.25%	18.52%	2.54%	2.45%
jc	91.97%	4.88%	45.15%	27.79%	5.52%	5.20%
jg	92.26%	4.85%	44.77%	25.62%	4.58%	4.27%
jf	91.95%	5.61%	46.17%	23.98%	4.62%	4.66%
jj	91.03%	6.05%	51.01%	25.32%	5.06%	5.07%

By sample size

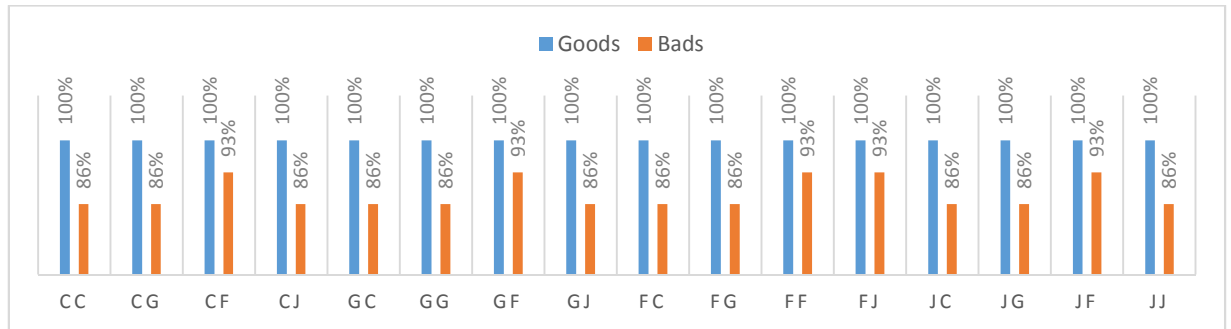
Set	Hit Rate	HR std. dev.	Type I Error	T1 std. dev.	Type II Error	T2 std. dev.
0.1k	90.29%	4.69%	36.41%	14.97%	5.83%	5.24%
0.5k	93.35%	4.63%	44.39%	26.44%	3.41%	3.26%
1k	94.10%	4.37%	46.00%	25.12%	2.94%	3.22%
5k	94.25%	4.26%	45.44%	23.94%	2.72%	2.95%
10k	94.32%	4.12%	44.05%	21.71%	2.53%	2.69%

By default rate

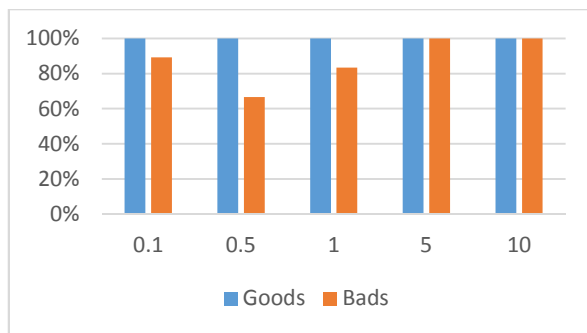
Set	Hit Rate	HR std. dev.	Type I Error	T1 std. dev.	Type II Error	T2 std. dev.
0.01	99.18%	0.19%	73.79%	16.89%	0.11%	0.10%
0.1	93.91%	1.64%	42.98%	9.18%	2.07%	1.36%
0.3	88.70%	2.78%	21.24%	5.27%	7.06%	3.12%

Annex 4. Performance characteristics of copula family selection method.

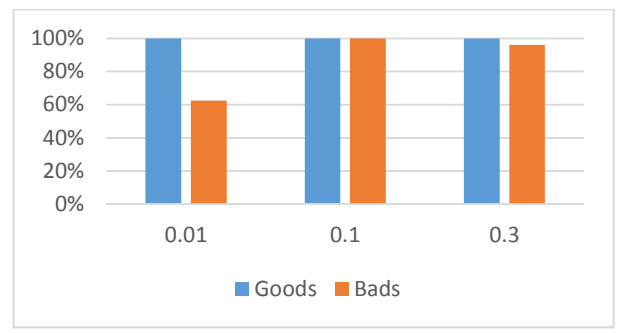
A. Percentage of estimations, where LL-CODA classifier had selected correct copula family for given class.



a) by copulas used for sample generation.

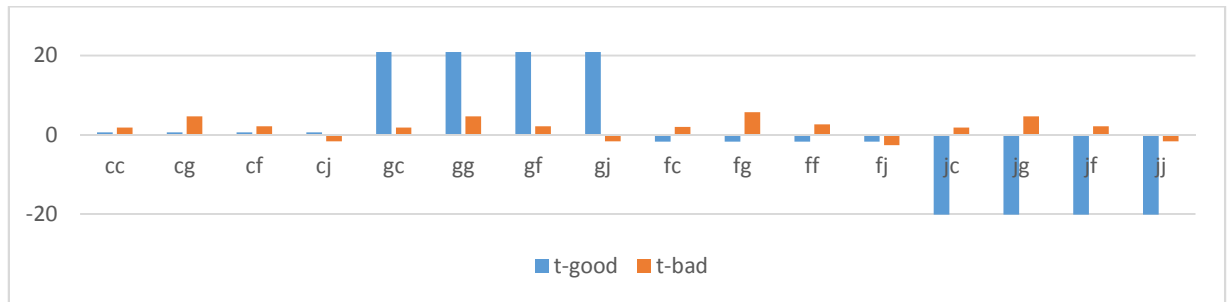


b) by sample size.

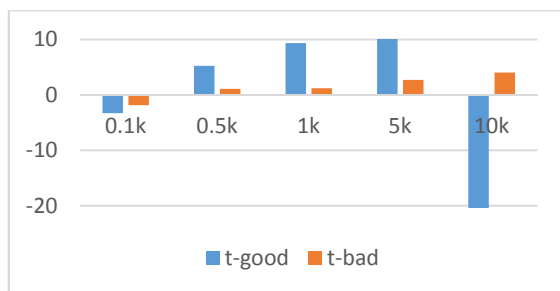


c) by default rate

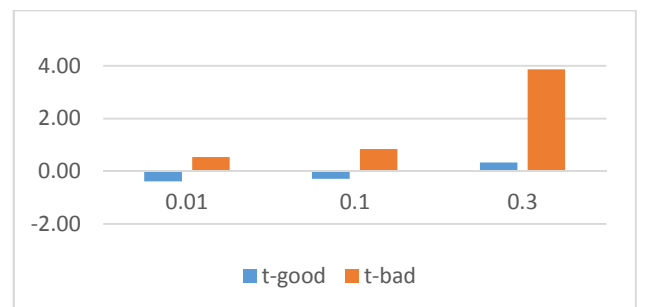
B. Averaged metrics of copula parameter estimations.



a) t-scores, by copulas used for sample generation.



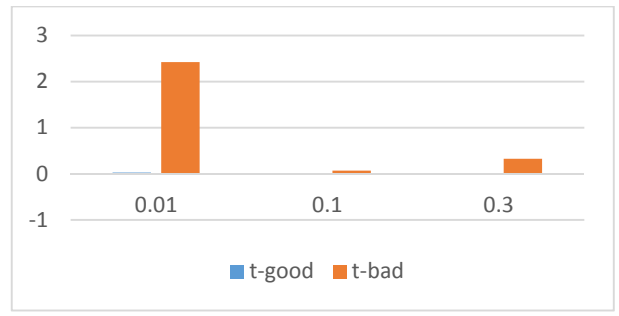
b) t-scores, by sample size.



c) t-scores, by default rate



d) deviations from true parameter, by sample size



e) deviations from true parameter, by default rate

C. Share of cases, where LL-CODA selected best combination of copula families in terms of given metrics.



a) by copulas used for sample generation.



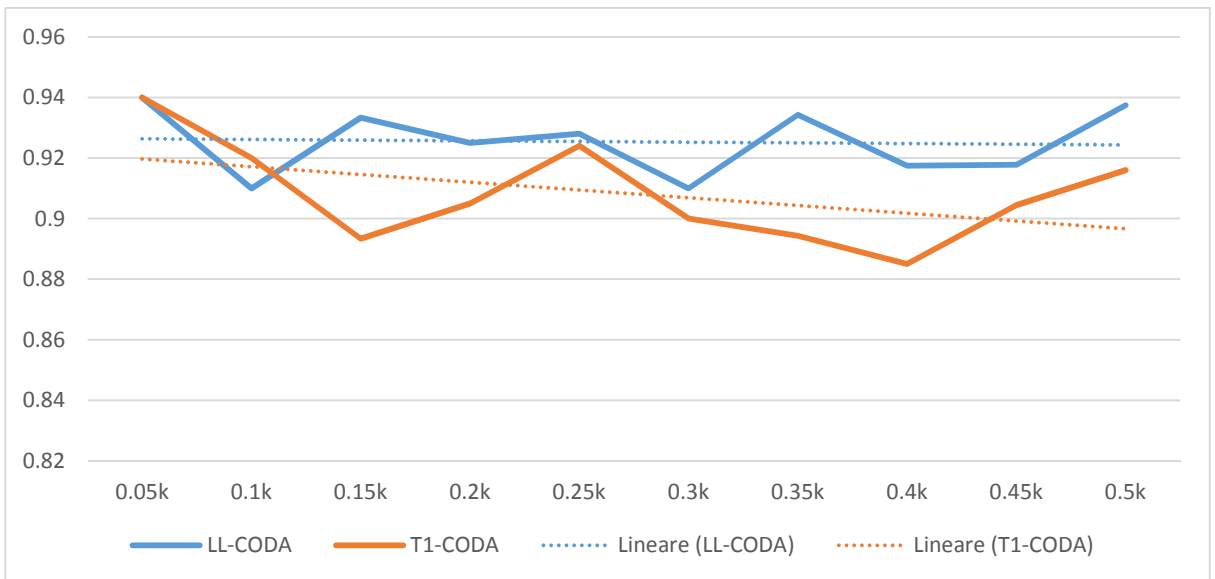
b) by sample size.



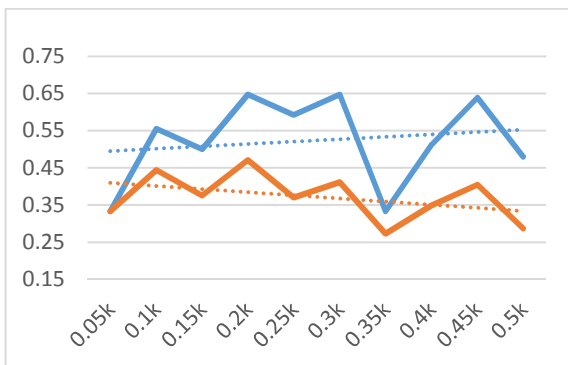
c) by default rate

Annex 5. Averaged on 4 distributions CODA performance characteristics on small samples.

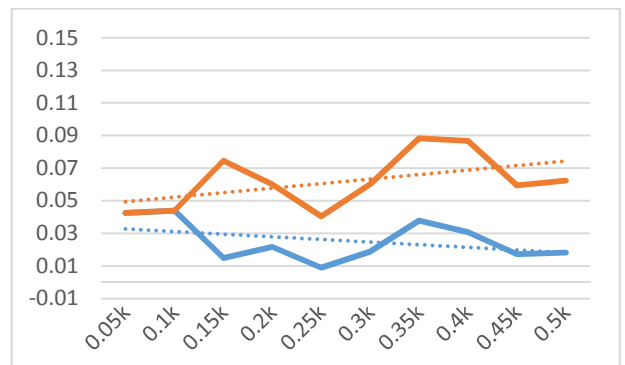
A. Dynamics of CODA performance metrics.



a) Hit Rate

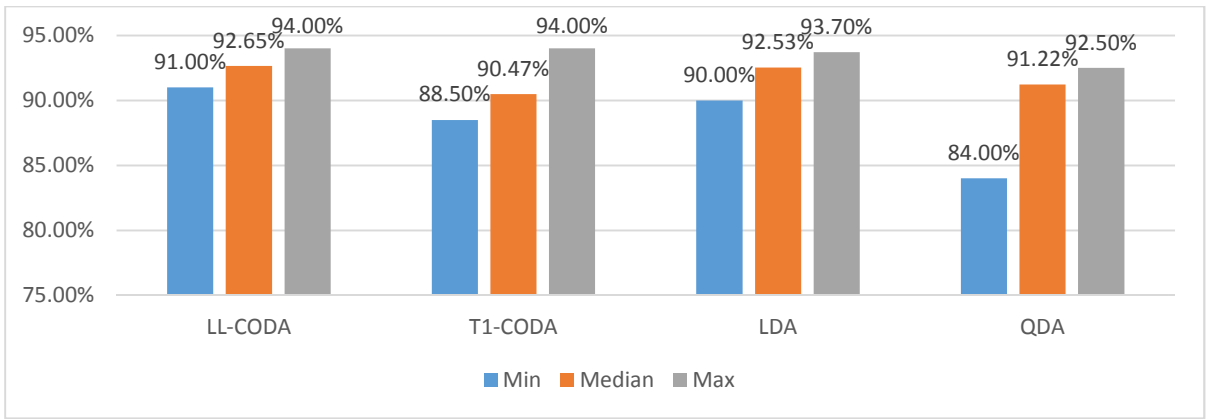


b) Type I Error

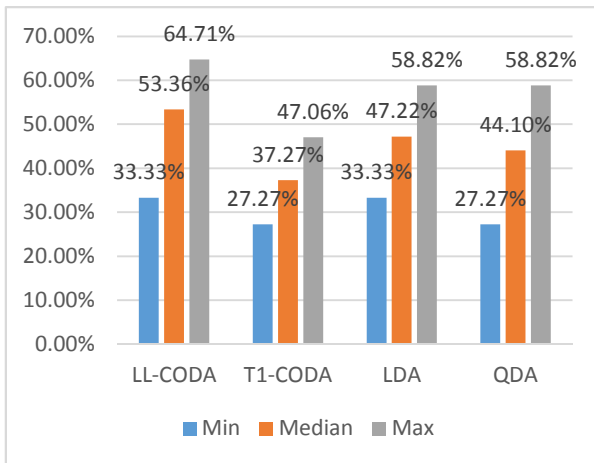


c) Type II Error

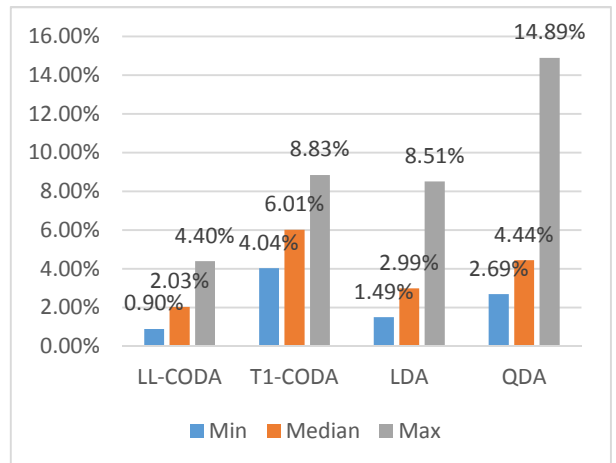
B. Aggregated results for CODA compare to LDA and QDA.



a) Hit Rate



b) Type I Error



c) Type II Error

Annex 6. Results of single factor analysis on banks sample.

Factor	AUC	Factor	AUC
BP_auc	66.09%	SRTS_auc	58.59%
NORM_KAP_auc	65.54%	PZS_F_auc	58.53%
CP_auc	65.38%	NORM_H4_auc	58.39%
ODB_auc	64.46%	KE_3_auc	58.24%
PKB_auc	64.39%	SO_LONG_auc	58.21%
PRIB_auc	64.09%	RPR_auc	58.20%
MBK_auc	64.06%	RUB1_auc	57.69%
PDMBK_auc	63.85%	SPC_auc	57.55%
PKK_auc	63.83%	RES_auc	57.36%
NMO_auc	63.79%	VDFL_90_auc	57.10%
PZS_auc	63.56%	KE_F_auc	56.82%
KSDB_auc	63.51%	VDFL_auc	56.76%
ORB_auc	63.24%	ORCB_auc	56.72%
RA_auc	62.17%	VDUL_auc	56.46%
NORM_LAT_auc	61.92%	SIP_auc	56.45%
LA_auc	61.80%	PNA_auc	56.17%
KE_auc	61.80%	PDPS_B_auc	56.11%
OKS_auc	61.56%	VB_auc	56.02%
RBP_auc	61.34%	VBD_auc	56.02%
NORM_SP_auc	61.33%	VBK_auc	56.02%
RUB2_auc	61.07%	DFL_auc	55.78%
PDK_auc	61.00%	DDBO_auc	55.67%
DK_auc	60.96%	KE_F_12_auc	55.07%
CA_VAL_auc	60.71%	NORM_H3_auc	54.82%
OV_auc	60.69%	MP_auc	54.61%
RUB_auc	60.48%	DUB1_auc	54.29%
CAB_auc	60.48%	RDOS_auc	53.48%
RPFL_auc	60.39%	RPD_auc	53.21%
SO_auc	60.24%	PD_auc	52.80%
RK_auc	60.20%	KSCB_auc	52.20%
NORM_OVM_auc	60.18%	DAC_auc	52.18%
KE_Prom_auc	59.95%	SK_auc	51.71%
SBP_auc	59.89%	KE_F_3_12_auc	50.93%
DUB2_auc	59.72%	OS_auc	50.70%
RSA_auc	59.70%	NORM_H1_auc	50.64%
VDFL_30_auc	59.35%	SP_auc	50.63%
DUB_auc	59.23%	CA_auc	50.00%
RKK_auc	59.20%	NORM_H2_auc	47.88%
PDFL_auc	58.87%	UF_auc	47.02%
KE_LONG_auc	58.74%	KE_3_12_auc	44.22%